

Use of molecular markers in biosystematics

Is it that simple..?

Gabriela Šrámková

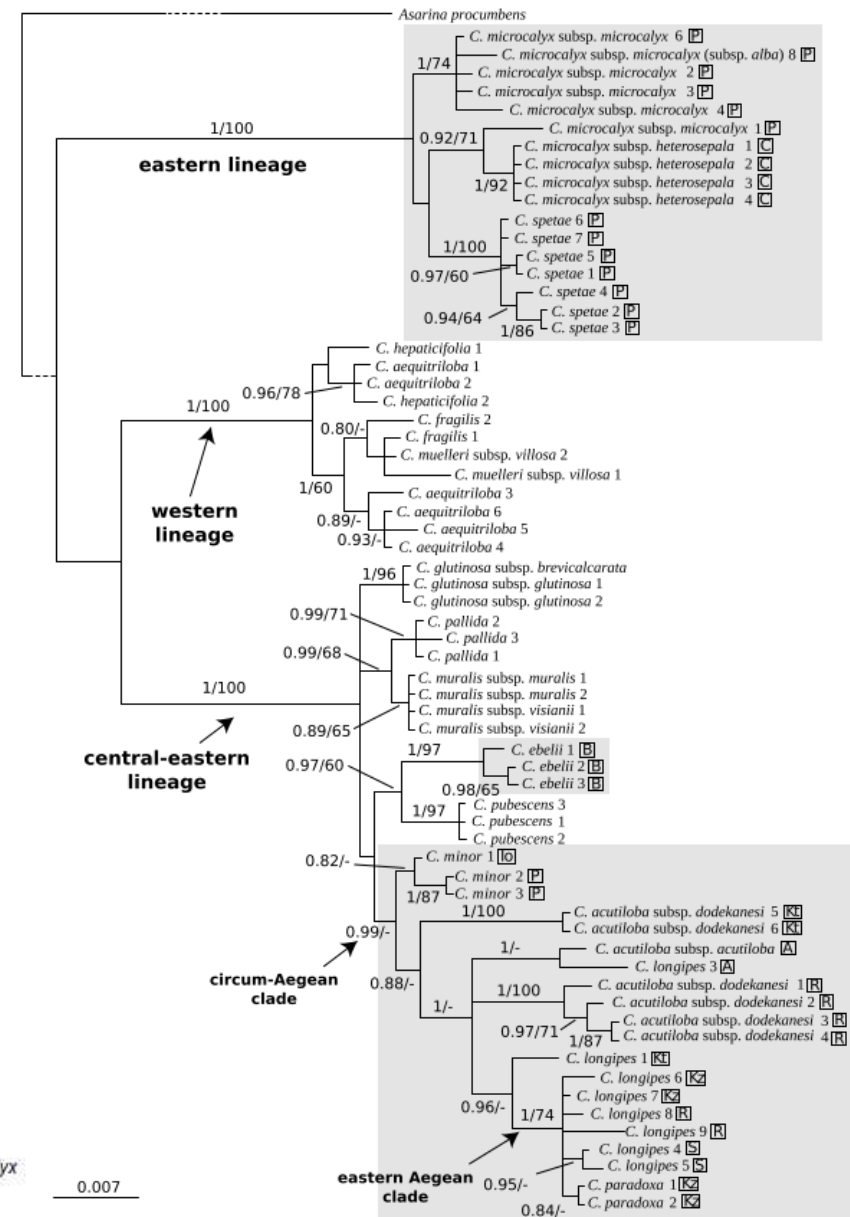
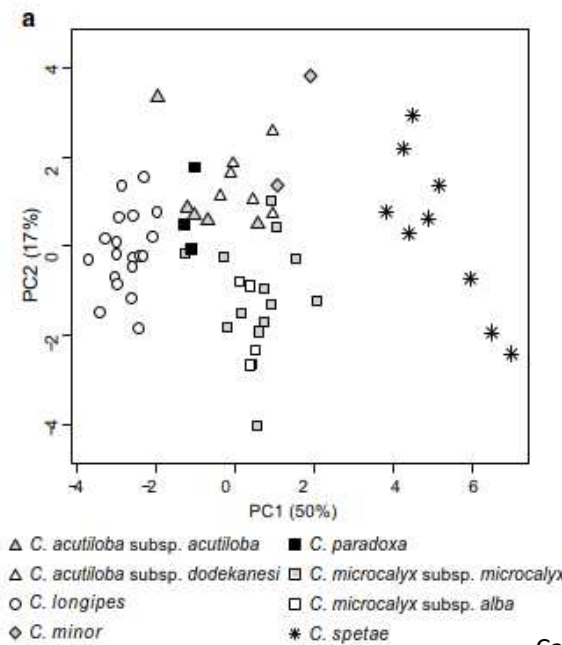


FACULTY OF
SCIENCE
Charles University



Biosystematics

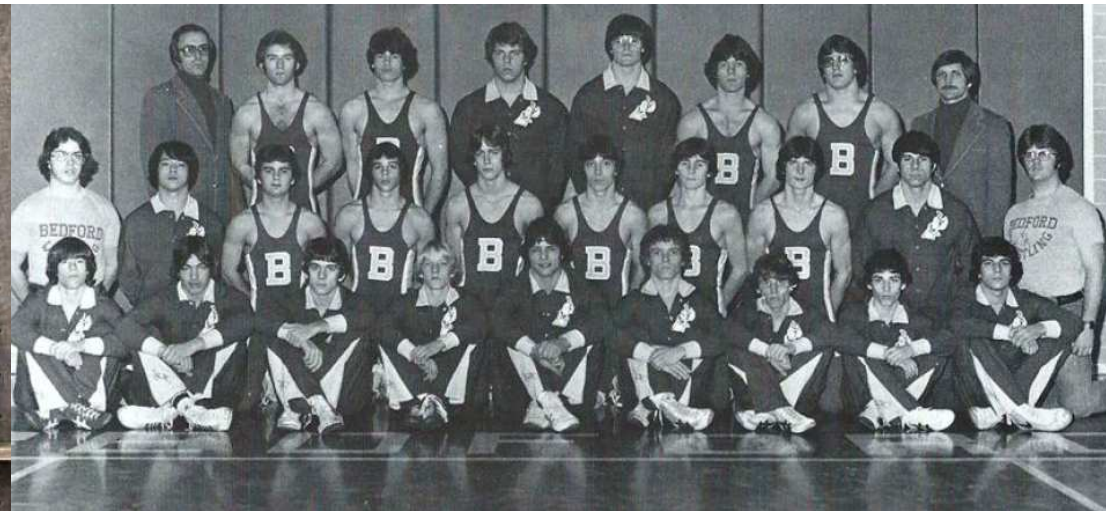
- complex study of taxonomy of organisms based on comparisons of their ecological, cytological or genetic characteristics



Carnicero et al. 2021 *Cymbalaria*

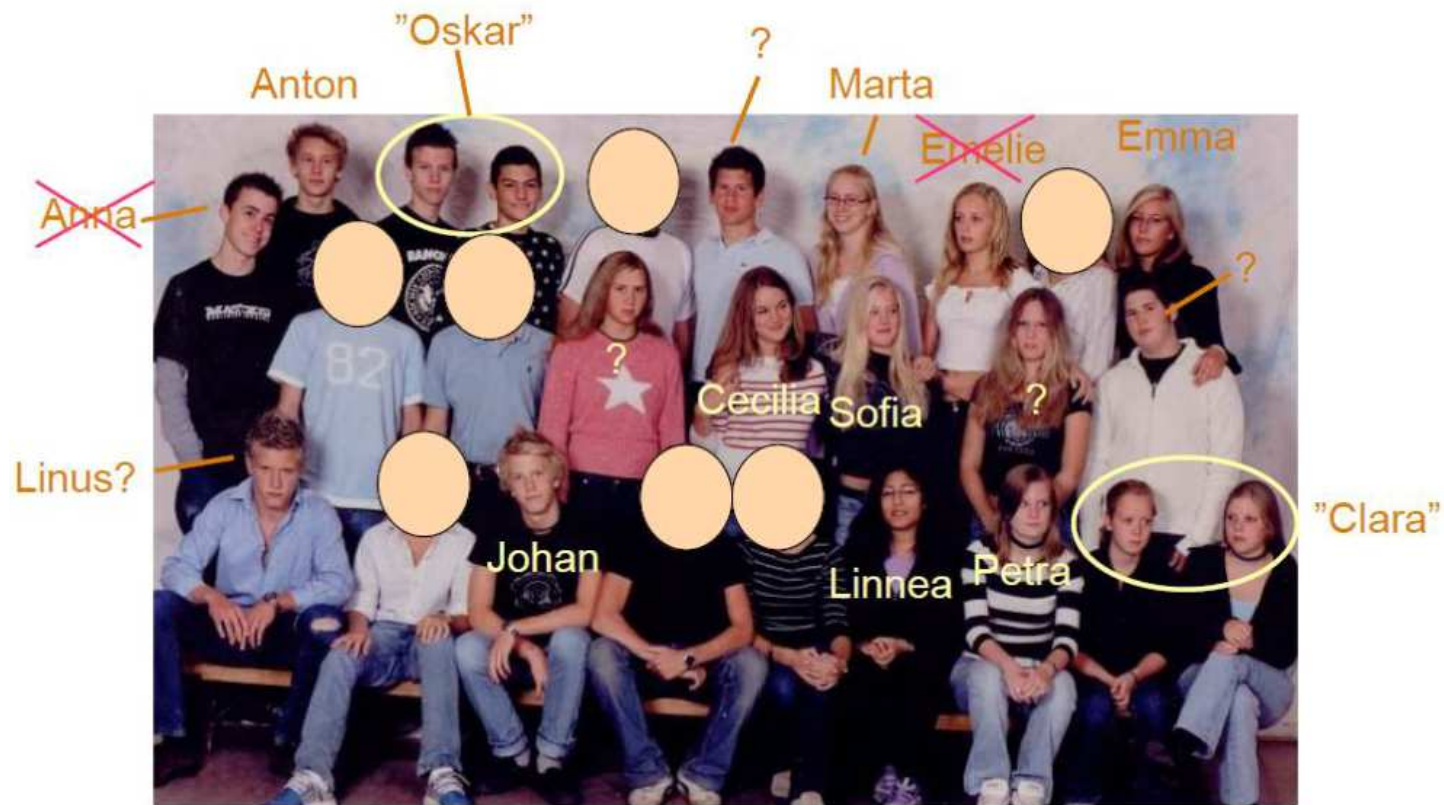
Biosystematics – what for?

- Imagine people or things do not have names...



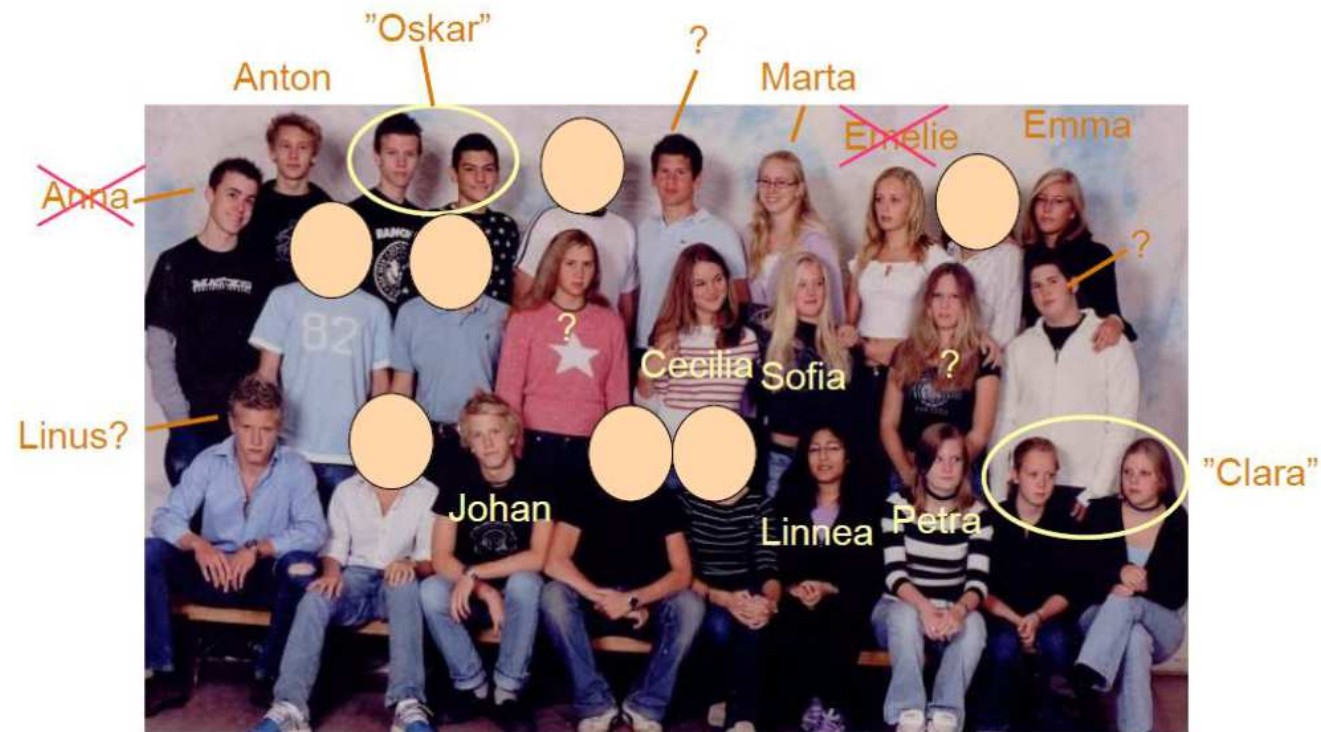
Biosystematics – what for?

- Or a lot of names would be wrong...

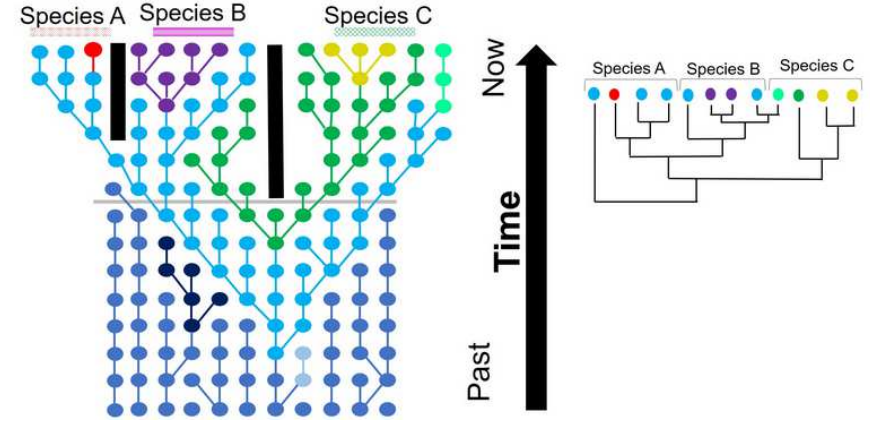
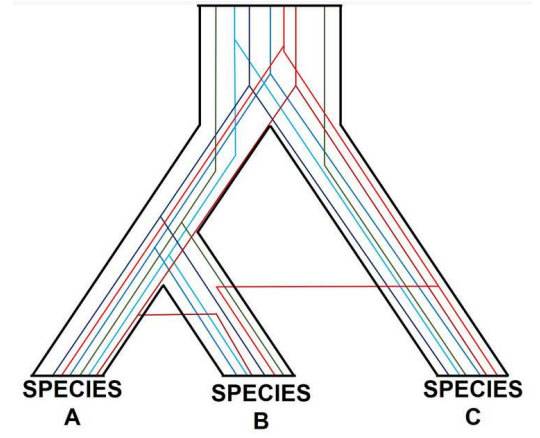
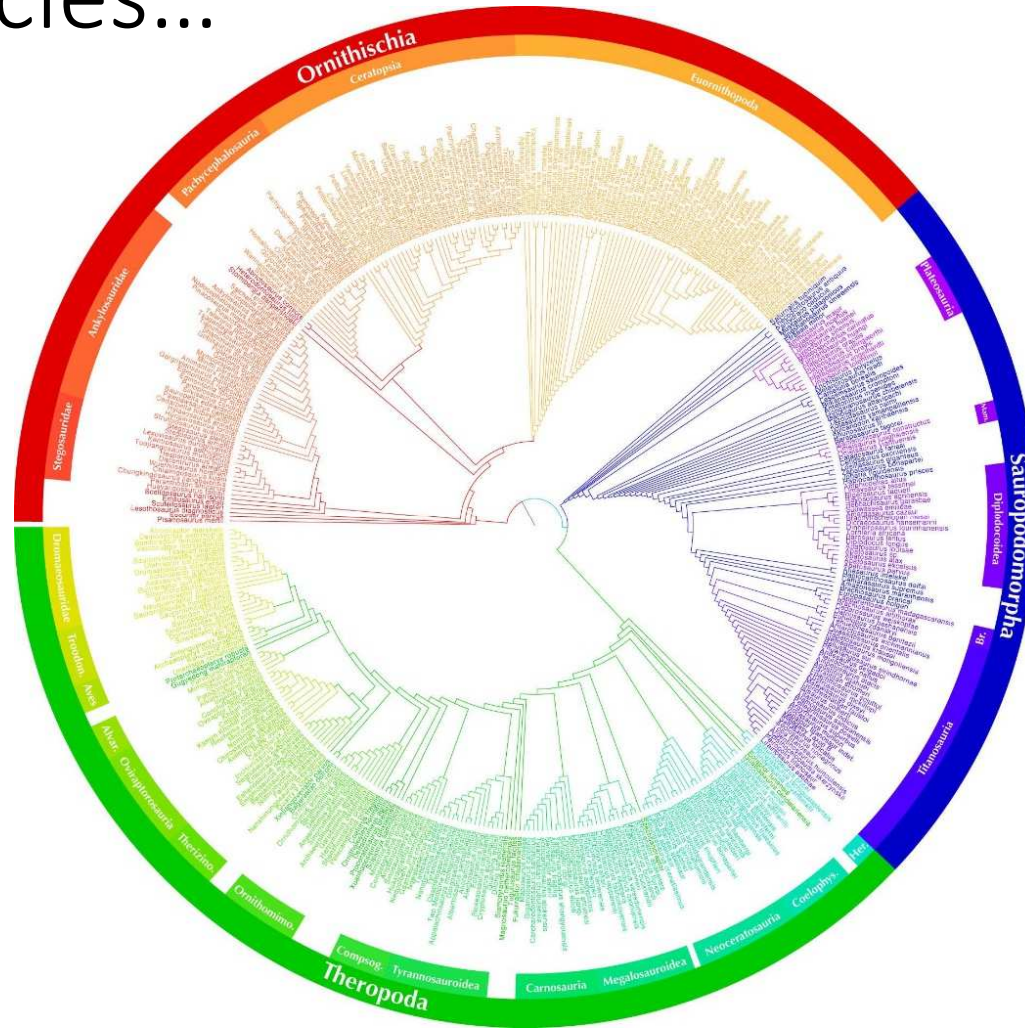


Biosystematics – what for?

- Biosystematics creates "language" and interprets it. Without it, many other disciplines (such as protection, evolutionary biology etc.) could not exist...



Species...





Species...

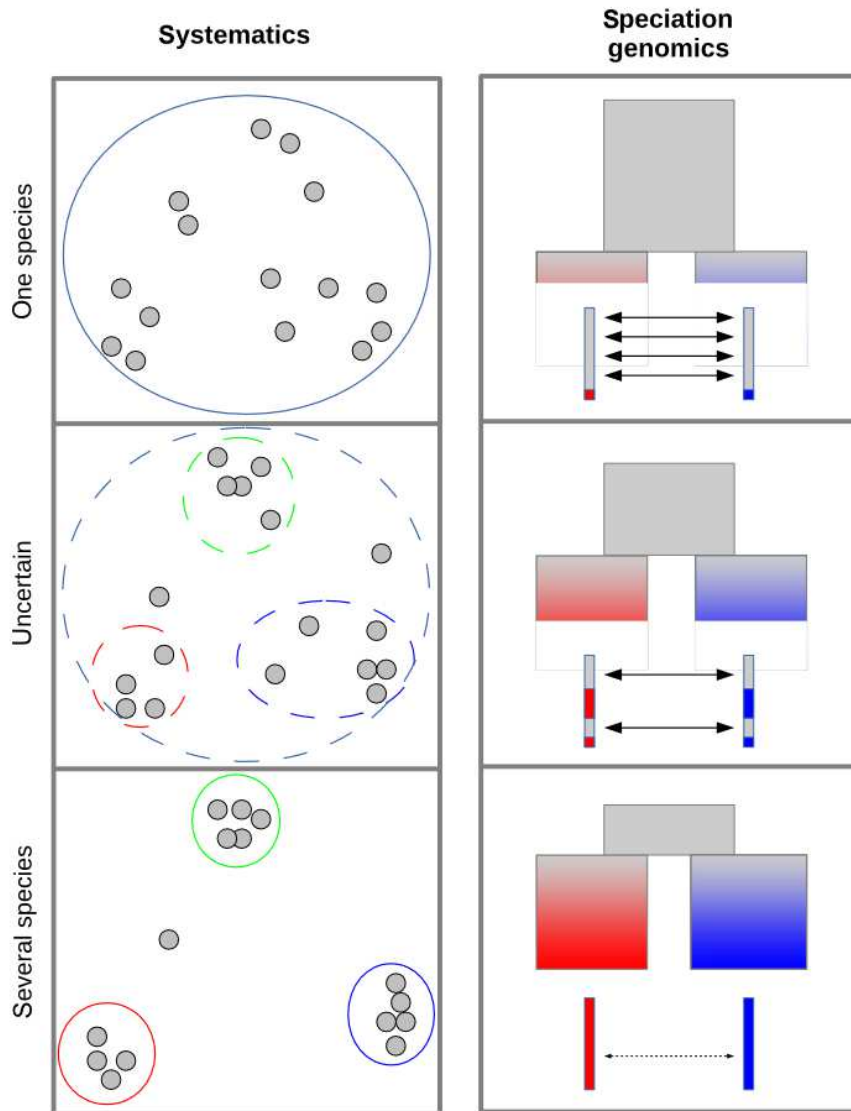


FIGURE 1 Two views on the continuum of speciation. Left: Species are defined as groups of organisms resembling each other according to an arbitrary set of variables. Right: Species are defined as entities sufficiently diverged such that gene flow (arrows) is very rare or inexistent. Top: unambiguous single-species situation. Bottom: unambiguous multiple-species situation. Intermediate: ambiguous situation. Ambiguous situations appear when groups can be identified but intermediate individuals are common (left) and when gene flow exists but is limited to a fraction of the genome (right)



Species...

Species concepts

- biological, morphological, phylogenetic...
- there are many and definitions vary
- the definition is based on breaking the continuum of variation
- the assumption of reproductive isolation

- the species is understood differently - by taxonomists, conservationists, in different countries, in
- different times, different groups of organisms
- (zoologists/botanists)

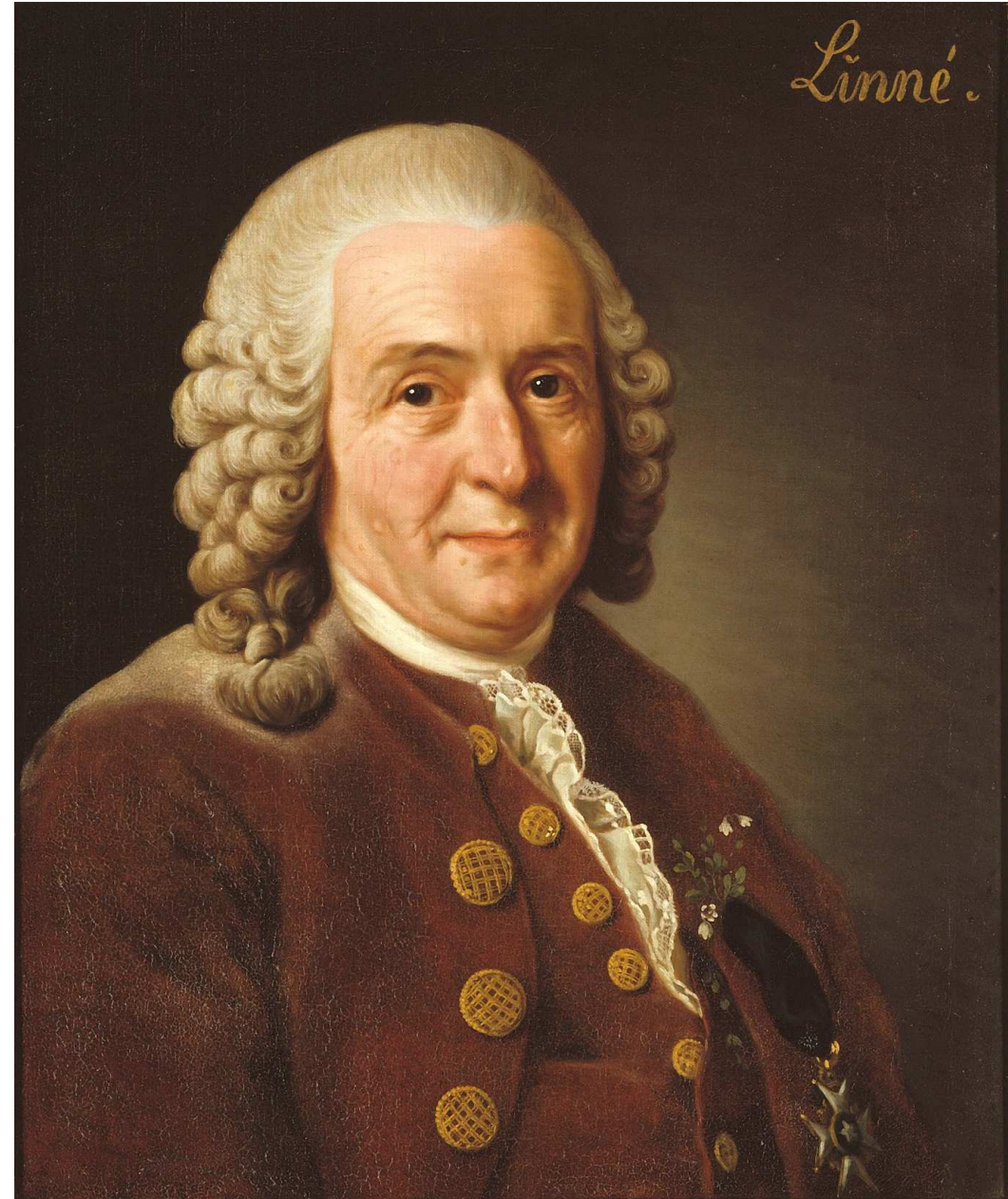


Where to begin?



Where to begin?

- Thorough study of the complex
- Sampling
- Phylogenetic relationships
- Ecological relationships
- Morphological analyses
- Taxonomic re-evaluation





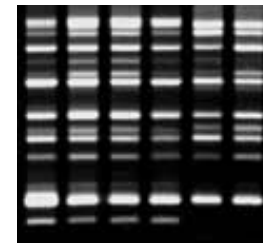
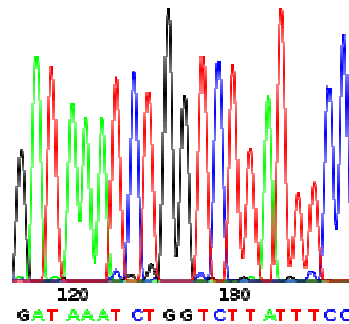
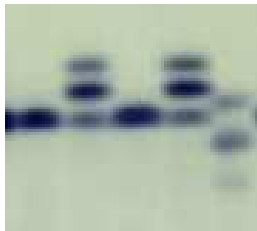
Where to begin?

- Thorough study of the complex
- Sampling
- **Phylogenetic relationships**
- Ecological relationships
- Morphological analyses
- Taxonomic re-evaluation



Molecular markers

- information about an organism obtained from the analysis of its molecules - proteins, DNA, RNA
- marker - character, unit of information - a purposefully or randomly selected part of the total information
- markers tell about the genetic similarity (relatedness) of individuals, populations or species
- electrophoresis of macromolecules





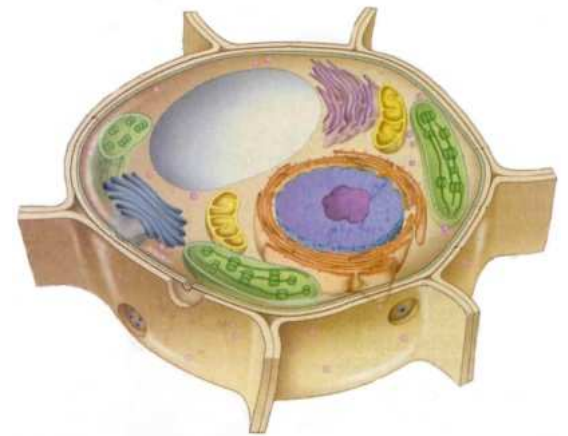
Molecular markers

The nature of molecular data:

- provide information about the genotype of an individual
- information independent of environmental conditions
- assumption of selective neutrality - no effect on fitness
- qualitative information - presence of fragment, allele, nucleotide
- unique information about the organism - clone identification
- changes during generative reproduction – recombination

Storage of genetic information:

- nucleus, plastids, mitochondria





Molecular markers

Types of questions:

- identification of clones
 - genetic diversity of clonal plants
- genetic structure of populations
 - intra-population genetic diversity
 - H-W equilibrium test
 - relationships between populations, distribution of genetic variability
 - gene flow
- study of plant migration
 - phylogeography
 - study of invasions
- type of reproductive system
- systematic studies at all levels, reconstruction of phylogeny
- hybridization, polyploidization

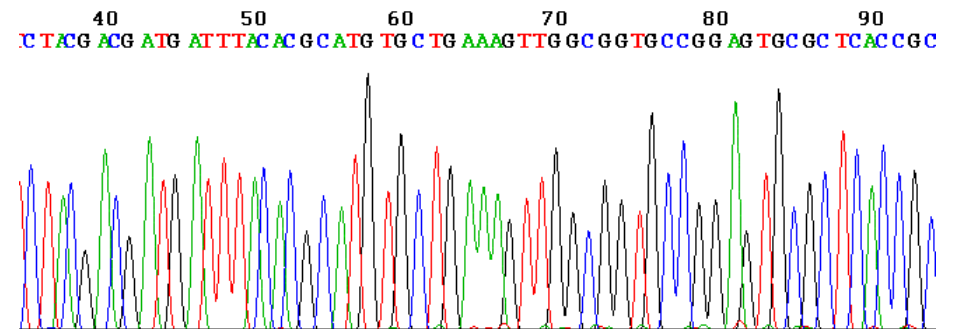


Sequencing

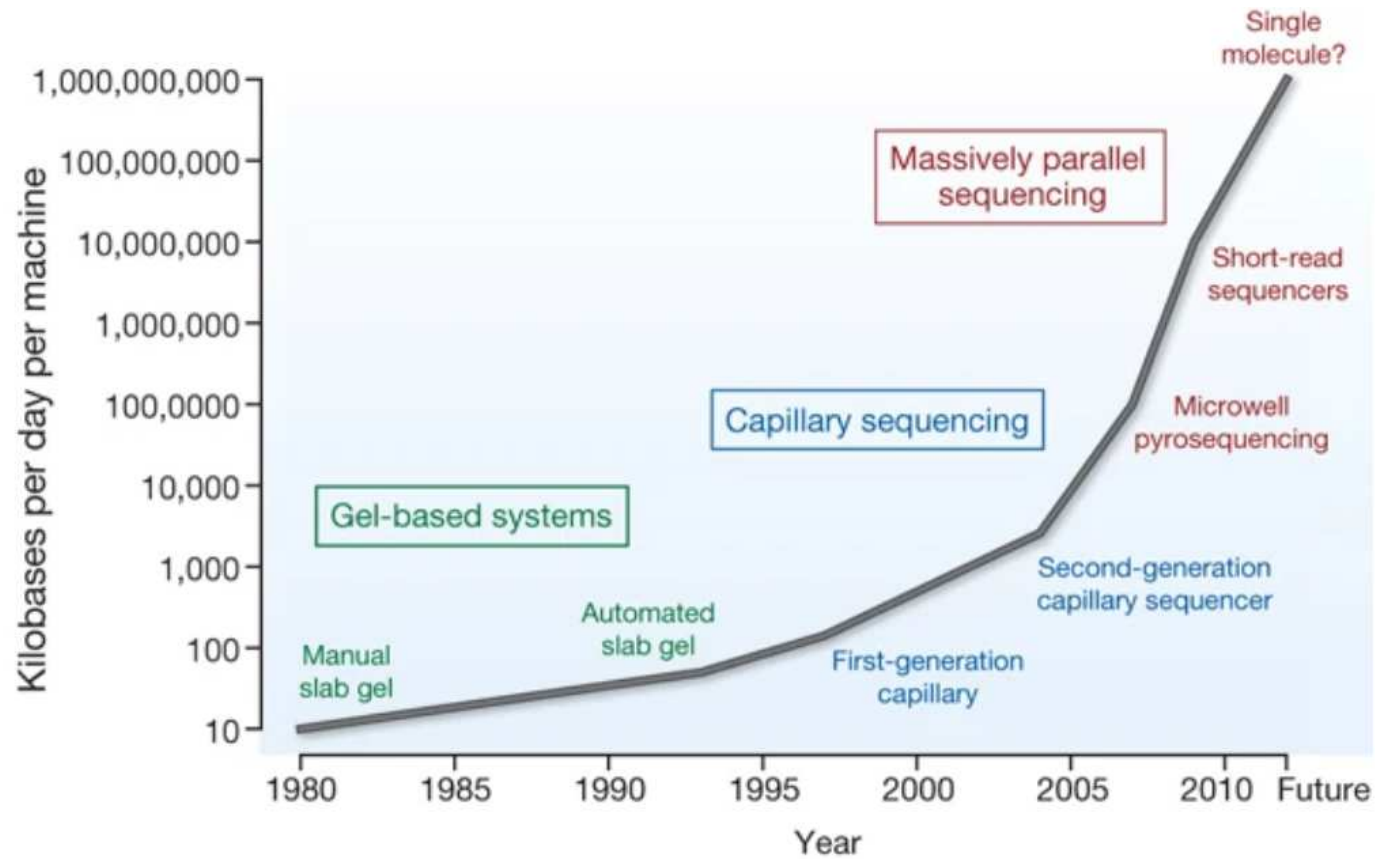
- finding the order of nucleotides in a DNA strand

...ATATATAGGCAAGGAATCTCTATTATTAATCATT...

- using information to determine the course and rate of evolution
- determining the similarity and relatedness of taxa
- capillary sequencing vs next generation sequencing



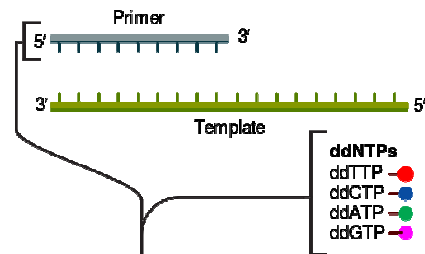
The evolution of sequencing



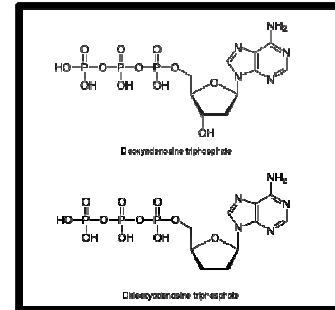
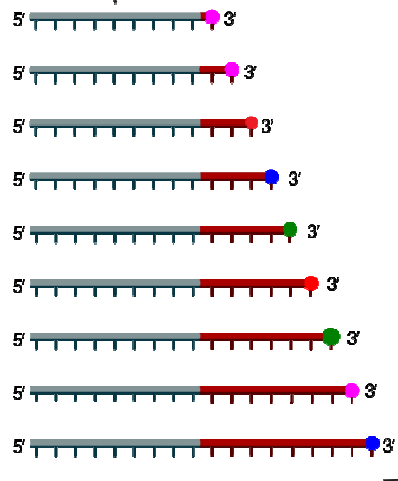
Sanger sequencing

① Reaction mixture

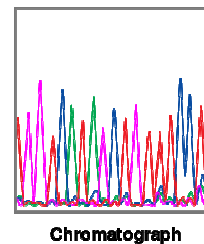
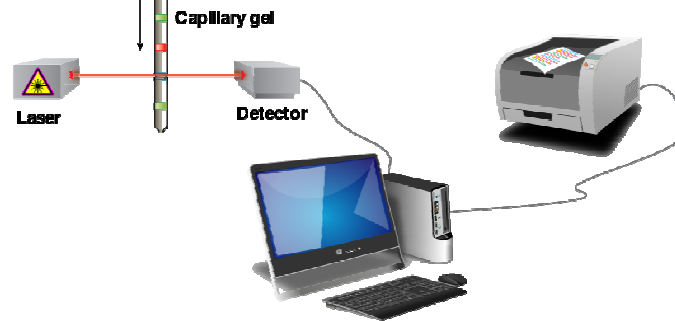
- ▶ Primer and DNA template
- ▶ DNA polymerase
- ▶ ddNTPs with flouorchromes
- ▶ dNTPs (dATP, dCTP, dGTP, and dTTP)



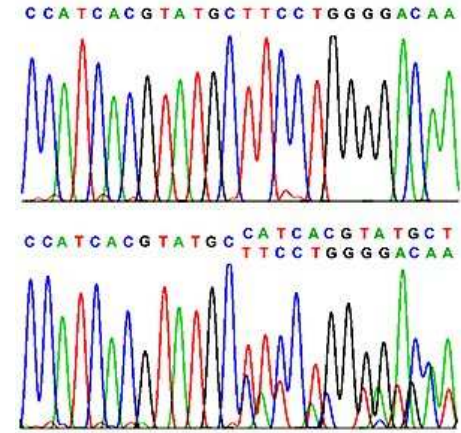
② Primer elongation and chain termination



③ Capillary gel electrophoresis separation of DNA fragments



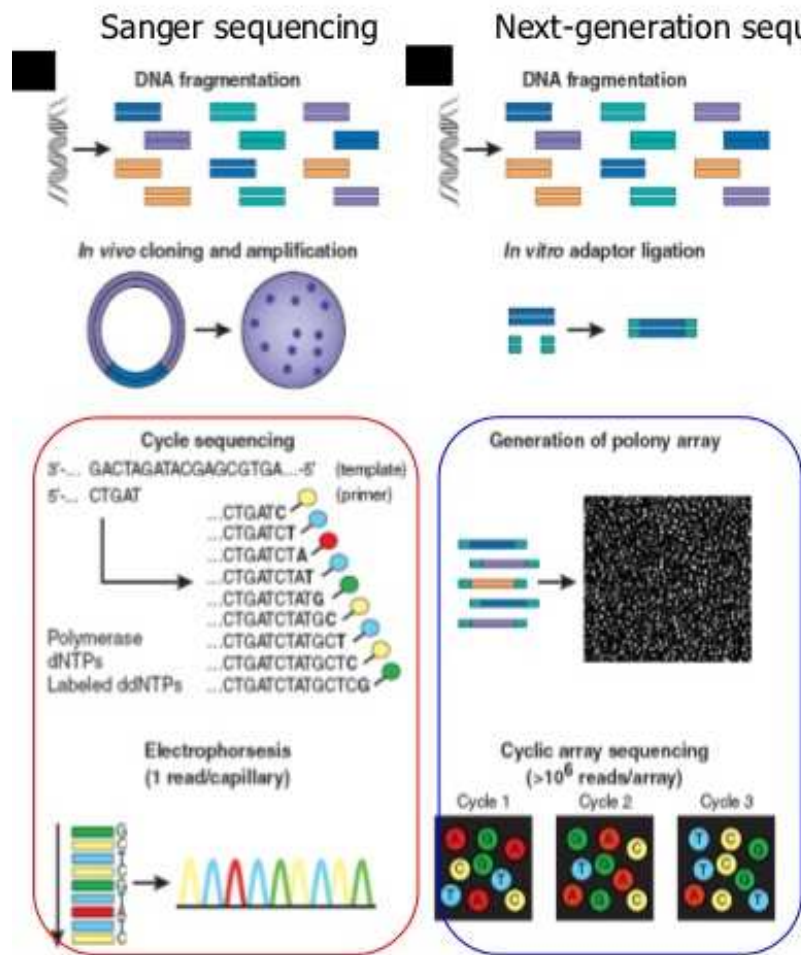
④ Laser detection of flouorchromes and computational sequence analysis



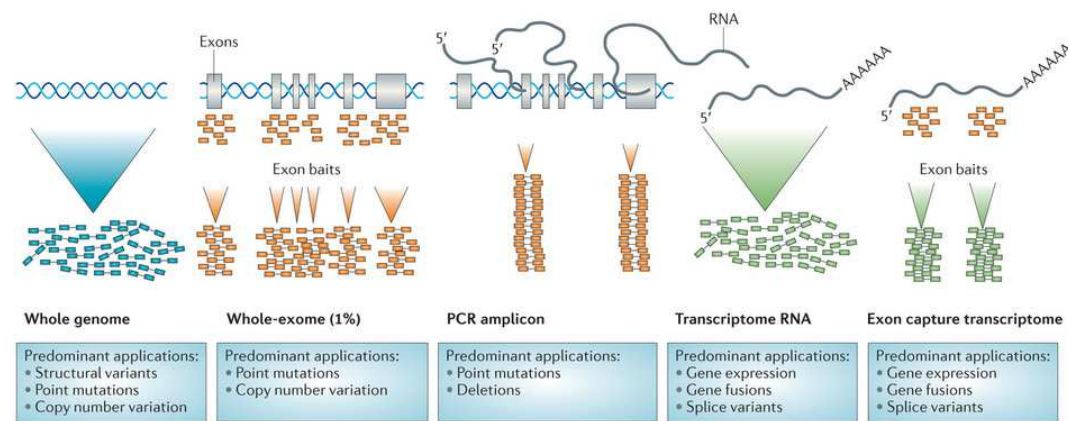
Disadvantages

- cloning
- low possibility of parallelisation

NGS (next generation sequencing)

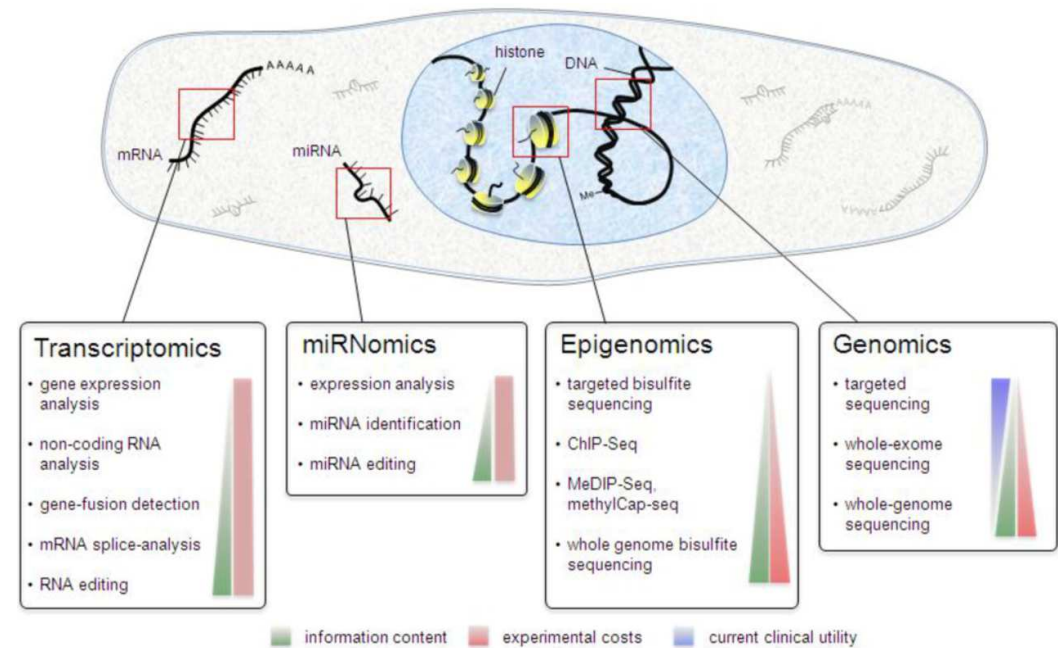


Massive parallelization



Methods

- Whole genome sequencing
- Libraries with reduced representation:
 - Transcriptomes
 - Exomes
 - Plastomes
 - Restriction-based methods
 - RadSeq, ddRadSeq
 - Target enrichment
 - HybSeq
 - Amplicons (PCR of multiple regions)



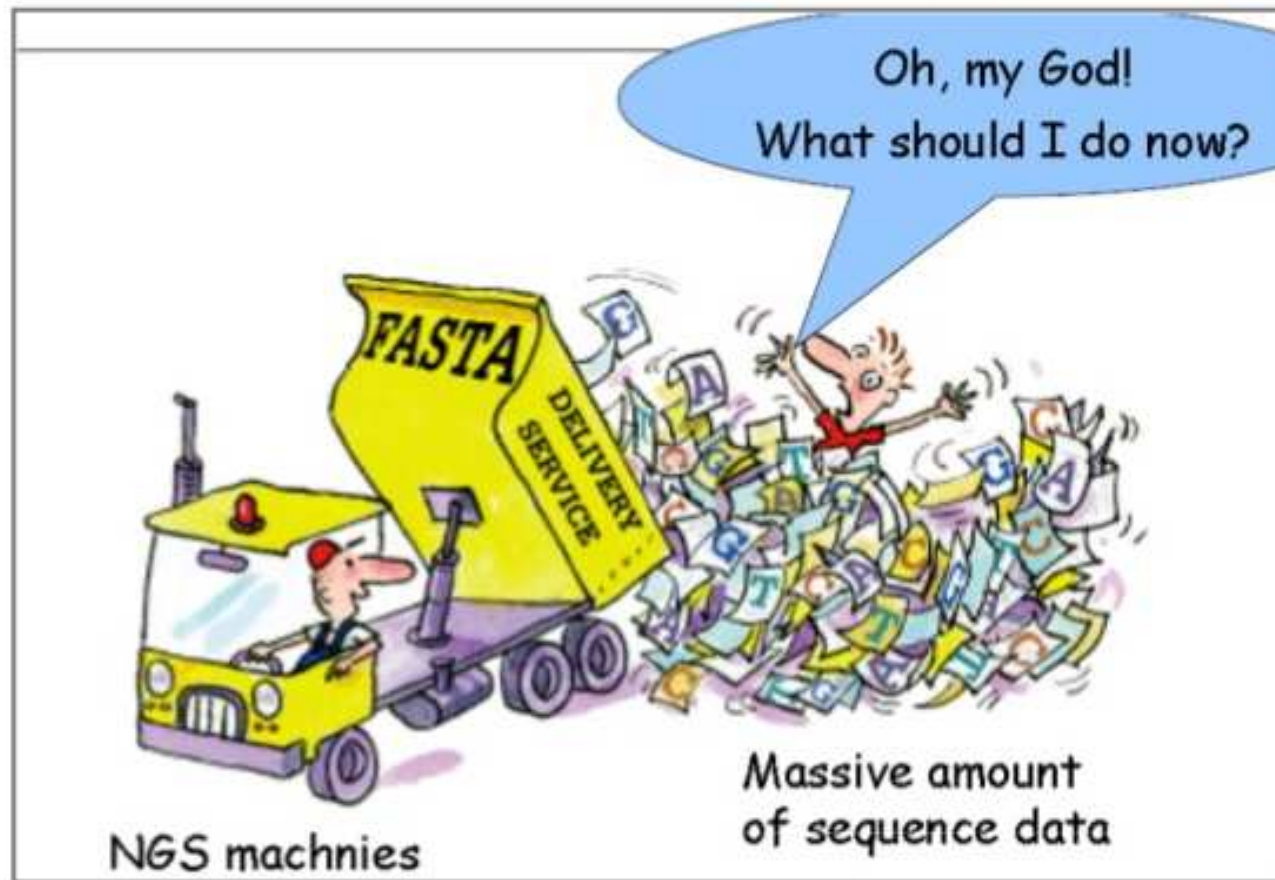


Bottle neck?

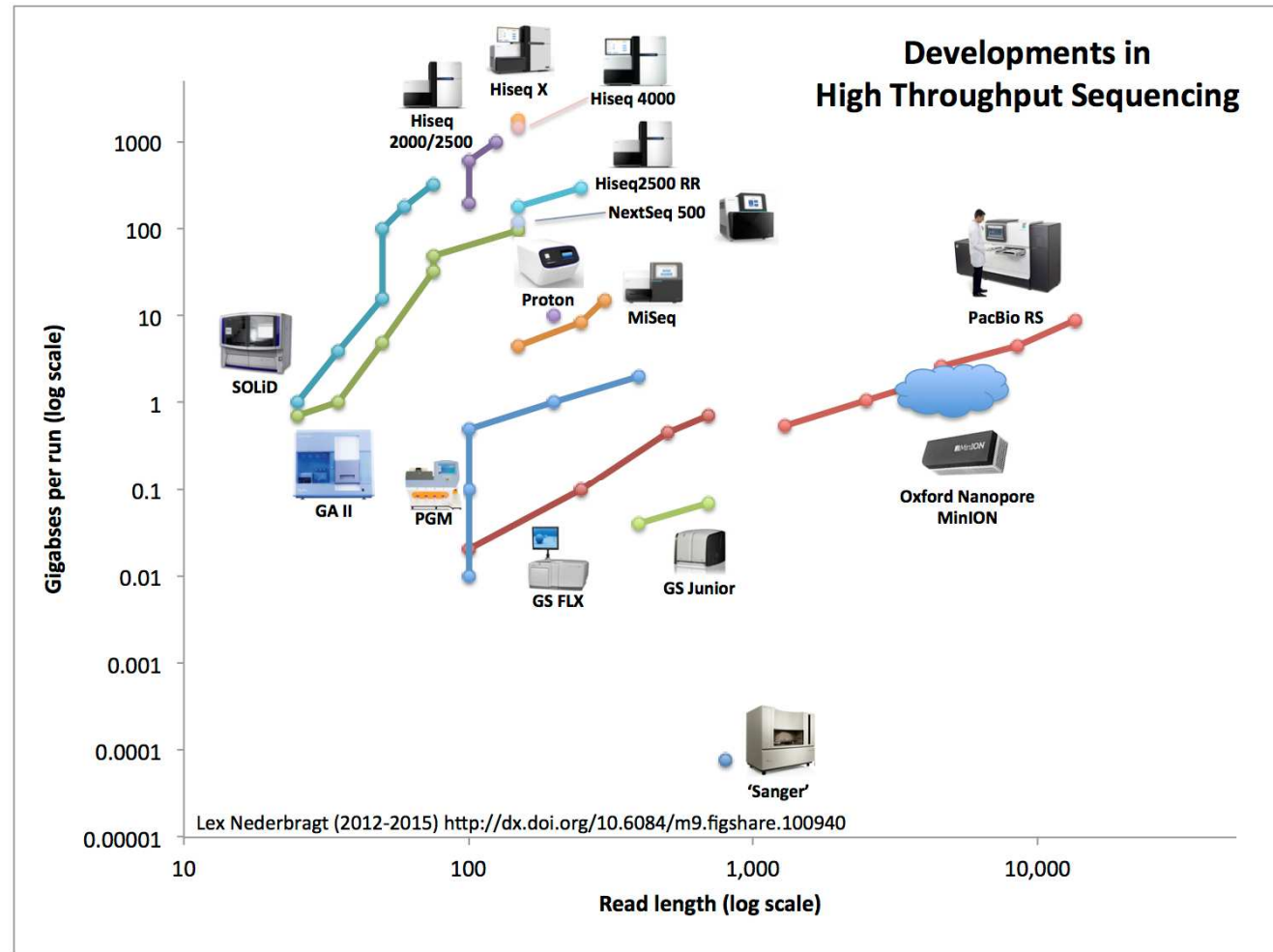
- Sanger
- NGS

- Sampling
- Wet-lab
- Sequencing
- Data procesing

Bottle neck



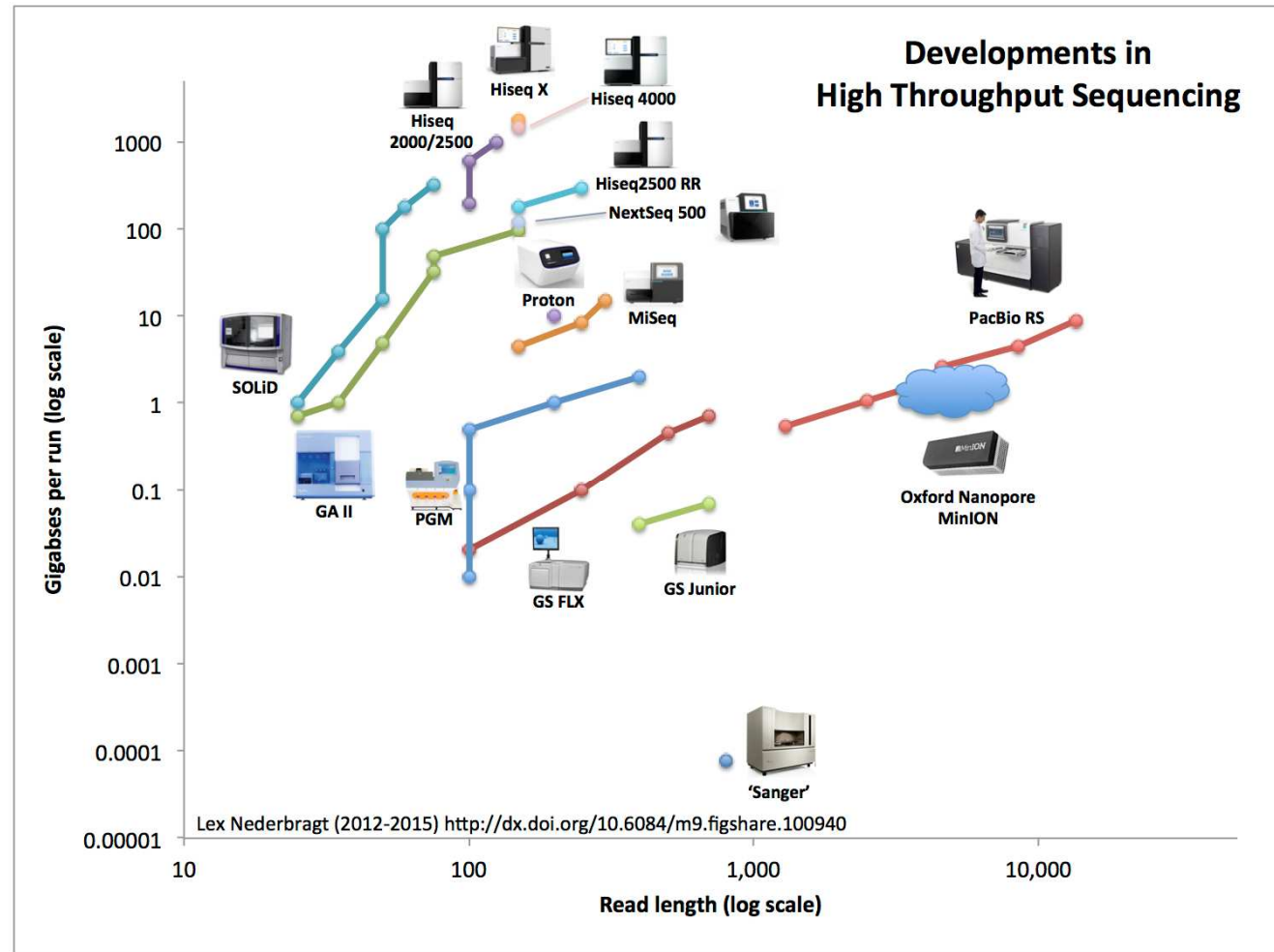
How to choose a method and procedure?



How to choose a method and procedure?

The question is #1

- related species, population, bulk segregant analyses, genomics, transcriptomics, phylogeography,...
- Equipment, protocol availability, data processing and data storage...





And what do you use?

- AFLP
- SSR
- cpDNA
- Single/low copy genes
- RNA
- RadSeq
- HybSeq
- WGS
- PoolSeq
- ...

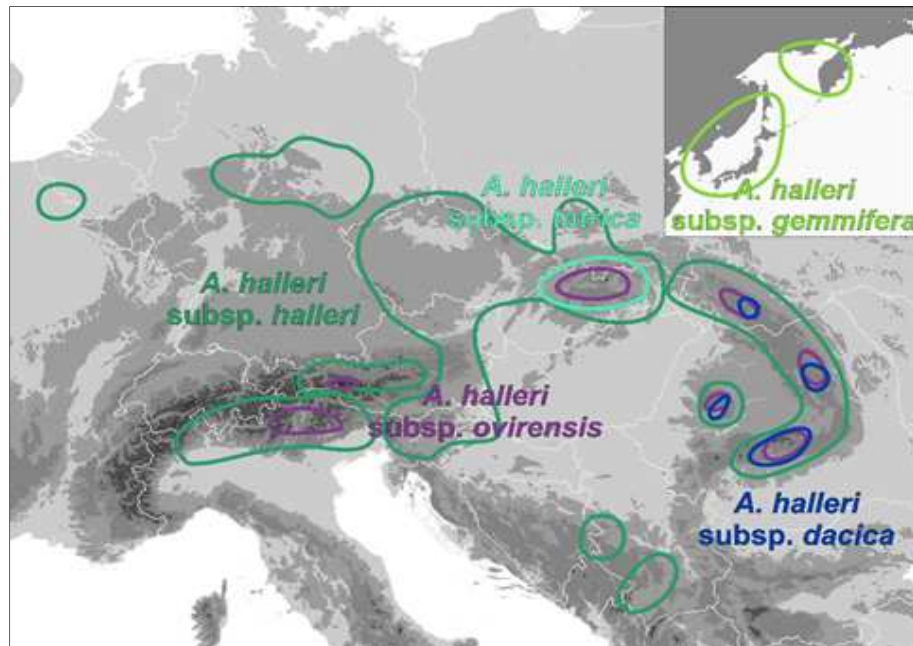
And why?





genus *Arabidopsis*

- ca 14 species worldwide – *A. thaliana* & “wild relatives”



Arabidopsis arenosa complex
distribution and ploidy levels of
described and undescribed taxa

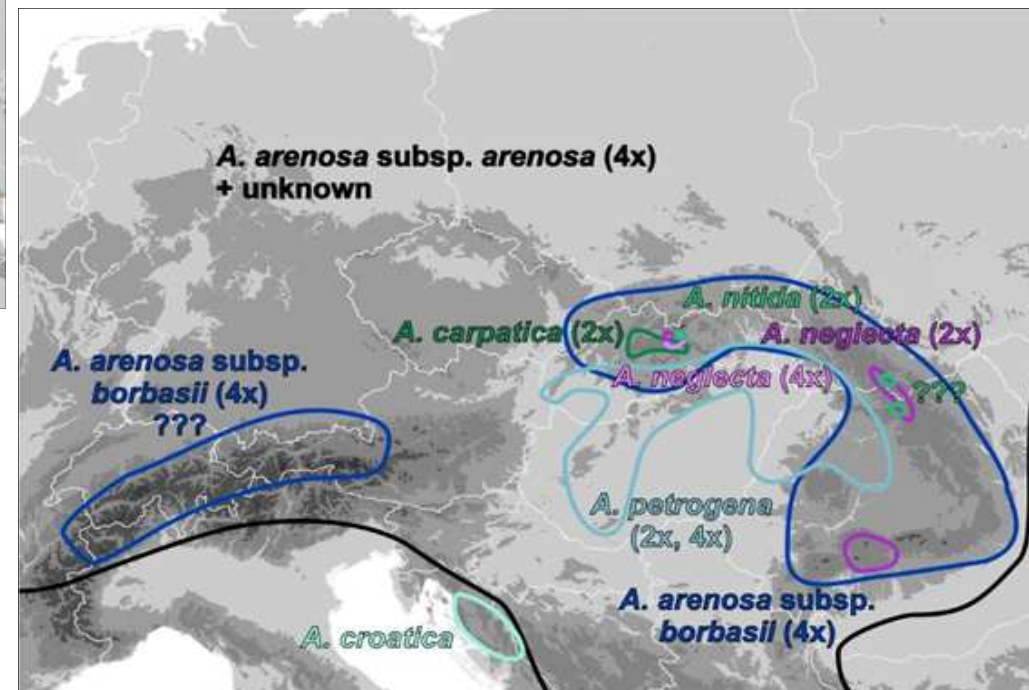
Based on Měsíček 1970

- polyploidy, parallel evolution

Arabidopsis halleri
distribution of described and
undescribed taxa

Based on Měsíček 1970, Kolník and Marhold 2006

- phytoremediation

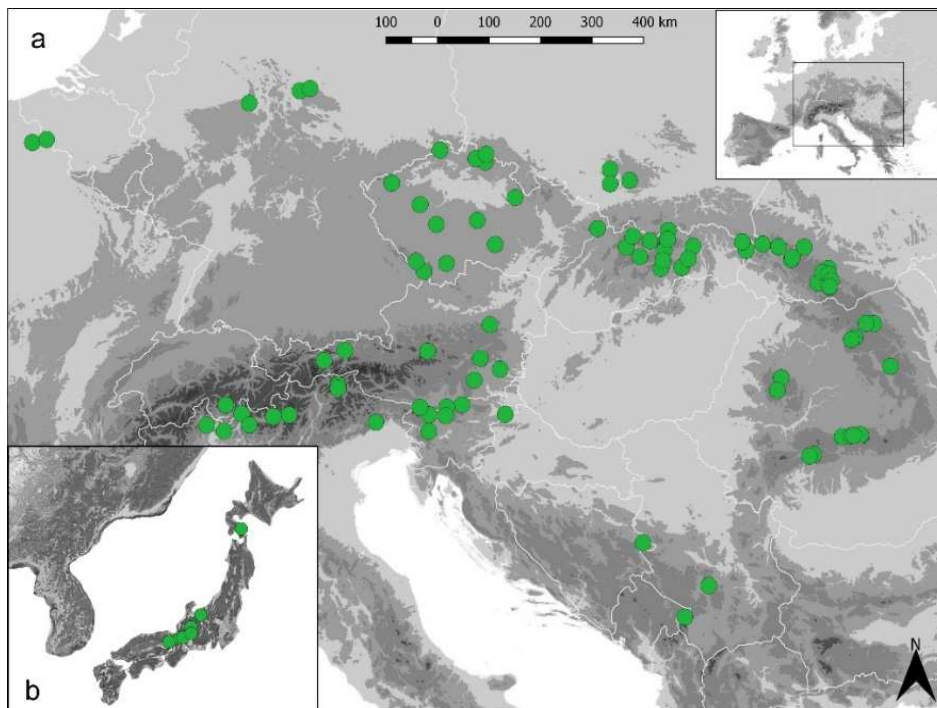




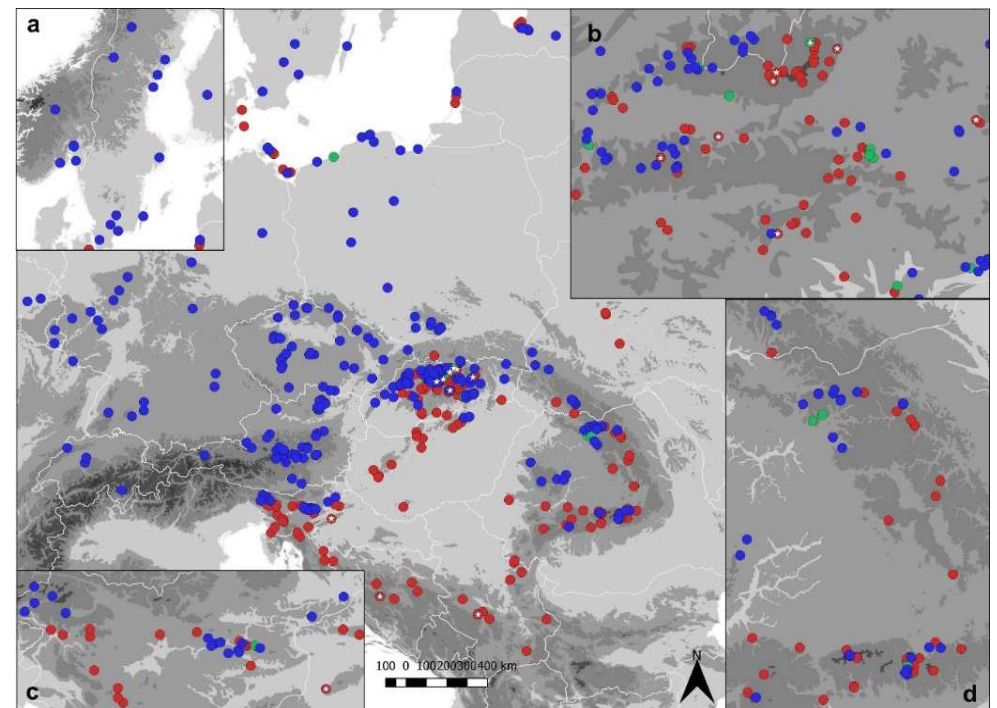
Materials and methods

- homogeneous sampling across the entire distributional range
- flow cytometry
- molecular methods: AFLP, SSR, cpDNA, single copy gene, dd RadSeq, WGS
- multivariate morphometrics

136 sampled populations – *A. halleri*

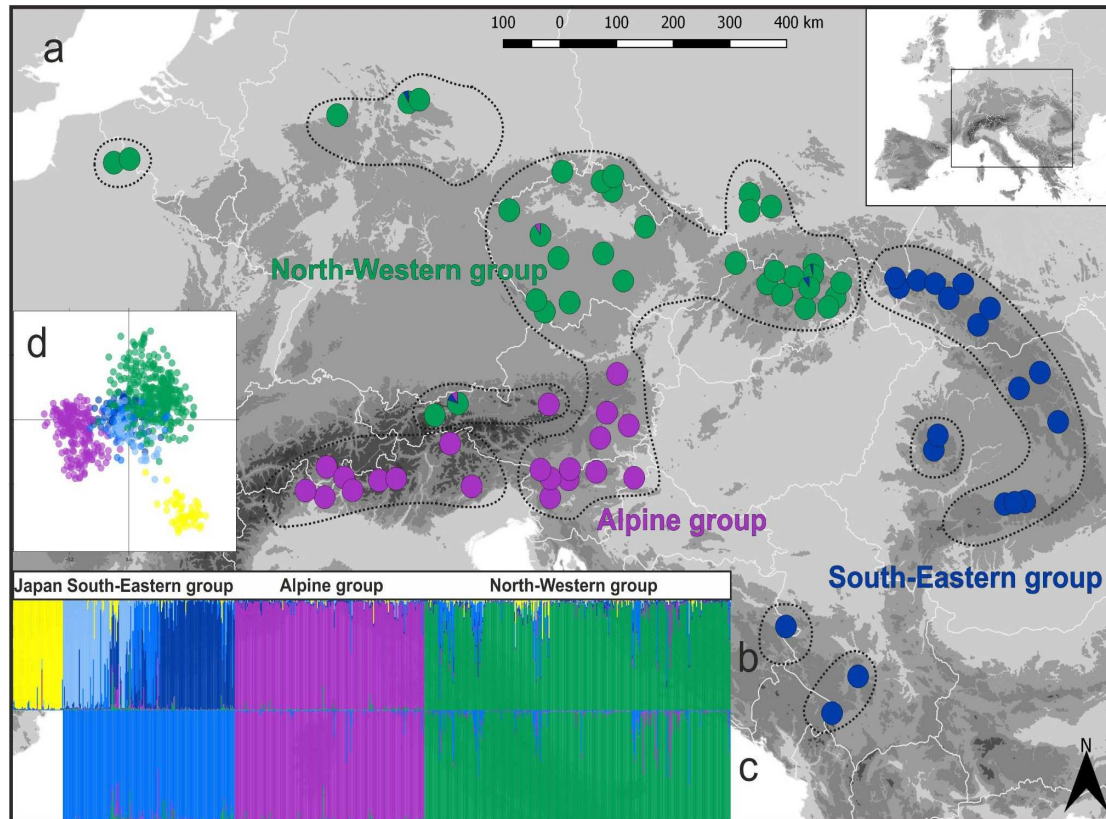


496 sampled populations – *A. arenosa*





Key results: *Arabidopsis halleri*



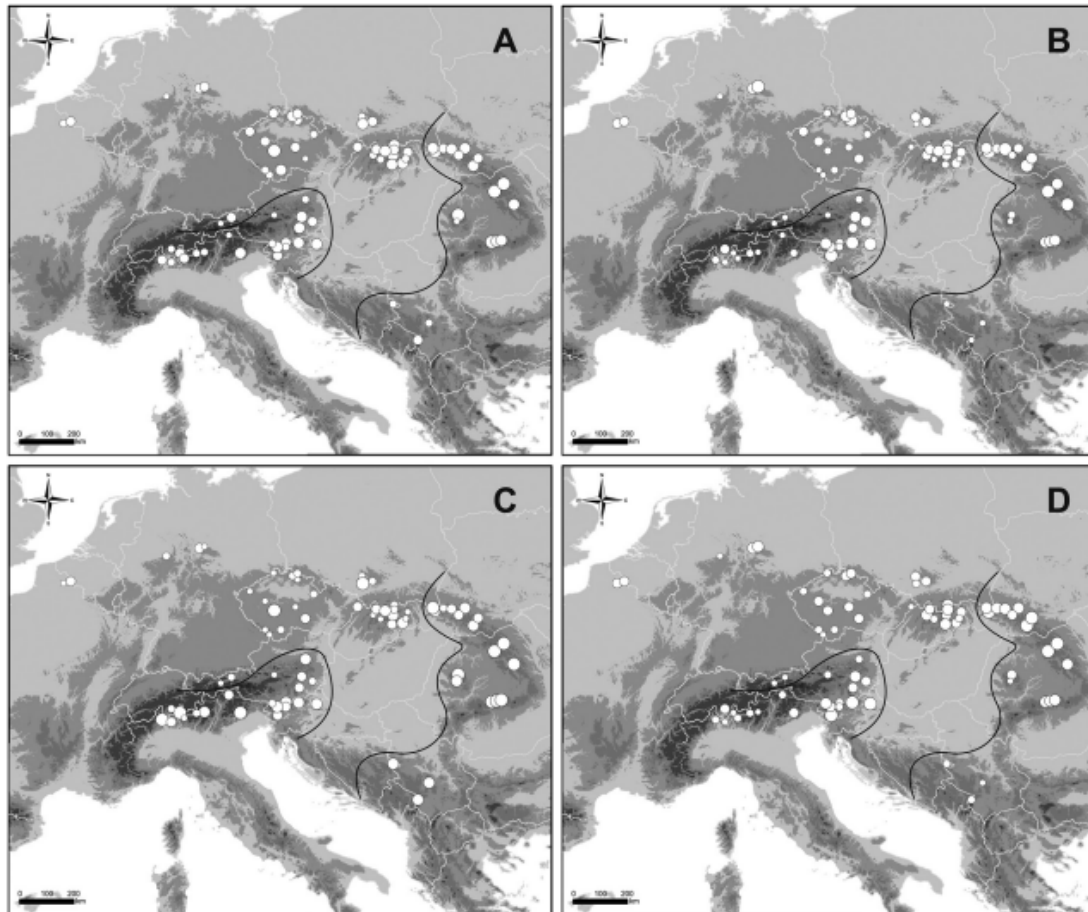
Genetic structure:

- a** geographic distribution (dotted line denotes borders of distribution)
- b** STRUCTURE of Eurasian dataset
- c** STRUCTURE of European dataset
- d** PCoA of Eurasian dataset

- northern hemisphere disjunction (Japanese populations)
- **three main European lineages** – Alpine, North-Western and South-Eastern



Key results: *Arabidopsis halleri*



Population-level diversity:

a gene diversity AFLP

b expected heterozygosity

SSR

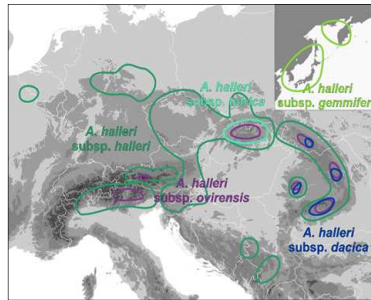
c proportion of rare
fragments AFLP

d allelic richness SSR

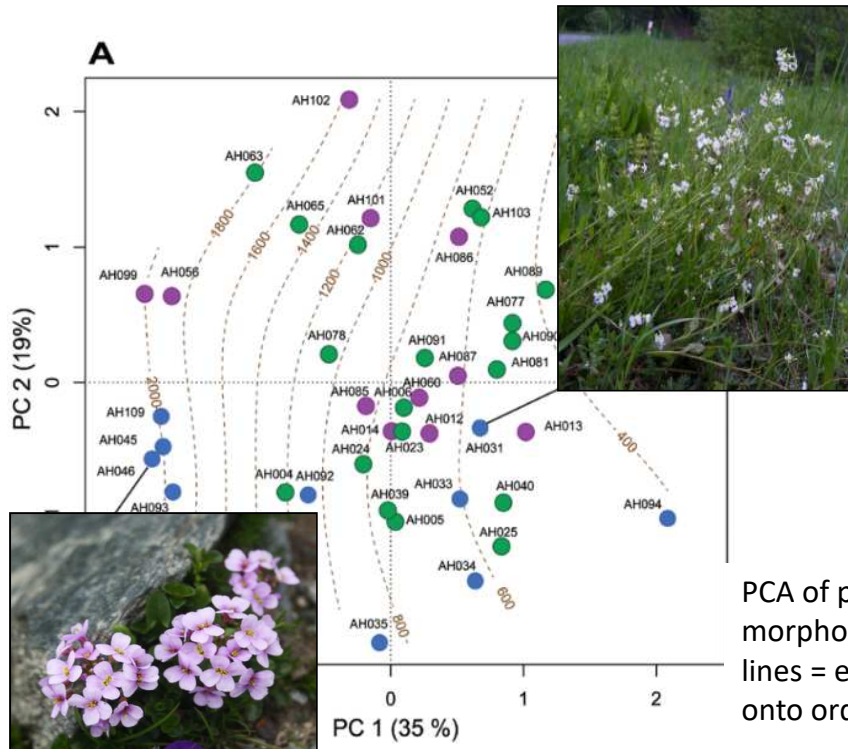
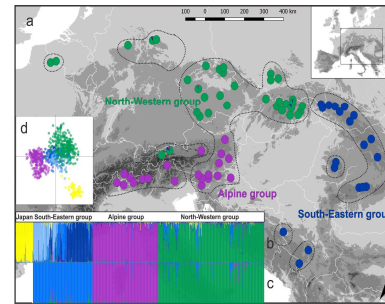
- the highest diversity and representation of rare alleles are in **W** and **SE Carpathians**



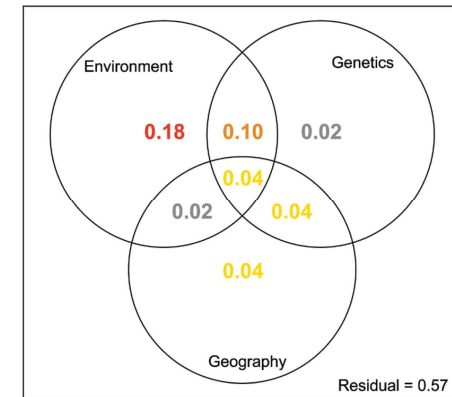
Key results: *Arabidopsis halleri*



≠



PCA of population morphology (brown lines = elevation fitted onto ordination plot)



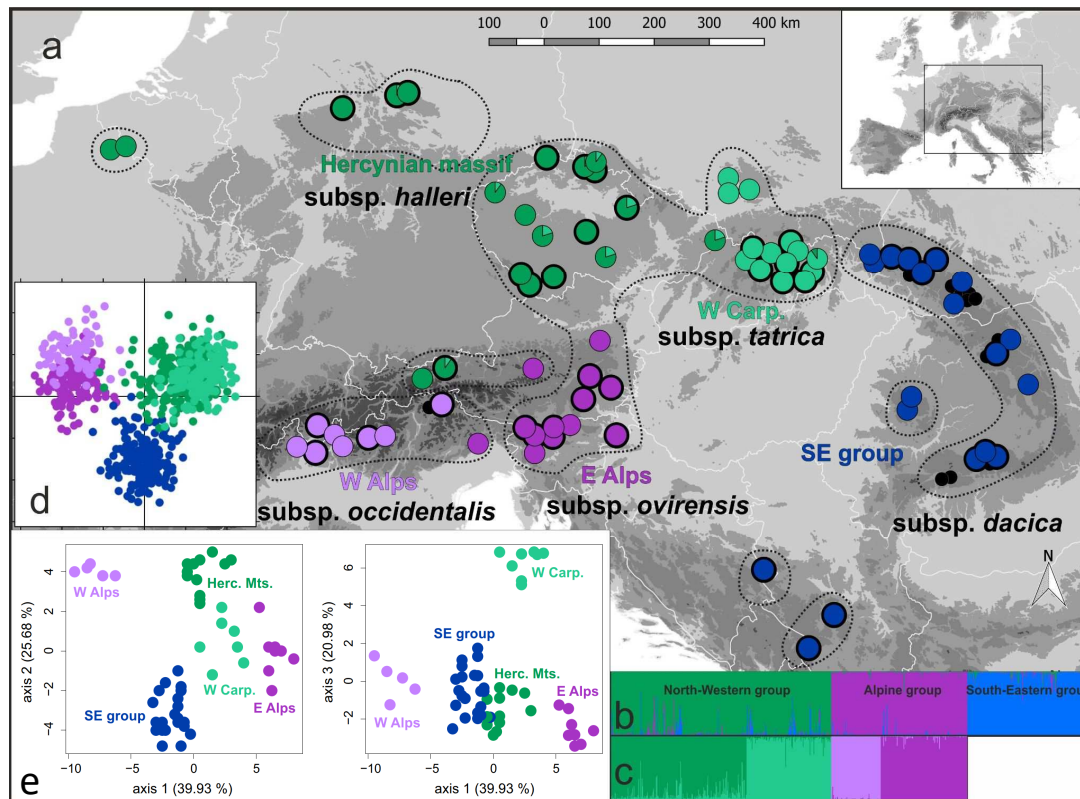
relative contributions of parameters to morphological variation (variation partitioning)

Genetic divergence is not a source of morphological variation...



Key results: *Arabidopsis halleri*

- detailed phylogeographic structure tested for morphological separation

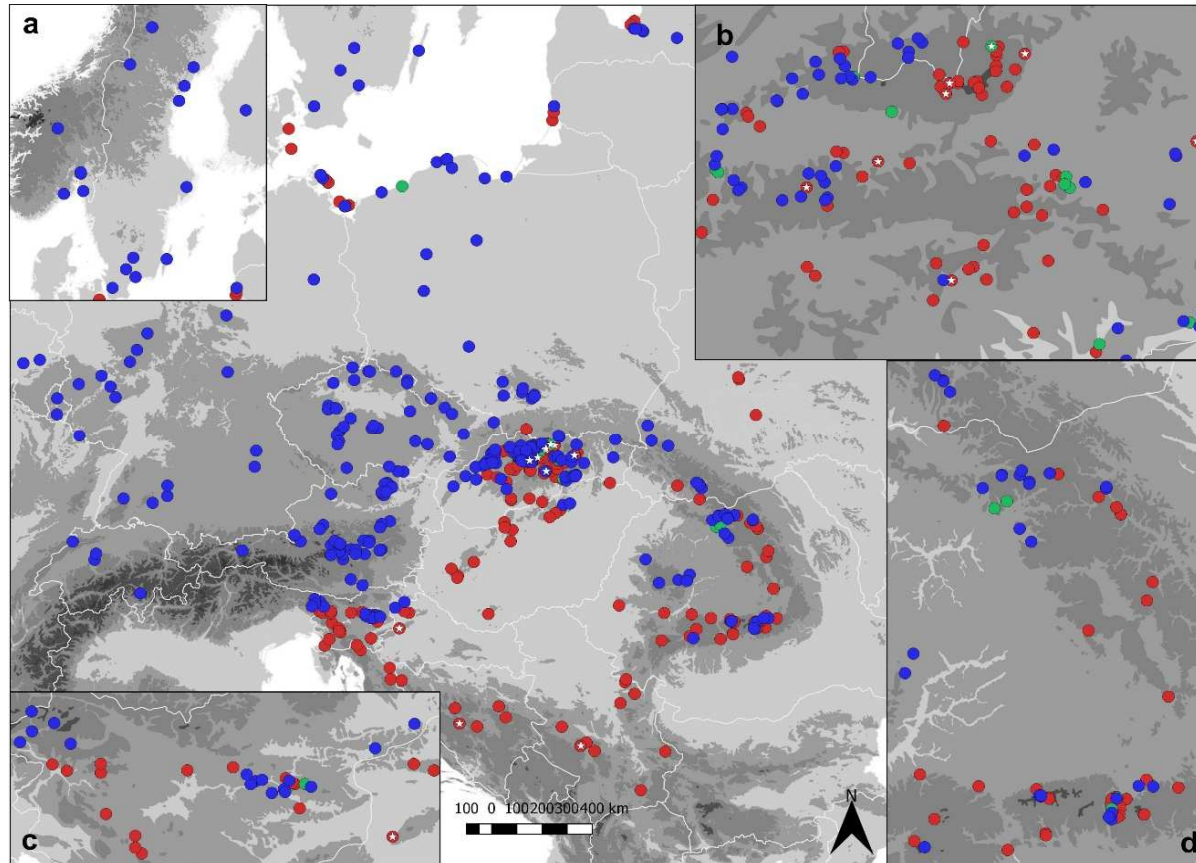


Taxonomic reassessment

- a subgroups geographic distribution
- b STRUCTURE EU dataset
- c separate STRUCTURE for lineages
- d PCoA of AFLP phenotypes
- e morphological separation of 5 subgroups (canonical discriminant analyses)

- five subgroups morphologically separated → taxonomic re-evaluation

Key results: *Arabidopsis arenosa* complex



Distribution and ploidy level

red – diploid

blue – tetraploid

green – mixedploidy

asterisk – triploid inds in population

a Scandinavia

b W Carp cont. zone

c Slovenia cont. zone

d SE cont. zone

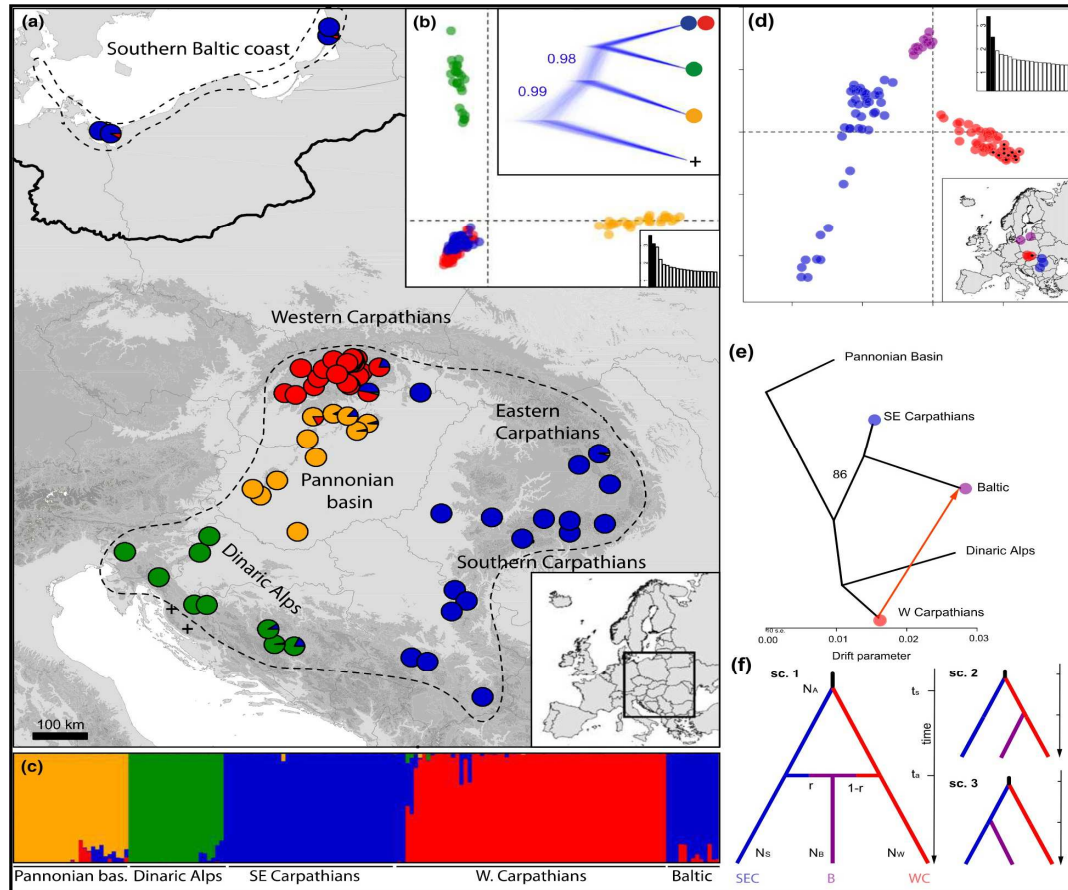
- 3 different ploidy levels – 2x, 3x 4x
- distinct 2x populations on Baltic Sea coast
- 3 contact zones



Kolář et al. 2016 BiolJLinnSoc + unpub.



Key results: *Arabidopsis arenosa* 2x



Genetic structure:

- a geographic distribution (dotted line denotes borders of 2x distribution)
- b PCoA and species tree (rooted with *A. croatica*)
- c STRUCTURE clustering

Reconstruction of Baltic-Carpathian relationships:

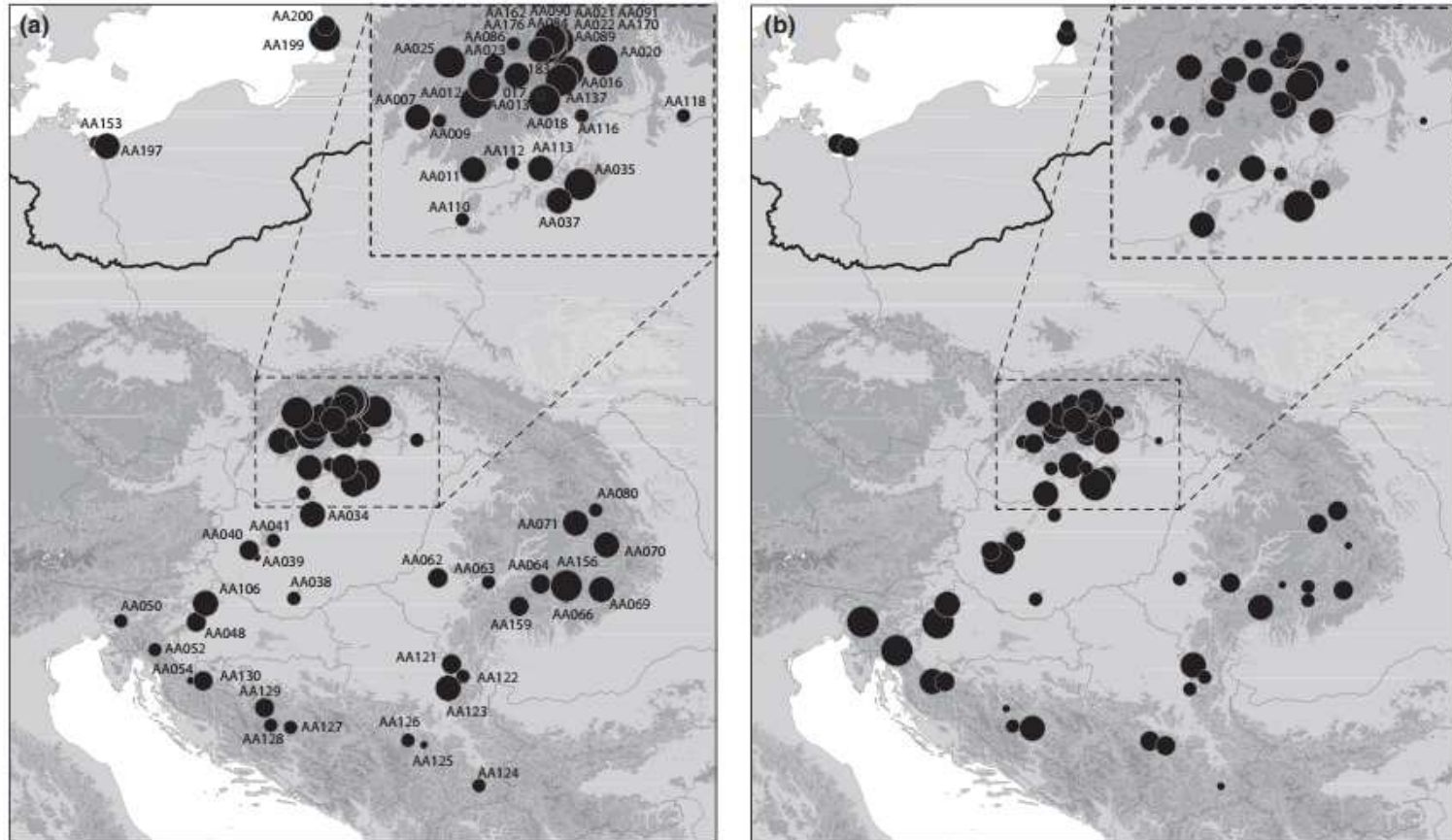
- d PCoA of Baltic and Carpathian individuals
- e graph of allele frequency covariance with admixture (Treemix)
- f ABC modelling of Baltic origin

Kolář, Fuxová, Závěská, et al. 2016 MolEcol

- 4 divergent European lineages – Dinaric, Pannonian, W Carpathian, SE Carpathian + spatially isolated Baltics
- Baltics originated from SE and W Carpathian lineages



Key results: *Arabidopsis arenosa* 2x

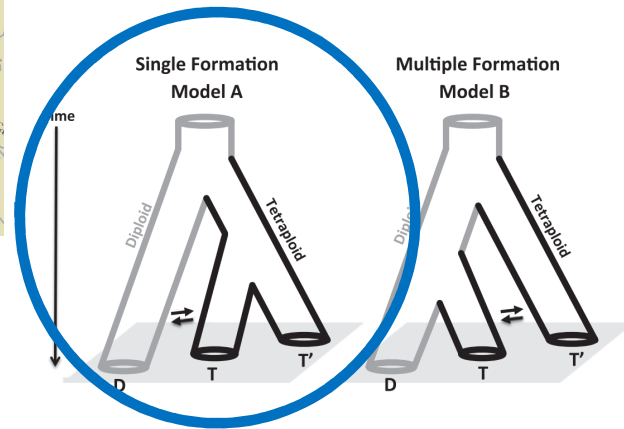
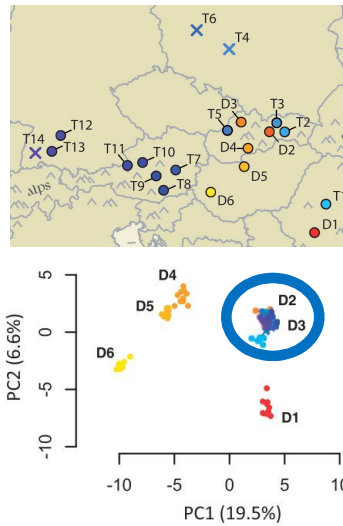
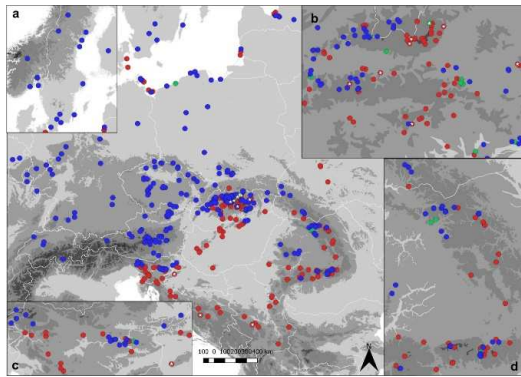


Population-level diversity (SSR): a expected heterozygosity b proportion of rare alleles (DW index)

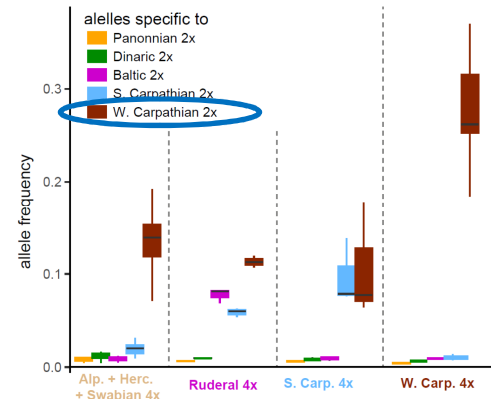
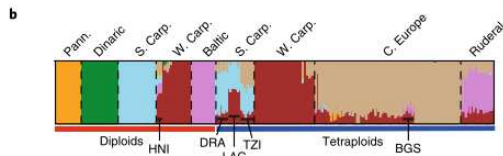
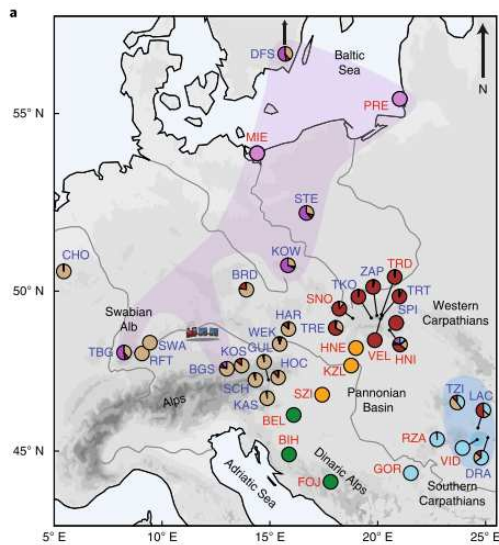
- the highest diversity and representation of rare alleles are in W Carpathians



Arabidopsis arenosa 4x



Arnold et al. 2015 Mol Biol Evol

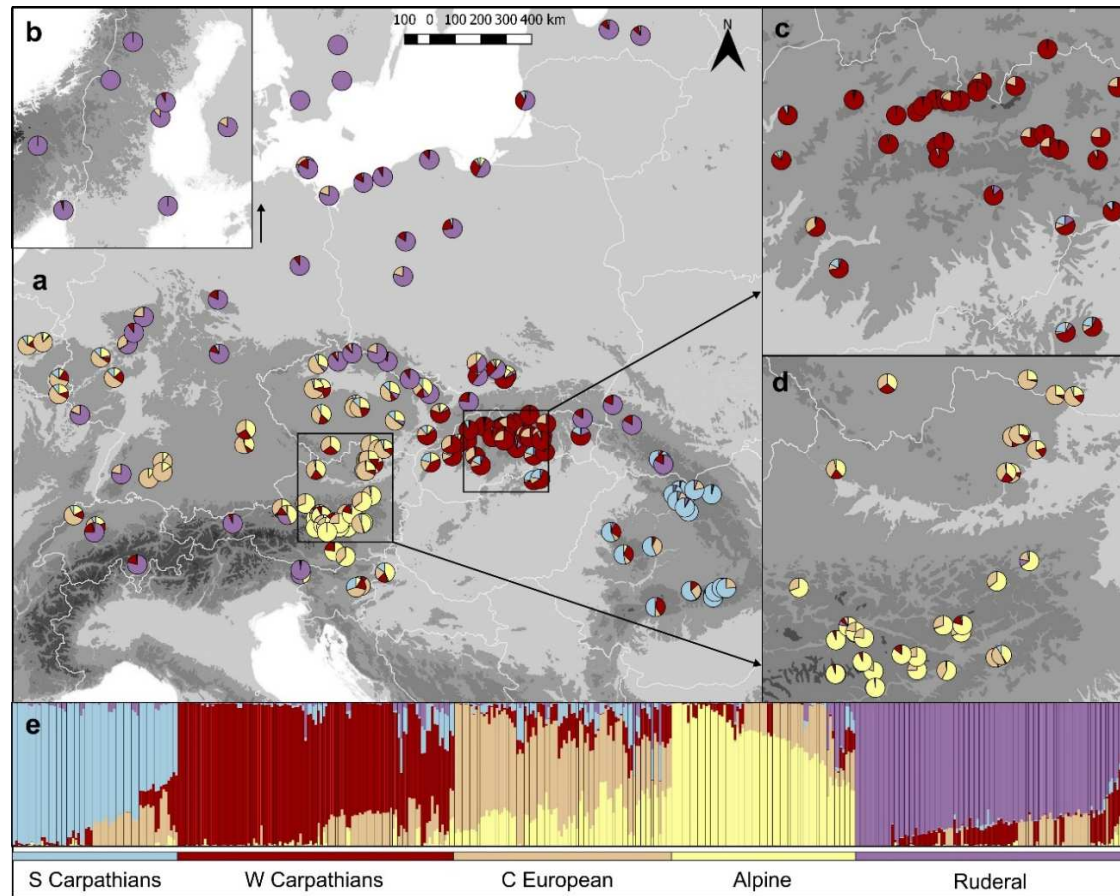


- single geographic origin of autotetraploid lineages

Monnahan et al. 2018 Nature Ecol Evo



Key results: *Arabidopsis arenosa* 4x



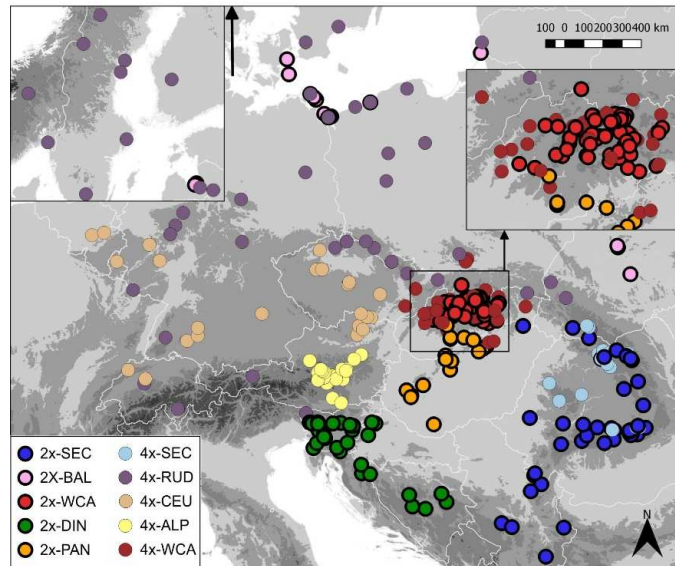
Geographic distribution of genetic structure of 4x populations

- a sampled populations
- b Scandinavia
- c W Carpathians
- d Austrian zone
- e individual STRUCTURE assignment

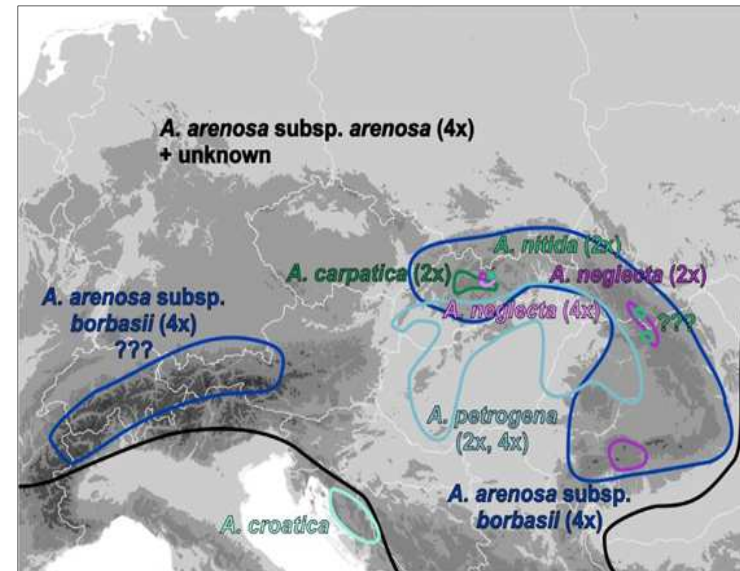
- five genetic clusters
- Ruderal lineage – anthropogenic niches, colonized higher latitudes



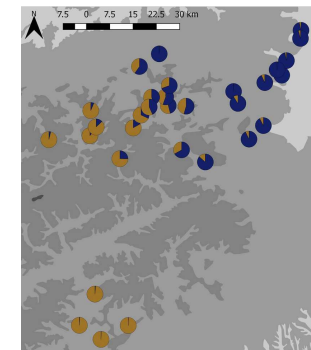
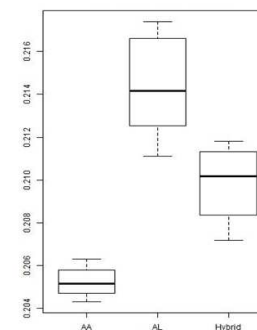
Future prospects



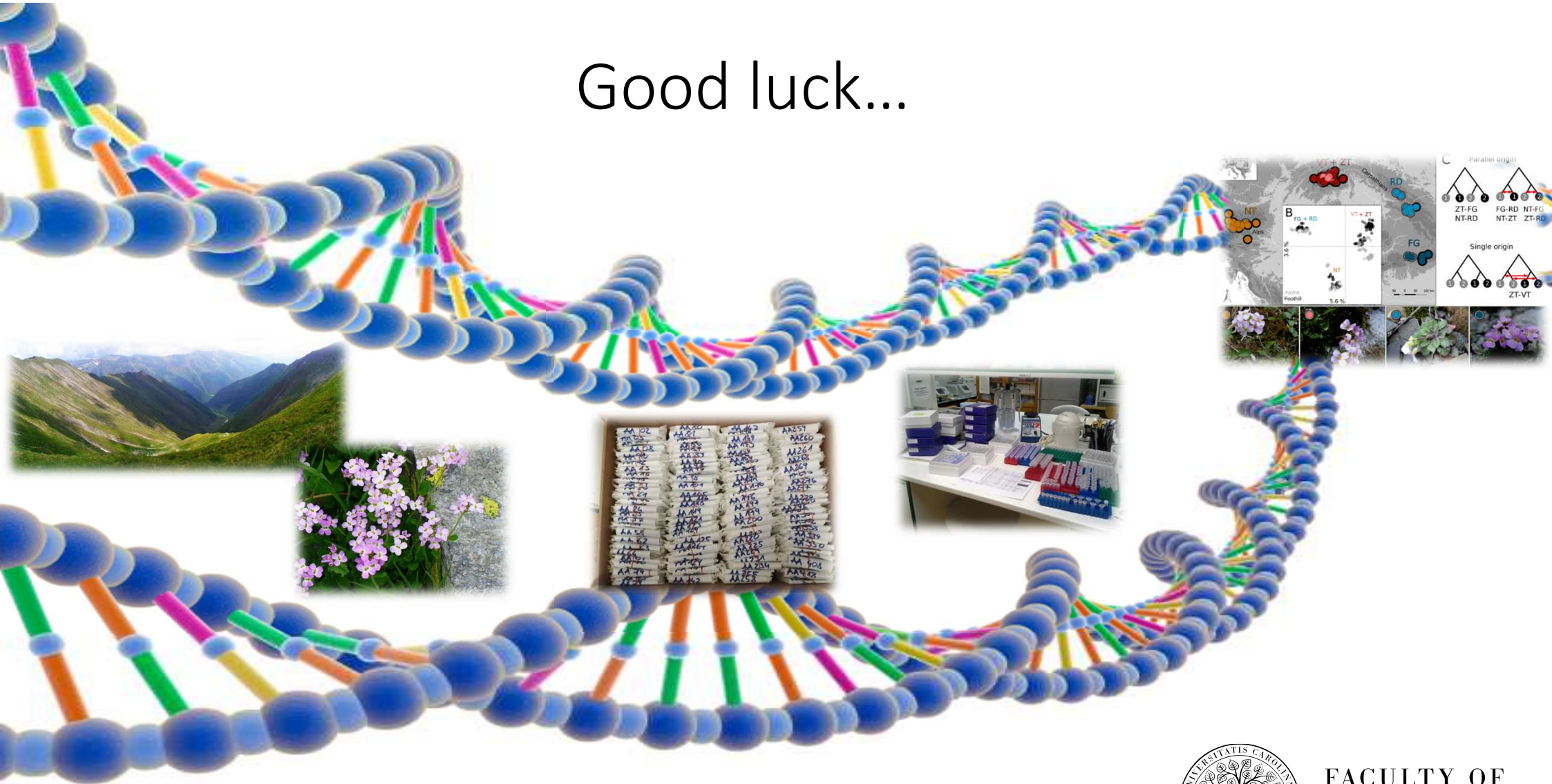
≠



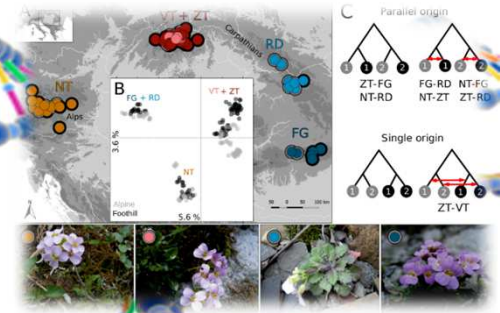
- taxonomical reassessment of *A. arenosa*
 - analyses of the whole dataset
 - morphometrics
- hybridization *A. arenosa* – *A. lyrata*
- ...



Good luck...



AA 02	AA 30	AA 42	AA 51
AA 03	AA 31	AA 43	AA 52
AA 04	AA 32	AA 44	AA 53
AA 05	AA 33	AA 45	AA 54
AA 06	AA 34	AA 46	AA 55
AA 07	AA 35	AA 47	AA 56
AA 08	AA 36	AA 48	AA 57
AA 09	AA 37	AA 49	AA 58
AA 10	AA 38	AA 50	AA 59
AA 11	AA 39	AA 51	AA 60
AA 12	AA 40	AA 52	AA 61
AA 13	AA 41	AA 53	AA 62
AA 14	AA 42	AA 54	AA 63
AA 15	AA 43	AA 55	AA 64
AA 16	AA 44	AA 56	AA 65
AA 17	AA 45	AA 57	AA 66
AA 18	AA 46	AA 58	AA 67
AA 19	AA 47	AA 59	AA 68
AA 20	AA 48	AA 60	AA 69
AA 21	AA 49	AA 61	AA 70
AA 22	AA 50	AA 62	AA 71
AA 23	AA 51	AA 63	AA 72
AA 24	AA 52	AA 64	AA 73
AA 25	AA 53	AA 65	AA 74
AA 26	AA 54	AA 66	AA 75
AA 27	AA 55	AA 67	AA 76
AA 28	AA 56	AA 68	AA 77
AA 29	AA 57	AA 69	AA 78
AA 30	AA 58	AA 70	AA 79
AA 31	AA 59	AA 71	AA 80
AA 32	AA 60	AA 72	AA 81
AA 33	AA 61	AA 73	AA 82
AA 34	AA 62	AA 74	AA 83
AA 35	AA 63	AA 75	AA 84
AA 36	AA 64	AA 76	AA 85
AA 37	AA 65	AA 77	AA 86
AA 38	AA 66	AA 78	AA 87
AA 39	AA 67	AA 79	AA 88
AA 40	AA 68	AA 80	AA 89
AA 41	AA 69	AA 81	AA 90
AA 42	AA 70	AA 82	AA 91
AA 43	AA 71	AA 83	AA 92
AA 44	AA 72	AA 84	AA 93
AA 45	AA 73	AA 85	AA 94
AA 46	AA 74	AA 86	AA 95
AA 47	AA 75	AA 87	AA 96
AA 48	AA 76	AA 88	AA 97
AA 49	AA 77	AA 89	AA 98
AA 50	AA 78	AA 90	AA 99
AA 51	AA 79	AA 91	AA 100



FACULTY OF
SCIENCE
Charles University



A few methods of library preparation...

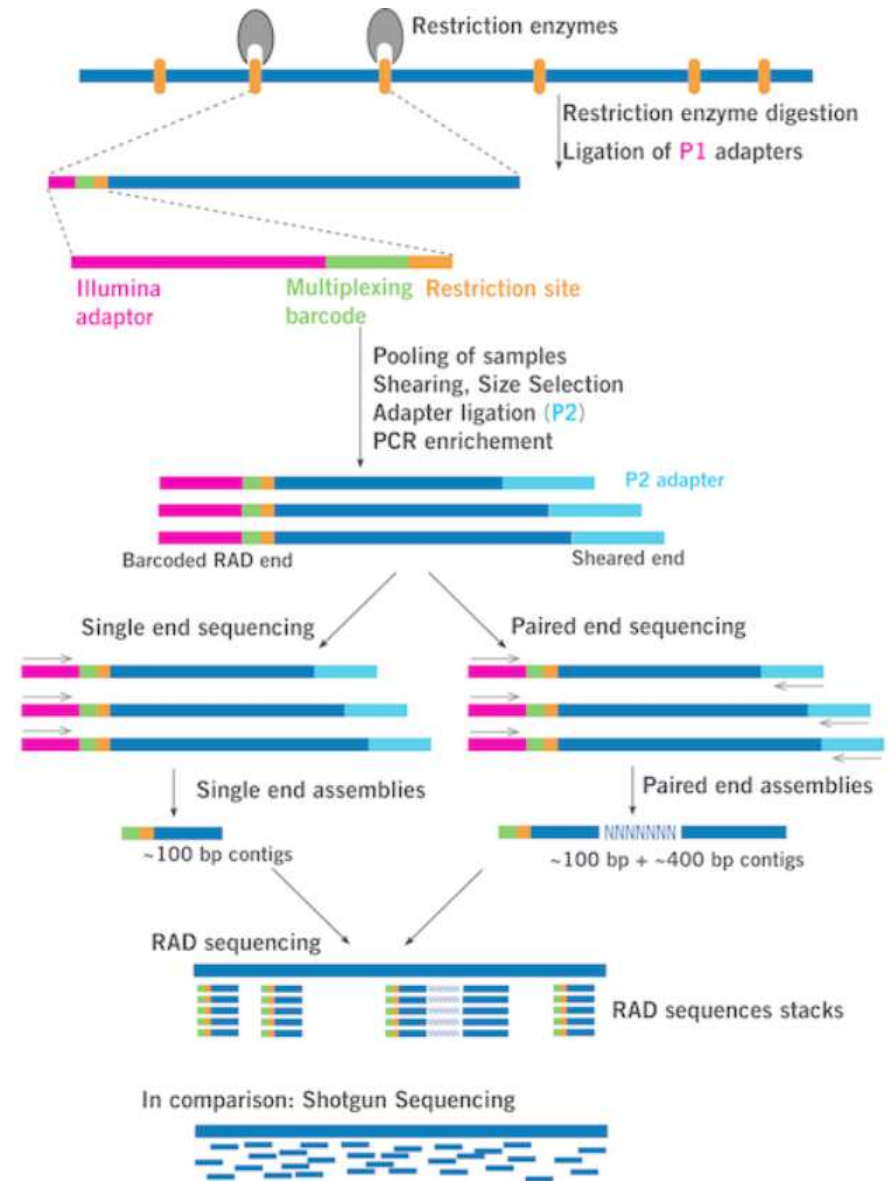
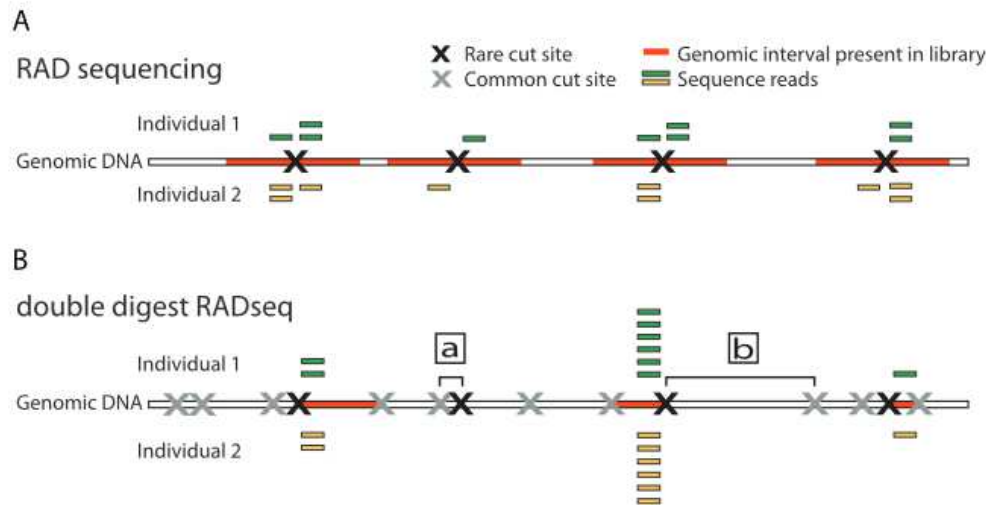
- RadSeq
- HybSeq
- WGS
- RNASeq
- HyRAD
- MigSeq
- SSR-GBS



RadSeq

Restriction (enzyme, sonication)

Radseq, ddRadseq



HybSeq

Custom-made proby (baits)

Genomická knihovna

Hybridizace

Sekvence obohacených fragmentů (+
neobohacených částí genomu)

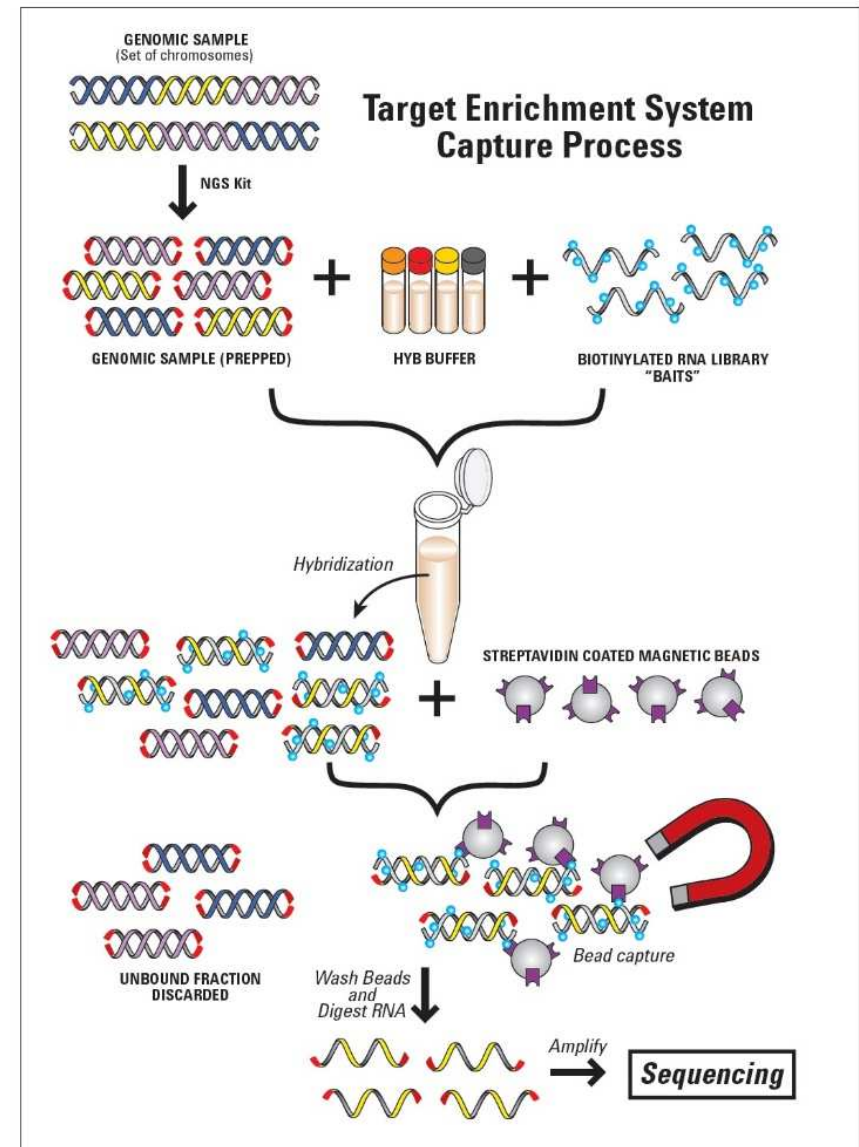
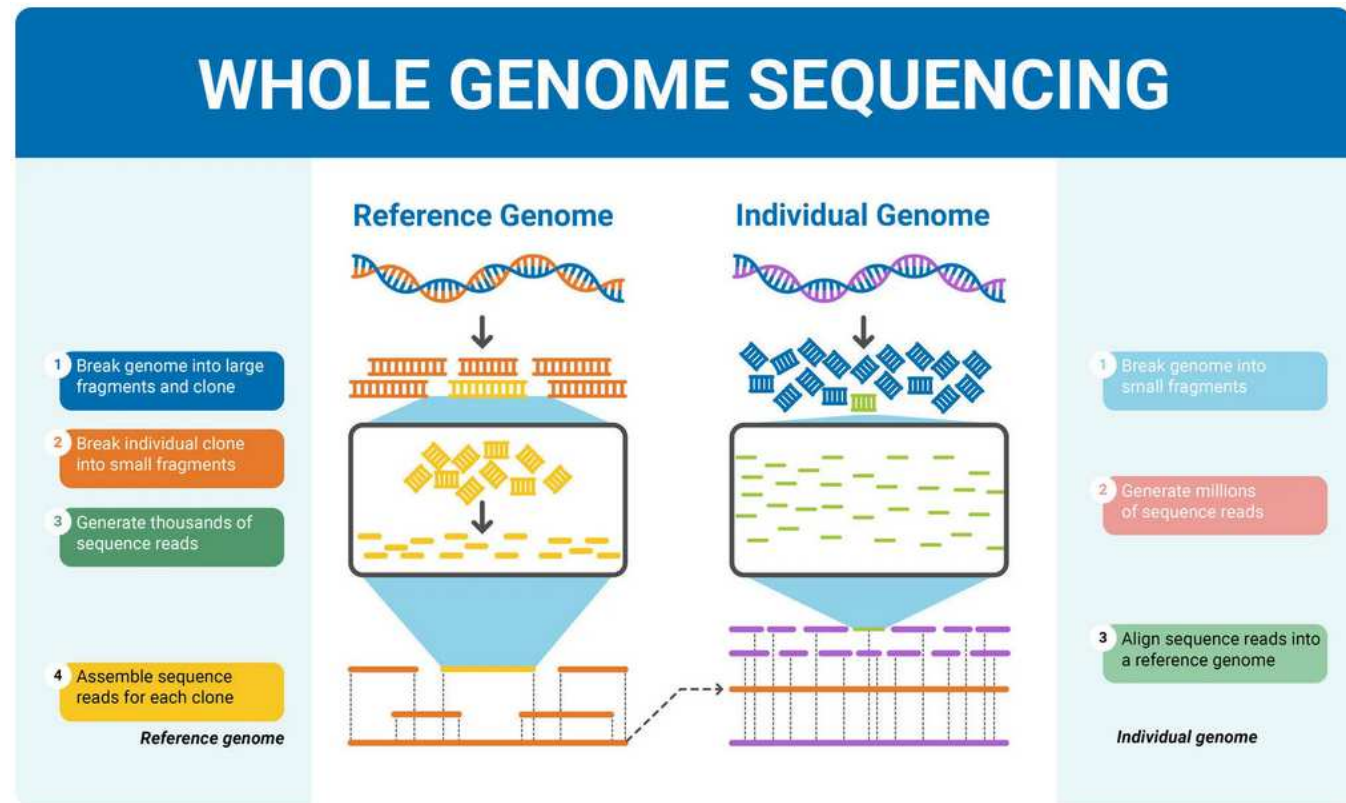


Figure 1 Target enrichment system workflow

Whole Genome Sequencing

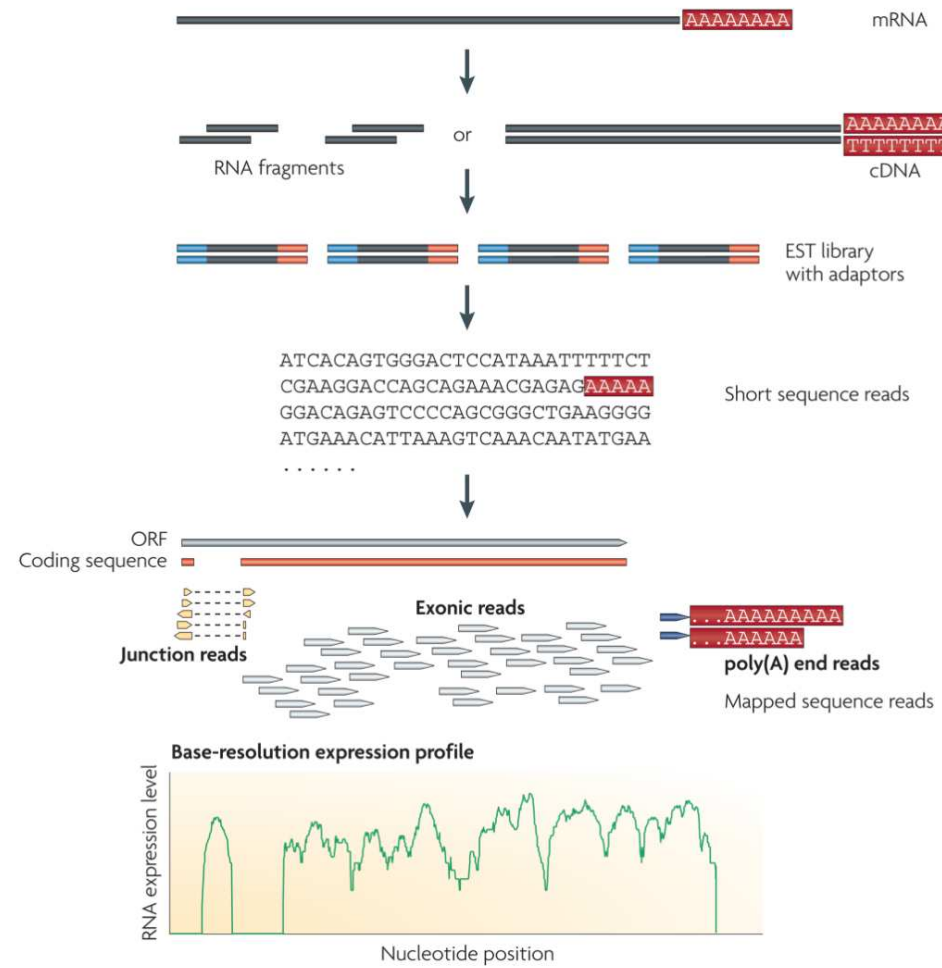
- Metodou služby (kity...)
- Home-made protokoly
 - LITE (Rowan et al. Genetics 2019) - transpozom





RNA Seq

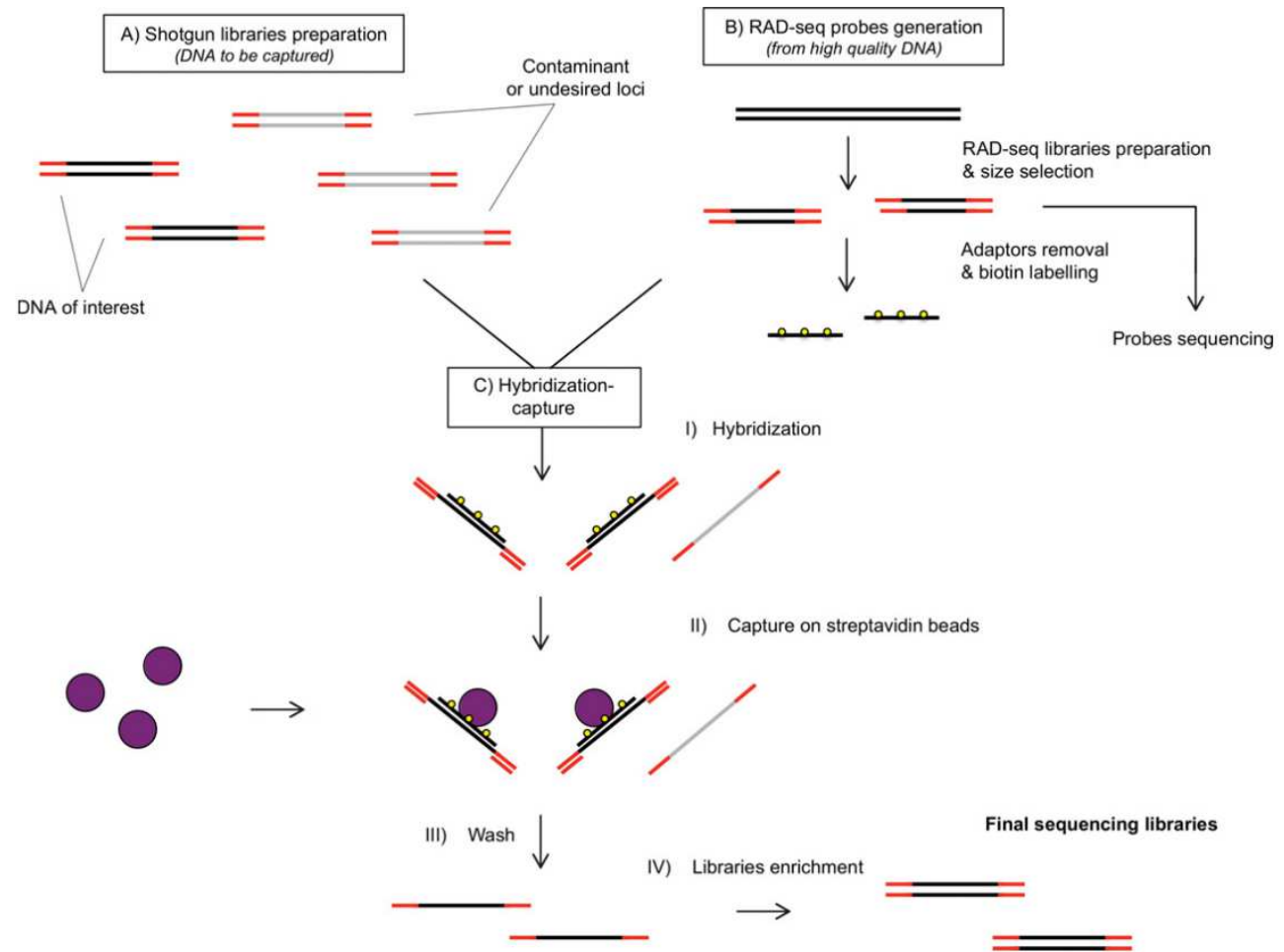
Izolace RNA





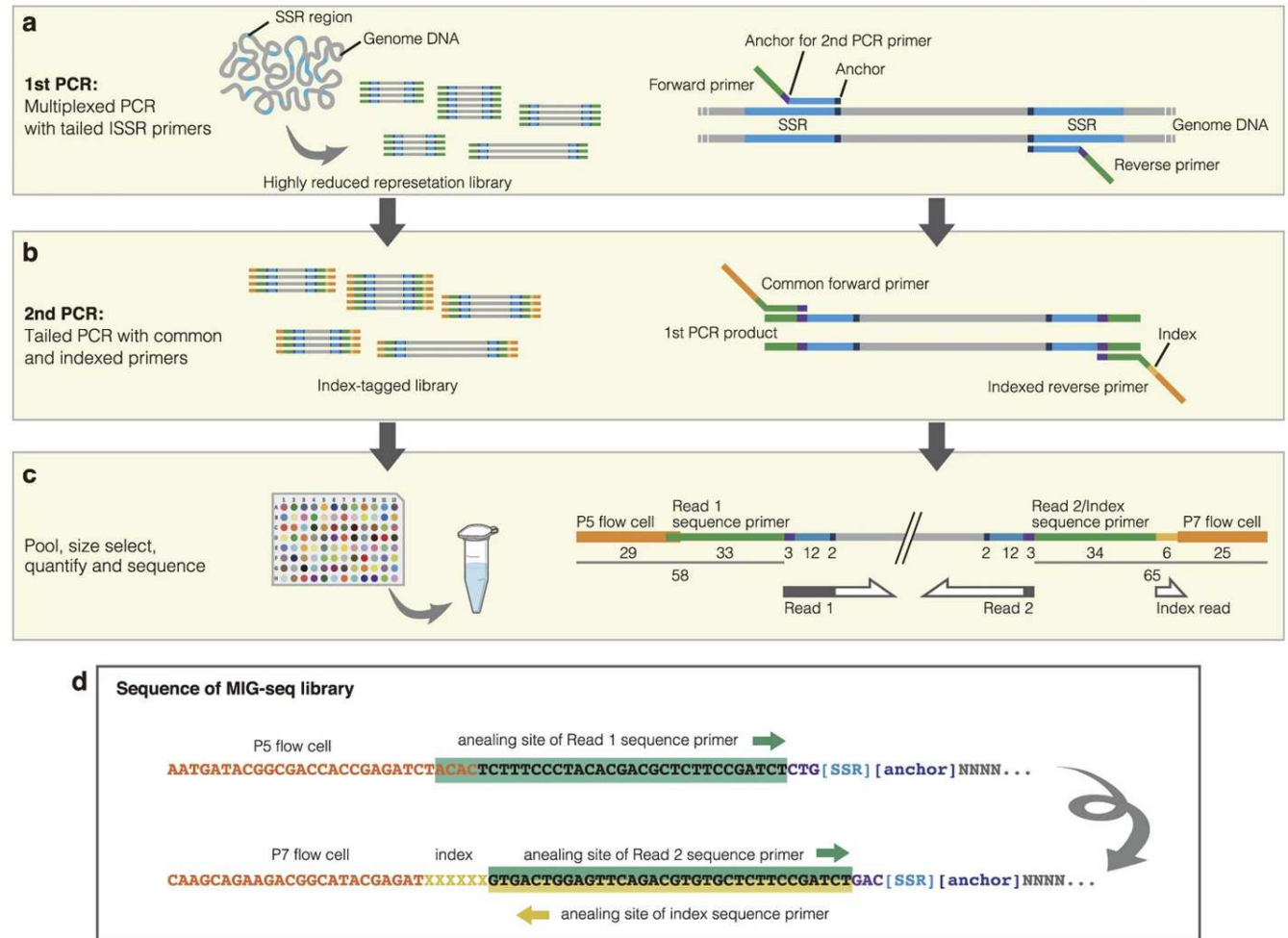
HyRad

Kombinace RadSeq + target enrichment



MigSeq

Restrikce (enzym - Inter-simple sequence repeats)





SSR-GBS

4-primerová PCR

