



Experimental Design

Joana Meier

Tree of Life Programme, Wellcome Sanger Institute

**THE
ROYAL
SOCIETY**

The
Branco Weiss
Fellowship
Society in Science

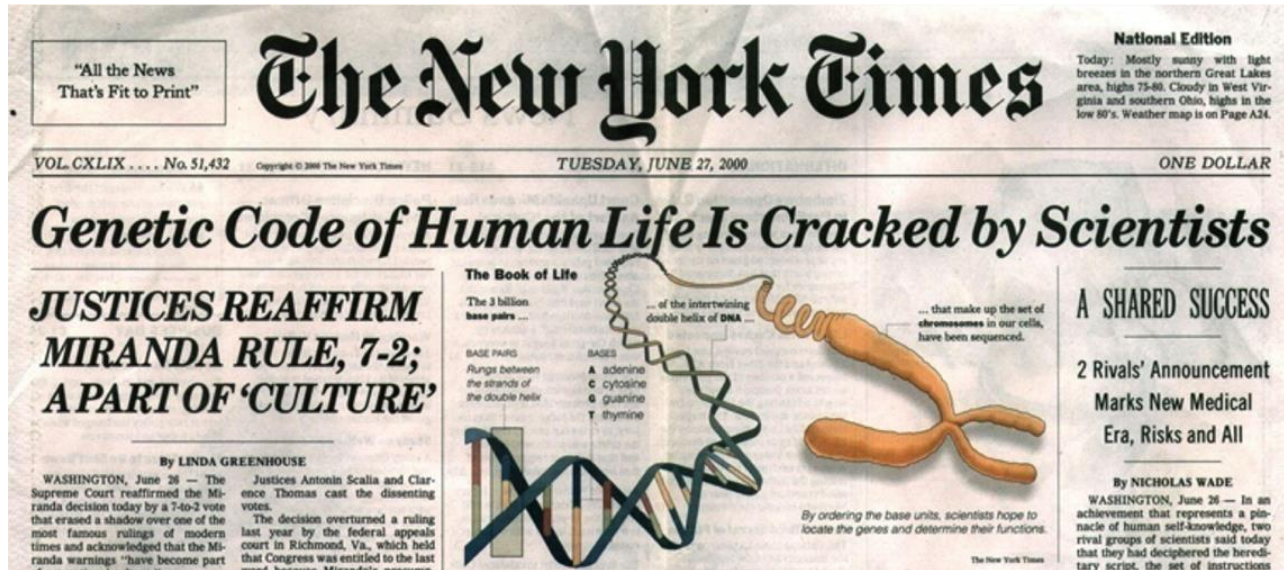


**UNIVERSITY OF
CAMBRIDGE**

Outline

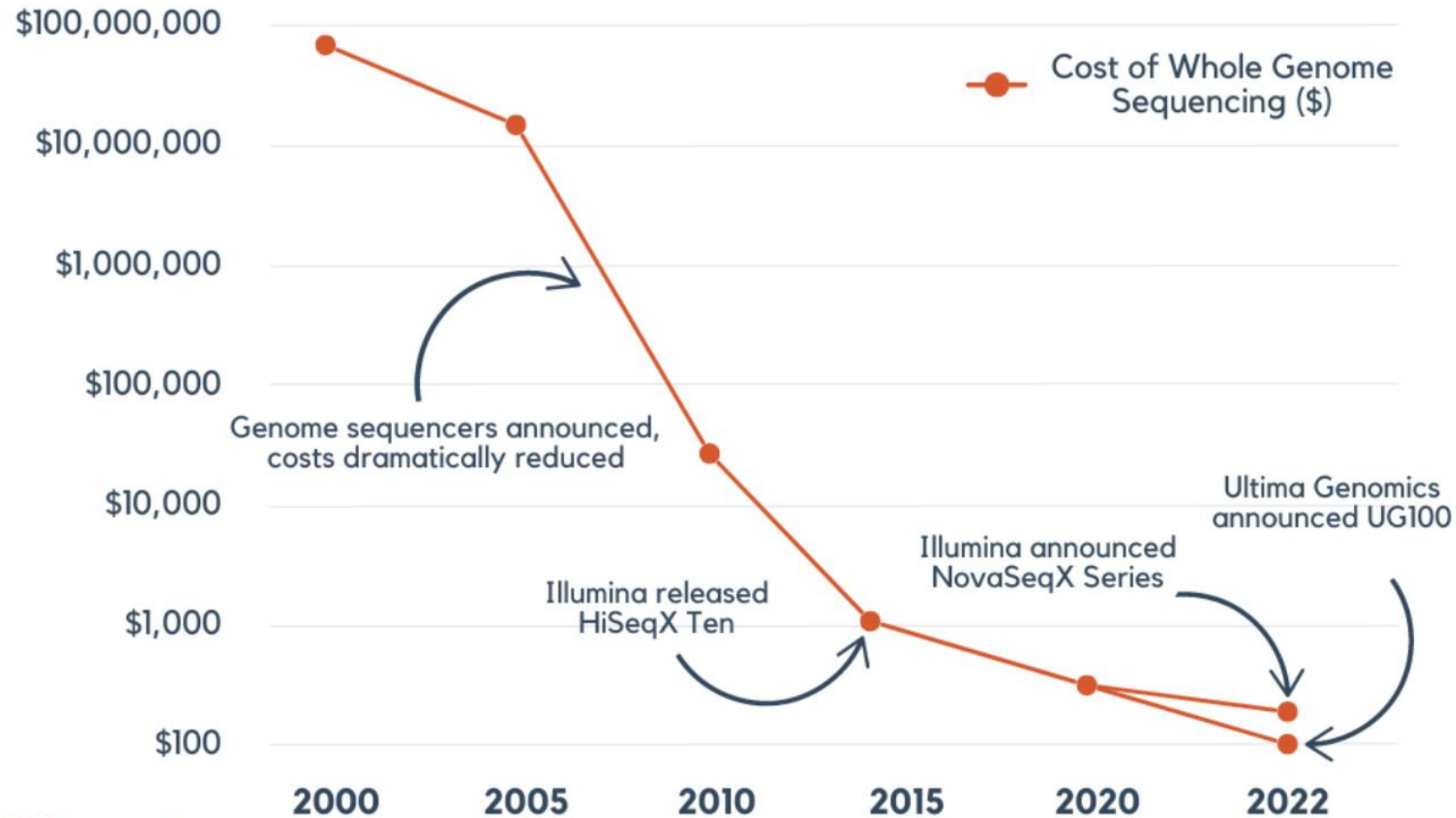
- Sequencing technologies
- Sampling and sequencing strategy
- Case studies
- Ensuring you trust your data
- Breakout groups to discuss your experimental design
- Reference genomes

Human Genome Project



- Human genome project – started in 1990, completed in 2003
- Sequenced across ~20 institutions worldwide
- Cost an approximate \$5 billion US dollars

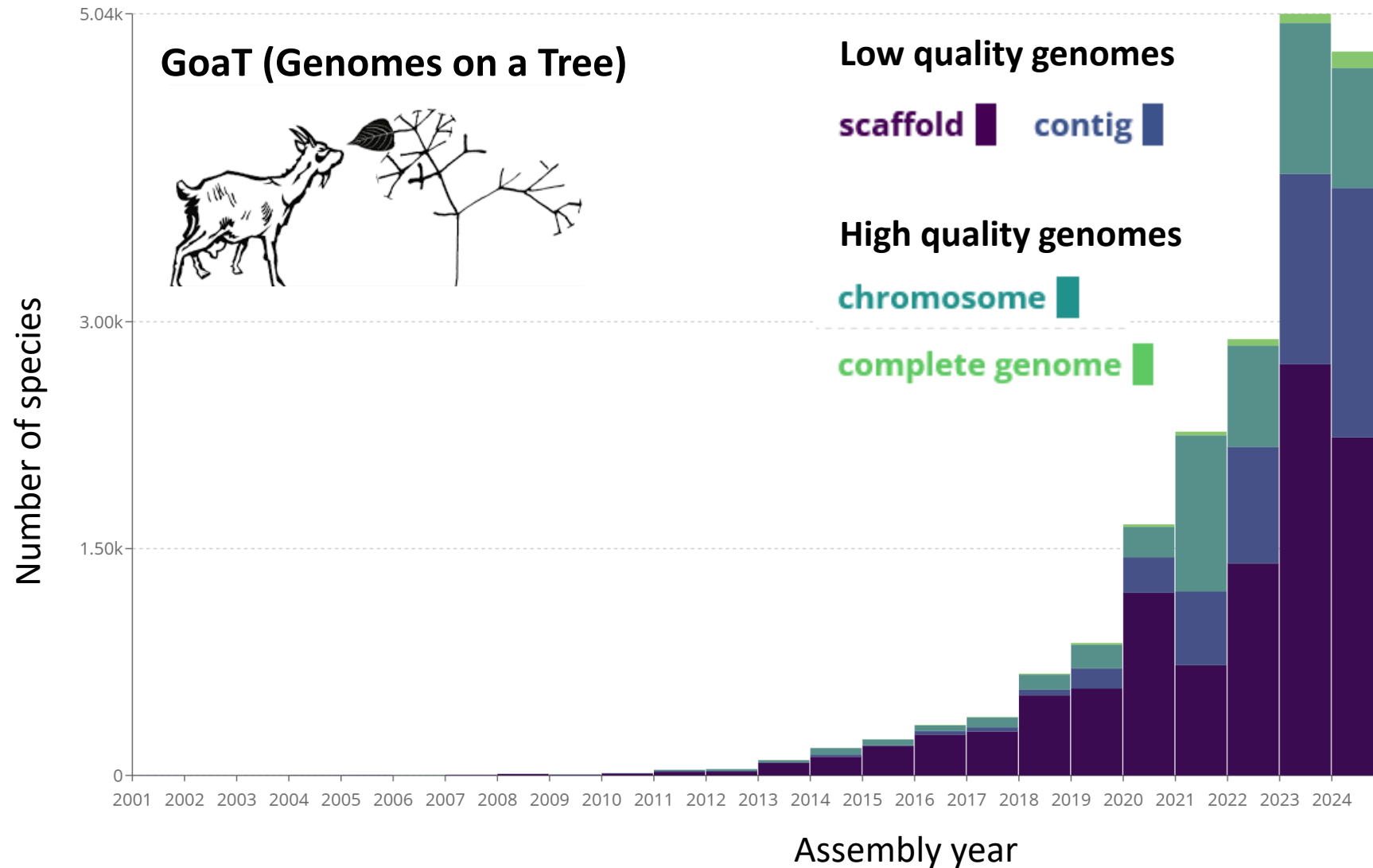
Sequencing costs are decreasing rapidly



Since Oct 2023
PacBio Revio
(66 Gbp per lane
in 15 kb reads)



More and more species have reference genomes



Two main types of high-throughput sequencing

- **Short-read sequencing**

- Reads are typically 150-300 bp long
- Cheaper than long-read sequencing
- E.g. Illumina, Ultima Genomics

- **Long-read sequencing**

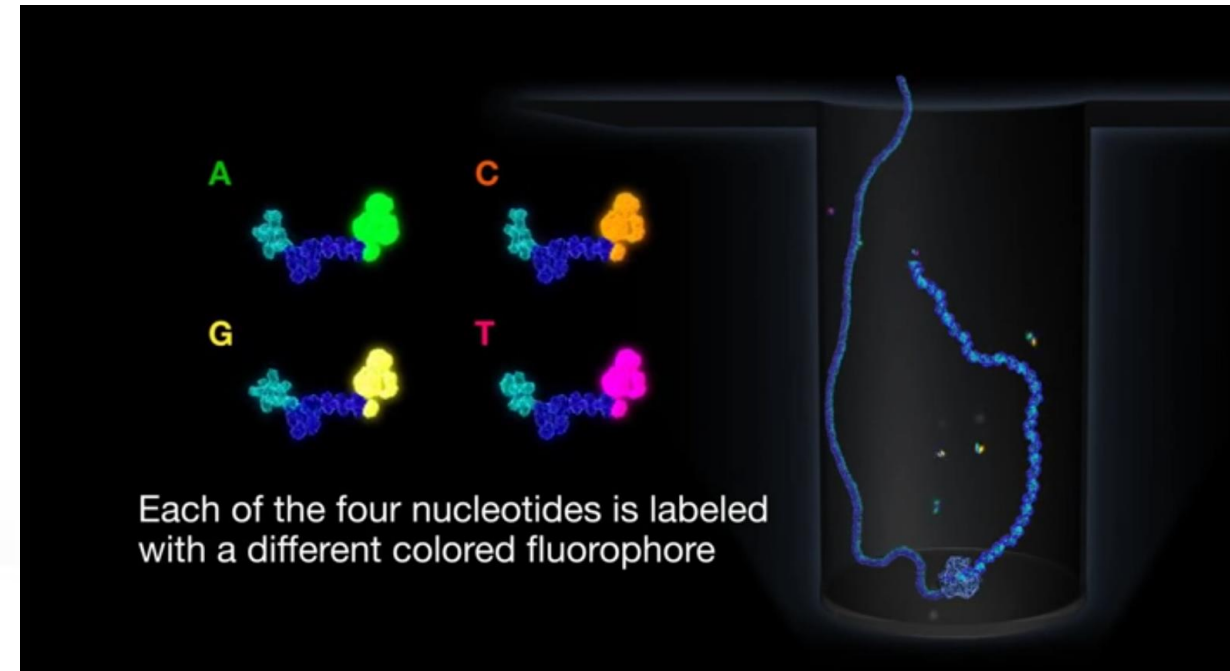
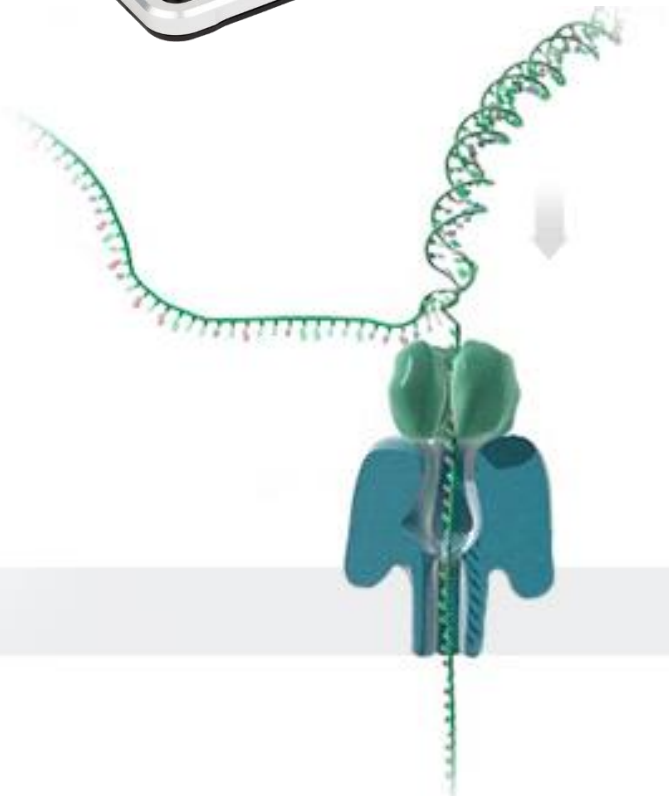
- Reads are typically >10 kb long (PacBio: 10-20 kbp, ONT/Nanopore: >100 kbp)
- More expensive than short-read technologies
- Required for making a reference genome
- E.g. PacBio or ONT/Nanopore

Long read sequencing technologies



ONT / Nanopore
(>100 kb reads)

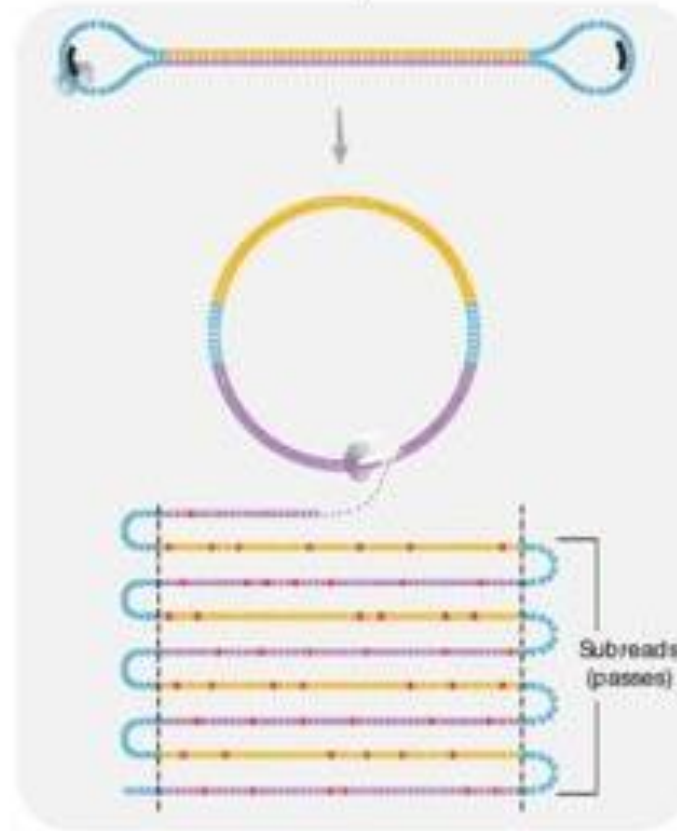
PacBio
~10-20 kb reads



PacBio HiFi reads (99.95% accuracy)

Each DNA-fragment is sequenced many times to get a high-quality consensus (=summary) read

Multi-pass sequencing on Sequel II System



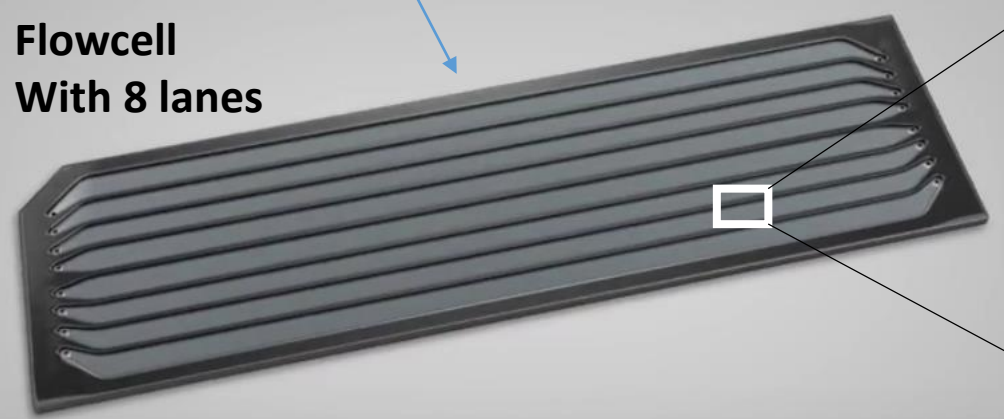
HiFi Read Base Calling



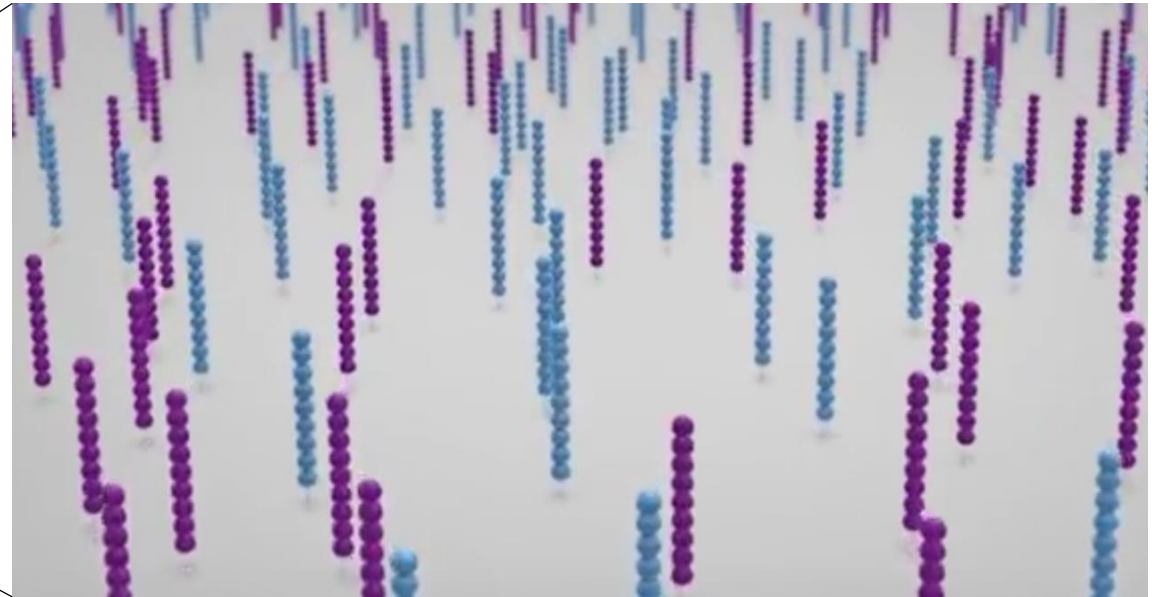
Illumina flowcell

DNA fragments with Illumina adapters

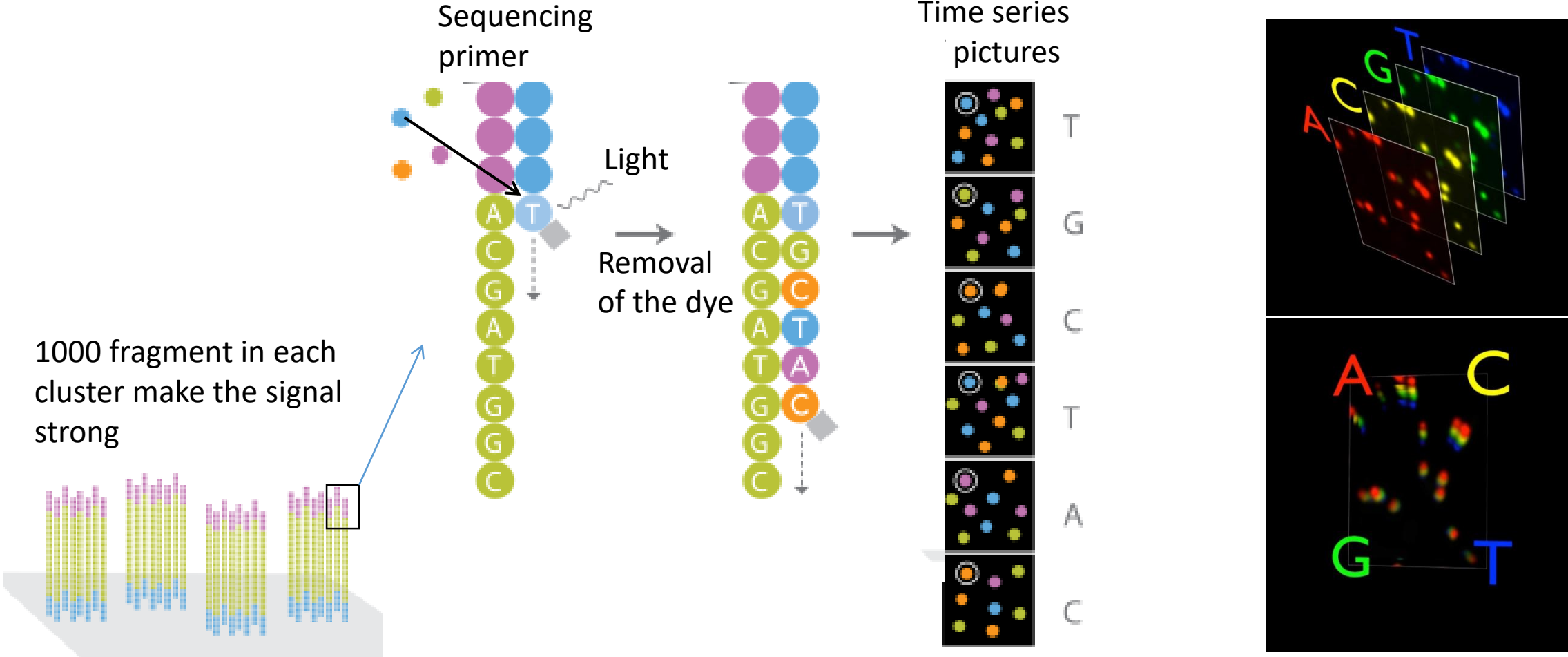
Flowcell
With 8 lanes



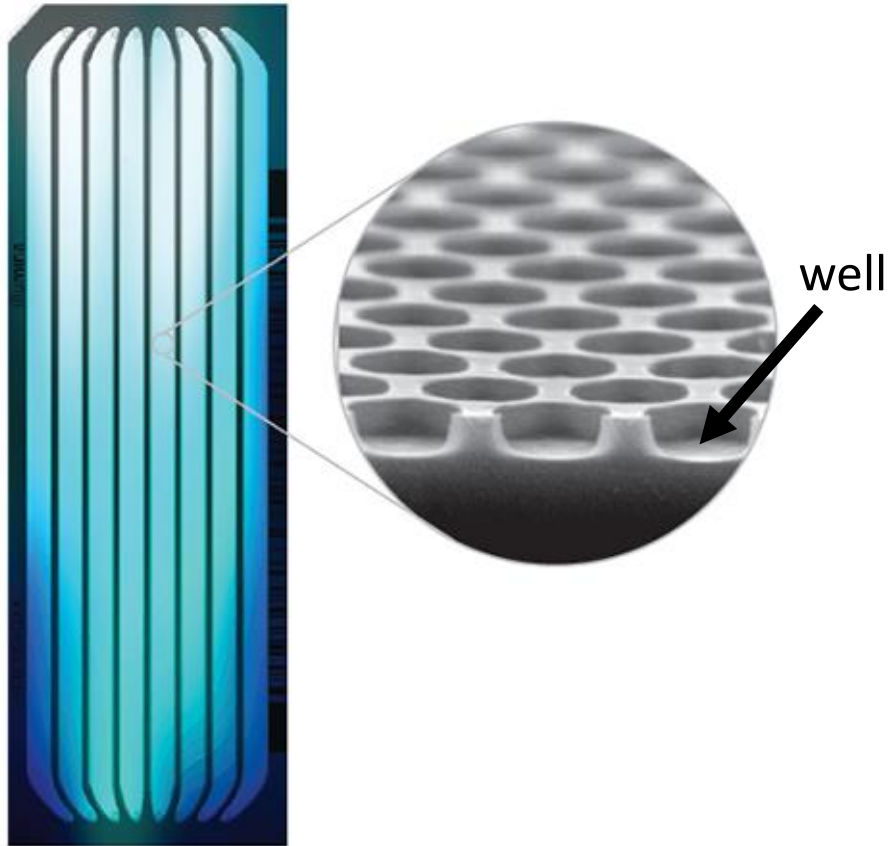
Each lane contains a dense lawn of Illumina primers










Short-read sequencing with Illumina



Newer Illumina machines use wells and only 2 colours
(e.g. Novaseq, Nextseq, MiniSeq. This makes it faster and cheaper)

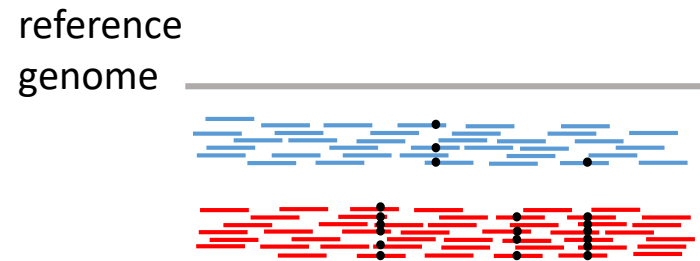


2-Channel Chemistry				
	 A	G	 T	 C
Image 1				
Image 2				
Result	A	G	T	C

Sequencing approaches for speciation genomics

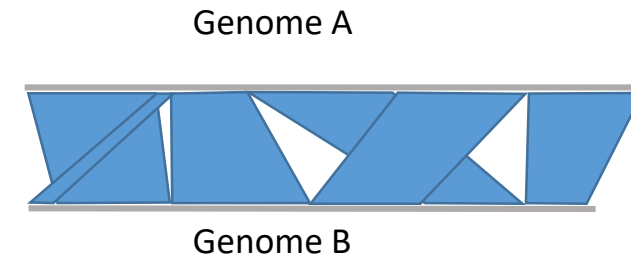
Whole-genome resequencing (short-read data)

- Requires a reference genome
- individuals need to be from the same or closely related species
- Complete genome sequenced



Genome assembly comparisons (long-read data)

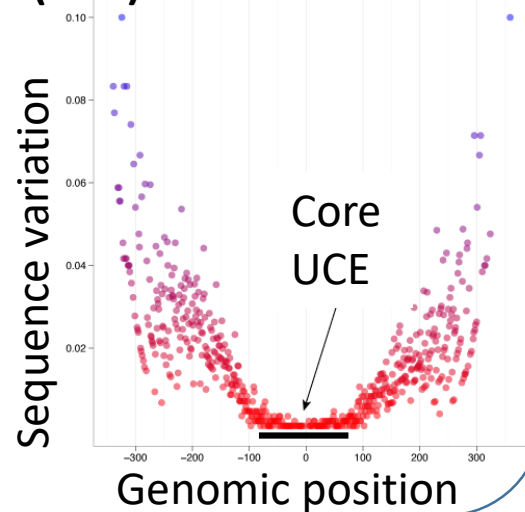
- Comparative genomics – studying structural variation between species, can be distantly related
- Gene expansions, transposable elements etc
- Phylogenomics across deeply divergent species
- Pangenomics – multi-genome assemblies to study within-species variation in structural variants



Reduced-representation techniques (only parts of the genome sequenced)

Ultra-conserved elements (UCE)

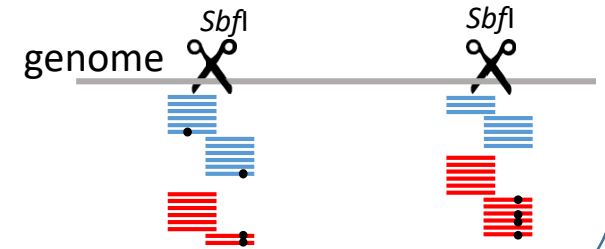
- Sequence capture with baits based on genomic regions that are conserved across many species
- Works with highly divergent species



Restriction Associated Sequencing (RAD)

(similar methods: GBS, ddRAD)

- does not require primers/baits or reference genome
- individuals need to be from the same or closely related species
- Information from thousands of loci distributed across the genome



Targeted or amplicon sequencing, e.g. barcoding

- Sequencing one or few genes
- requires primers
- e.g. CO1 (mitochondrial barcoding region), advantage: large database (BOLD) available to compare to for species identification

Environmental DNA (eDNA)

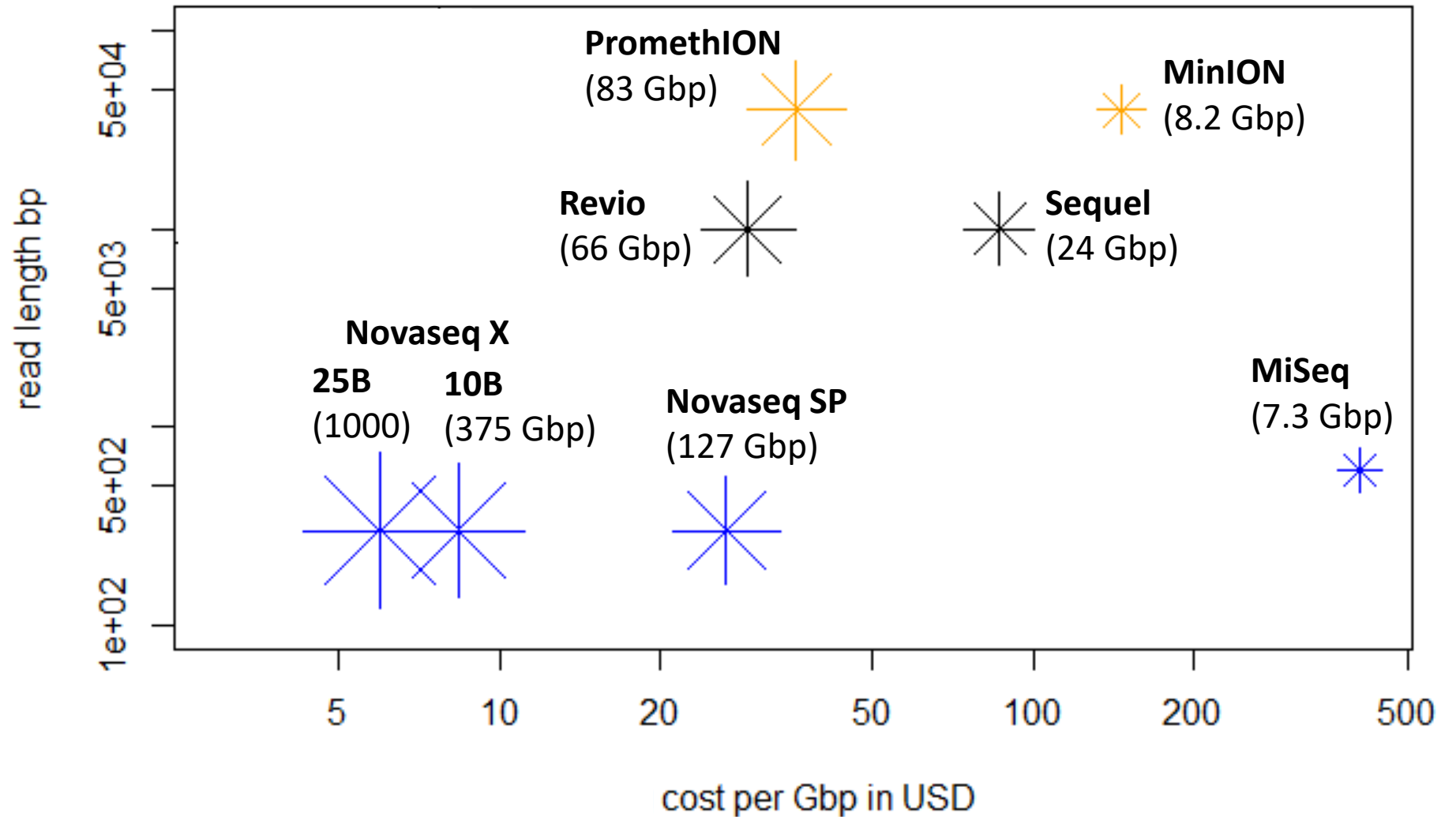
- Mostly CO1 sequencing from soil, water, air (spider webs)
- Studying local species richness



Read length versus per Gbp sequencing costs for different sequencing machines (note the axes are in logarithmic scale)

- * ONT
- * PacBio
- * Illumina

Star size shows the total throughput per lane, also given in parentheses ()



Sequencing strategy

Total sequencing data gets divided by:

- Number of sites to sequence
 - Depends on genome size and sequencing strategy, e.g. RAD versus whole-genome
- Number of specimens to sequence
 - and how many specimens per population, how many populations per species
- Sequencing depth (e.g. sequencing at 10x depth of coverage)

- Example: 1 NovaseqX 10B lane
~2.5 billion paired-end reads of 150 bp each -> 375 Gbp data
 - 100 whole-genomes of a species with 0.375 Gbp genome size at 10x coverage
 - 19 whole-genomes of a species with 1 Gbp genome size at 20x coverage
 - 375 individuals sequenced with a RAD sequencing approach resulting in 50 Mbp at a sequencing depth of 20x

Different questions require different experimental design

Resolving the taxonomy

- Placing potentially new species
- Species delineation



Adaptation to climate change

- Identifying genomic regions involved in adaptation

Pfenninger
et al. 2021



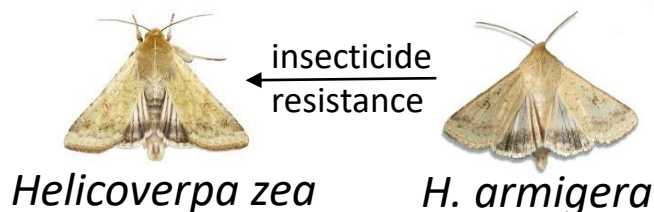
Are the species declining?

- Detecting past inbreeding
- Assessing genetic diversity



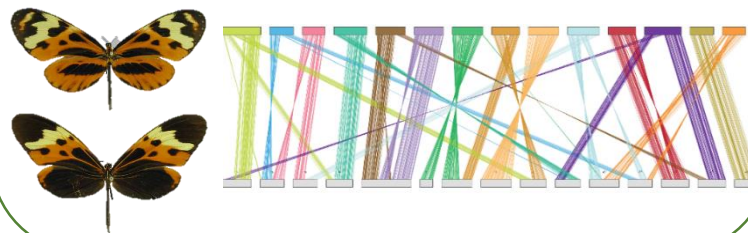
Detecting introgression

- Are populations/species hybridising now or in the past?
- Finding regions of adaptive introgression



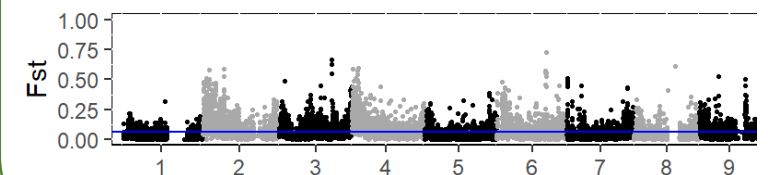
Studying genome evolution

- Gene expansions, e.g. olfactory
- Chromosomal rearrangements
- Genome size evolution (TEs, etc)



What genes are under divergent selection

- Short-term (recent sweep)
- Long-term (fixed substitutions)



Sampling considerations

- Knowing the natural history and morphologically inferred taxonomy if your helps massively
- If you want to test for introgression, you need outgroups
- Aim to get even sampling of populations/species to compare
- Sampling sites matter:
e.g. if you want to test for ongoing gene flow, you need specimens from sympatry or parapatry
- High number of specimens needed for:
 - Inferring recent demographic changes
 - Studying rare alleles
 - GWAS (inferring the genetic basis of traits, particularly if polygenic)

Case studies to discuss

Three examples with very different optimal sampling designs

- **Identifying barriers to gene flow and the genetic basis of relevant traits in a hybrid zone**

H. erato and *H. melpomene* in Ecuador



Heliconius erato



Highland



Lowland

Heliconius melpomene



Highland



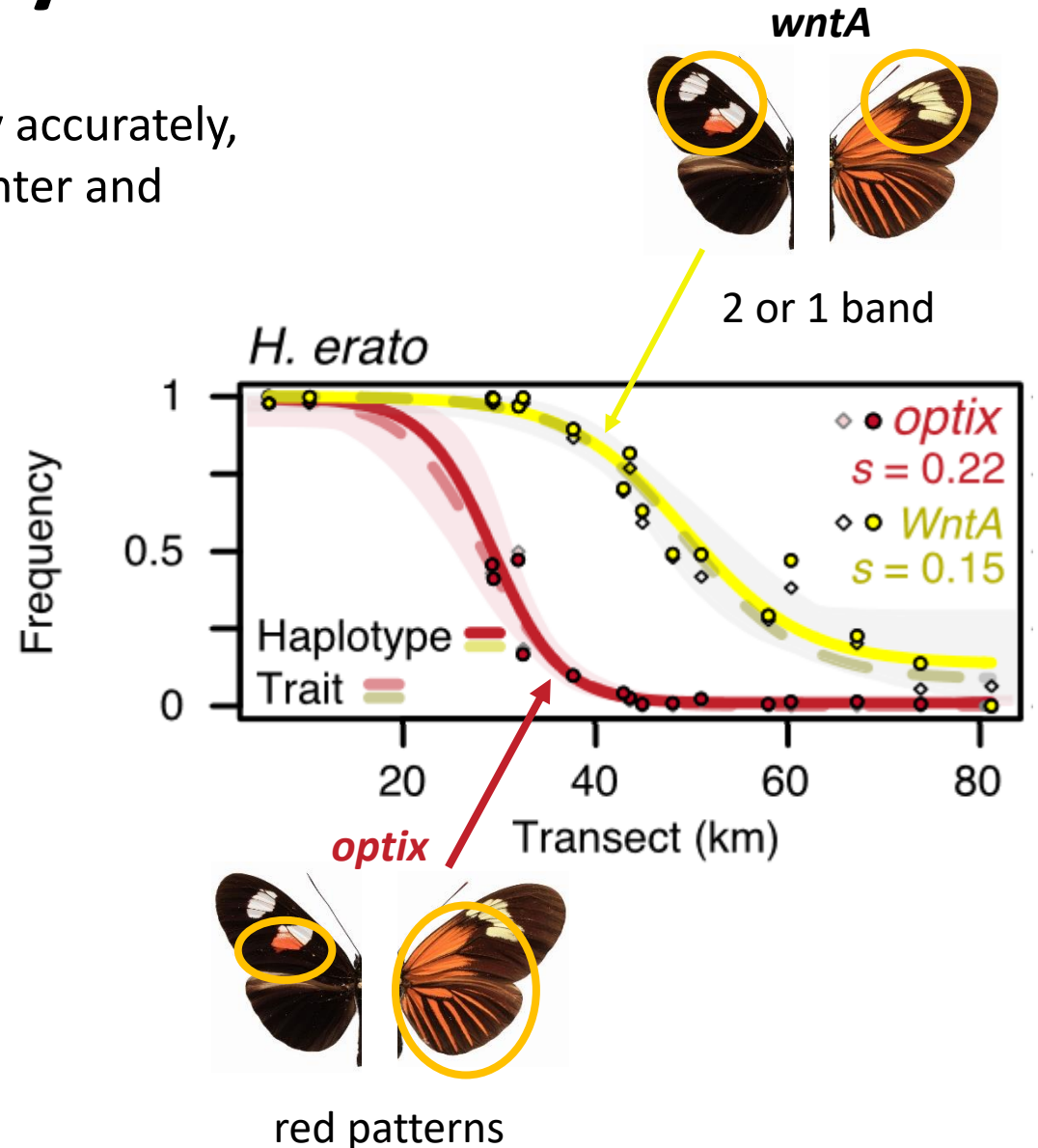
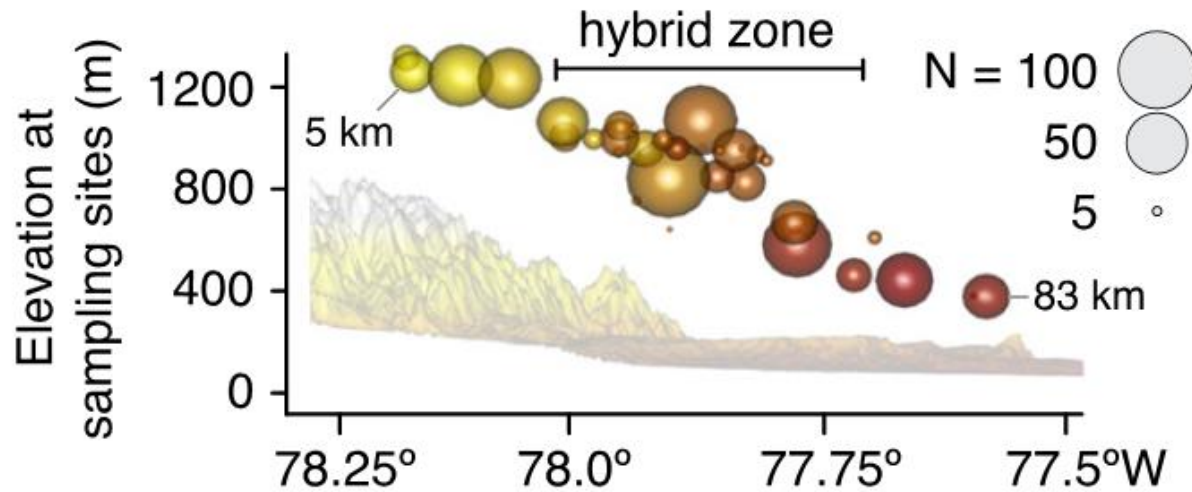
Lowland

Patricio Salazar Chris Jiggins

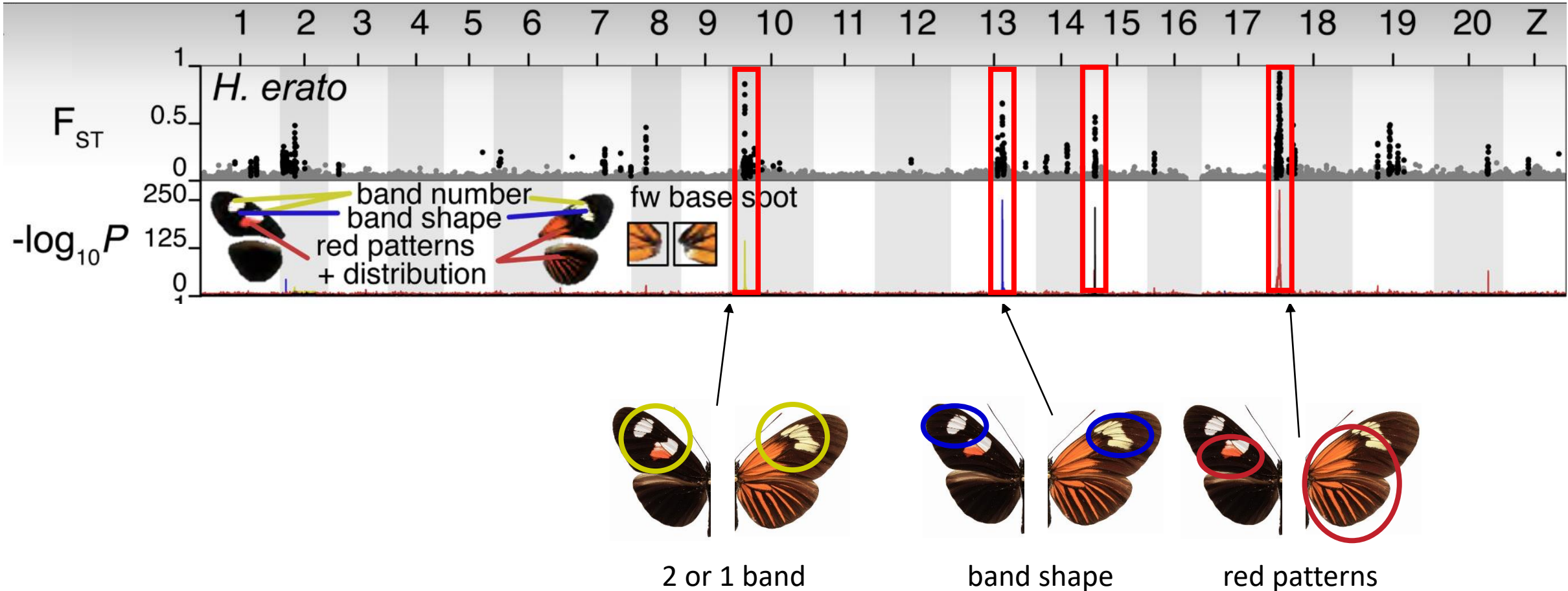


Heliconius erato hybrid zone

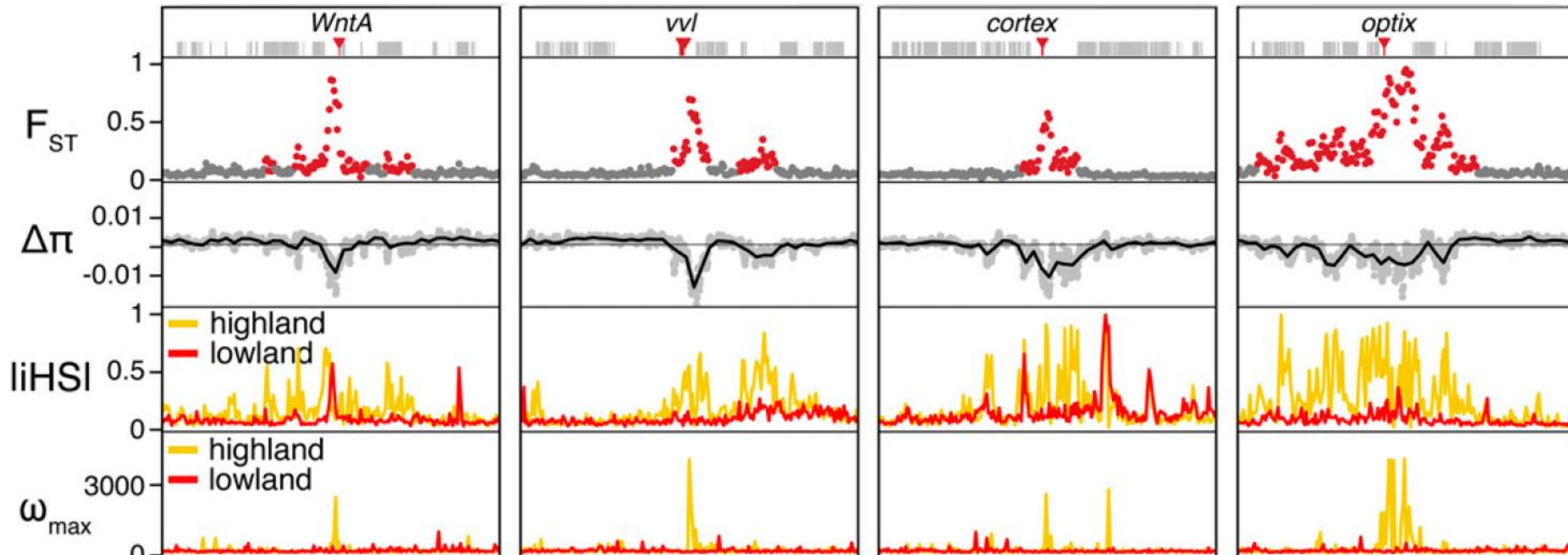
Many individuals per population to estimate allele frequency accurately, many populations along the hybrid zone to infer the cline center and shape accurately



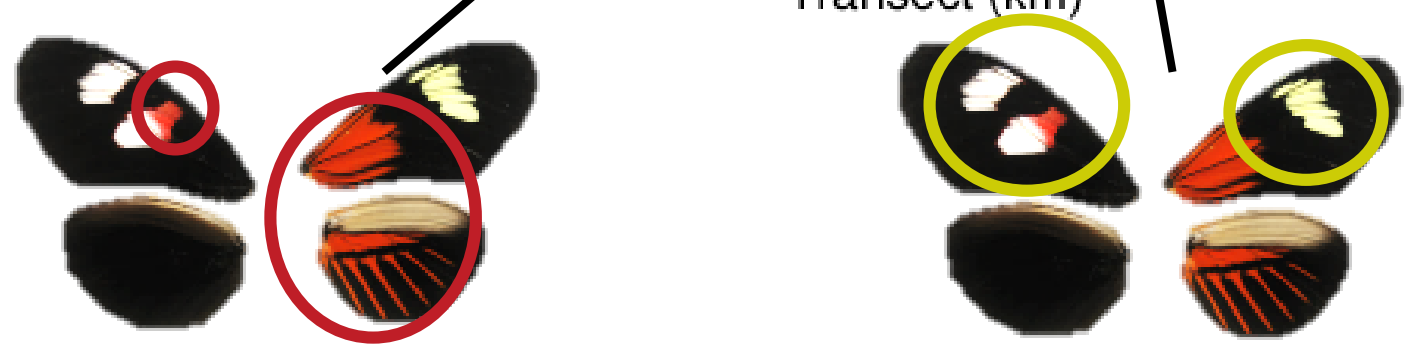
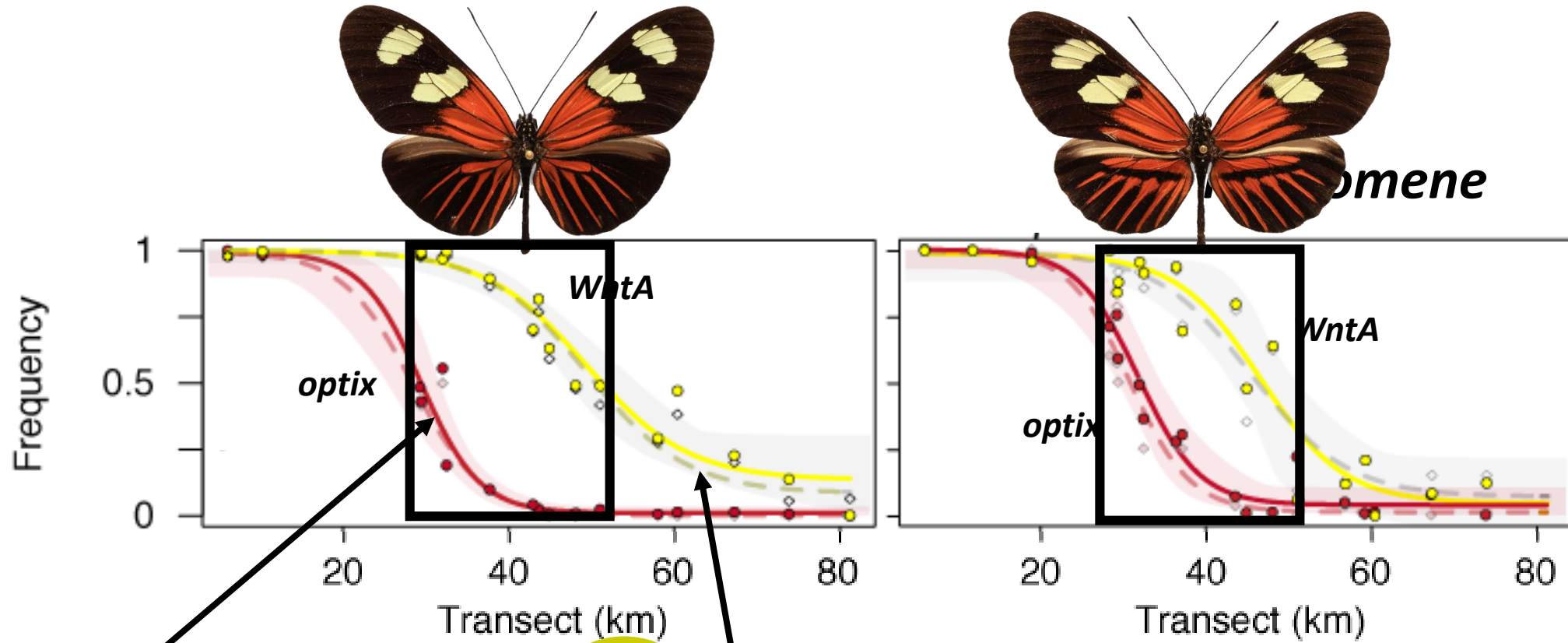
F_{ST} peaks coincide with GWAS peaks and show signatures of selective sweeps



F_{ST} peaks coincide with GWAS peaks and show signatures of selective sweeps



In both *Heliconius erato* and *H. melpomene*, the clines at the major colour loci are shifted by 18 km

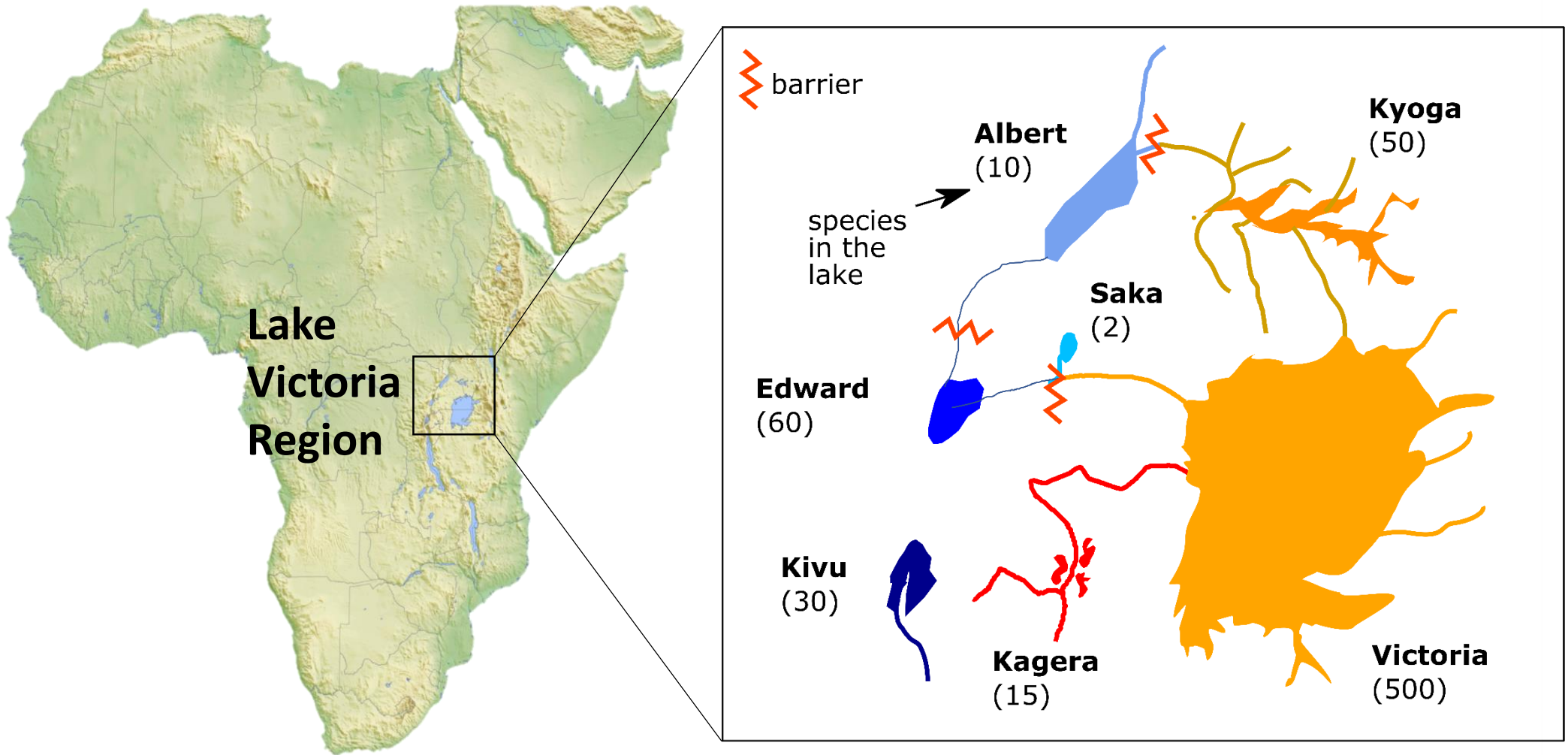


Case studies to discuss

Three examples with very different optimal sampling designs

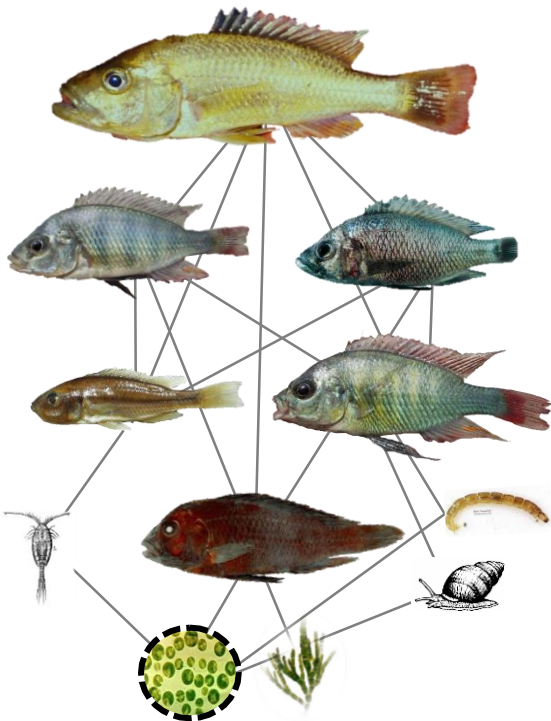
- Identifying barriers to gene flow and the genetic basis of relevant traits in a hybrid zone
- **Inferring if a fish species community evolved in a lake or if the different species independently colonised the lake.**

The same lineage of cichlid fishes diversified in all major lakes in the Lake Victoria Region in only 150,000 years

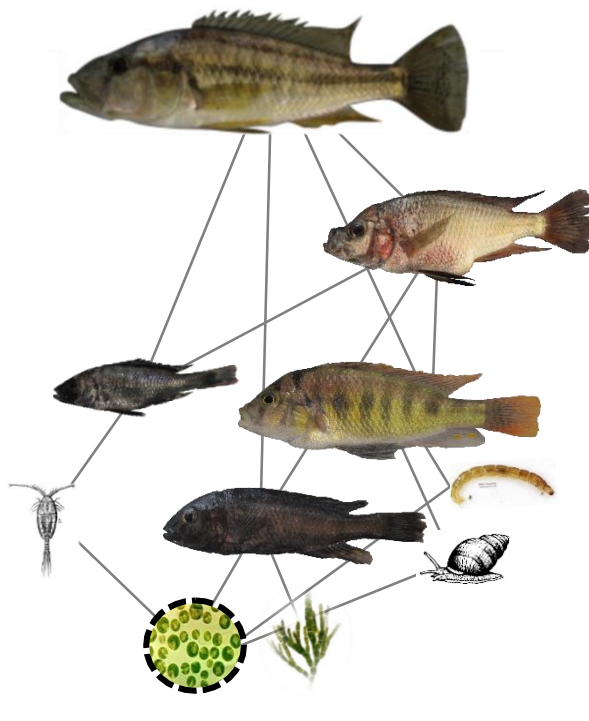


**In each lake, cichlids occupy a wide range of ecological niches.
But Lake Victoria is much younger! Only 16,000 years!**

Lake Edward

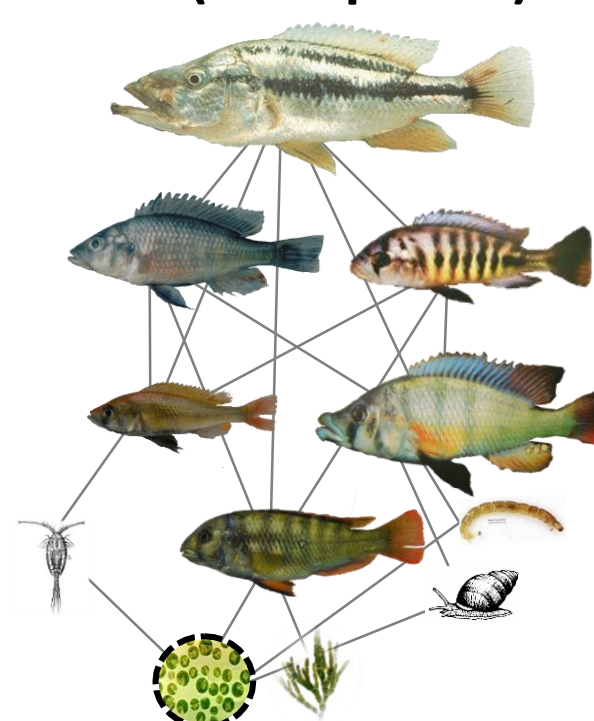


Lake Kivu

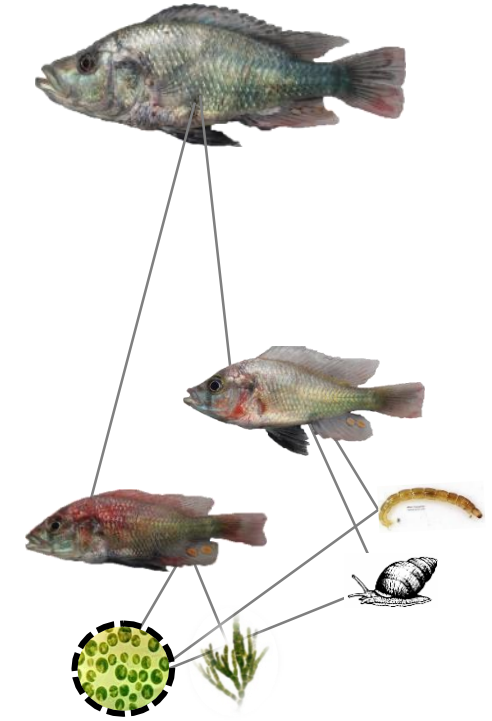


Less than 16,000 years old!

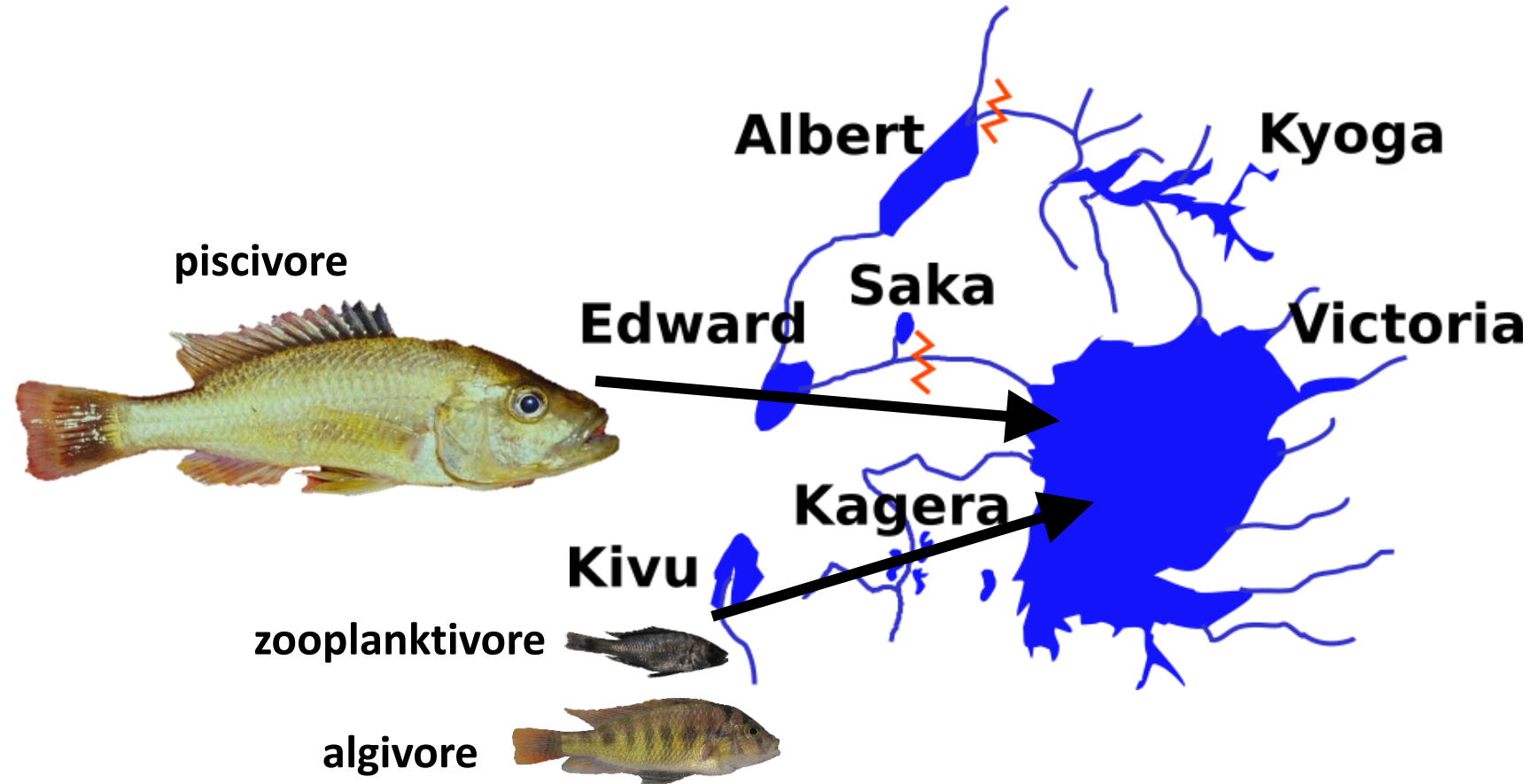
**Lake Victoria
(500 species)**



Kagera Lakes

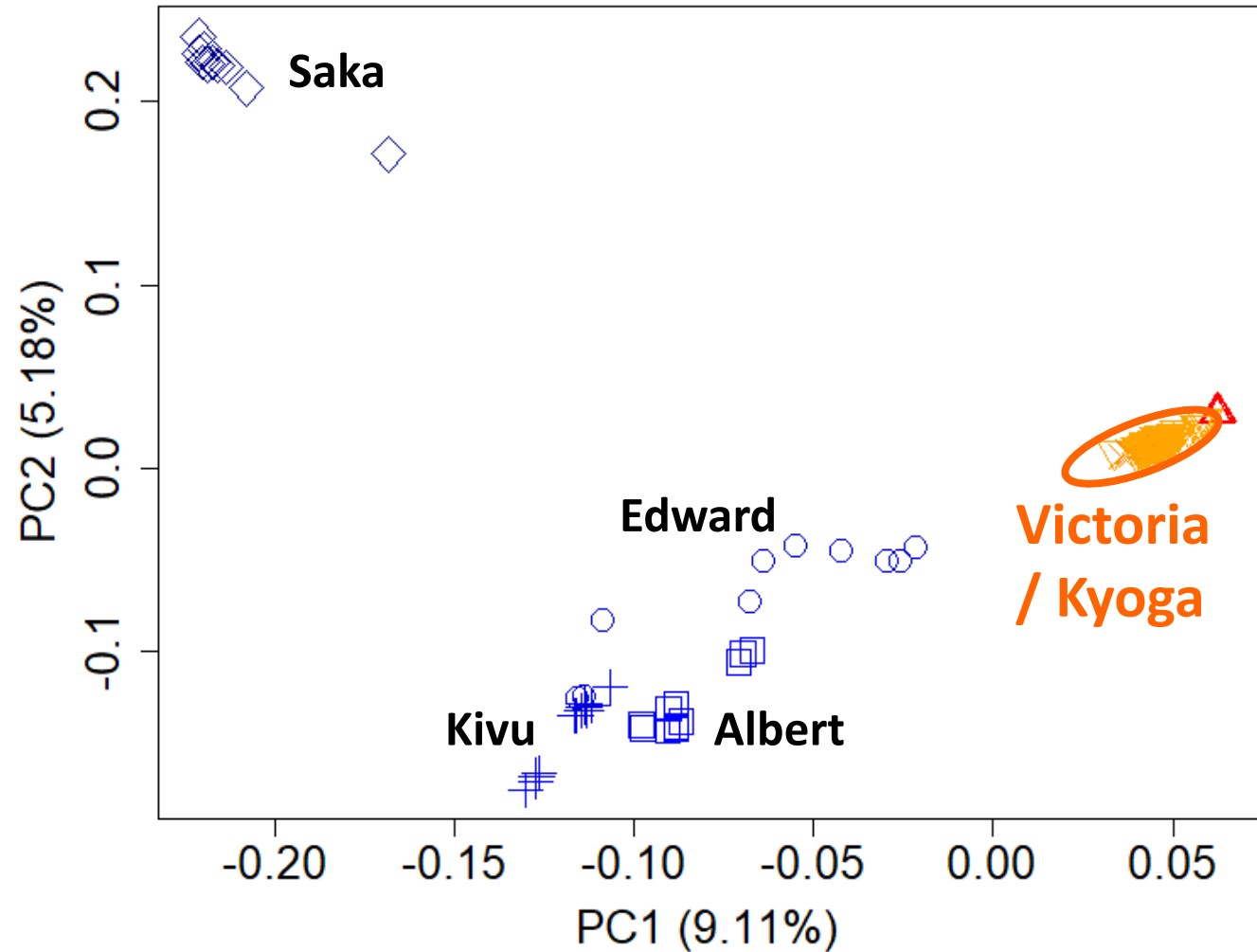
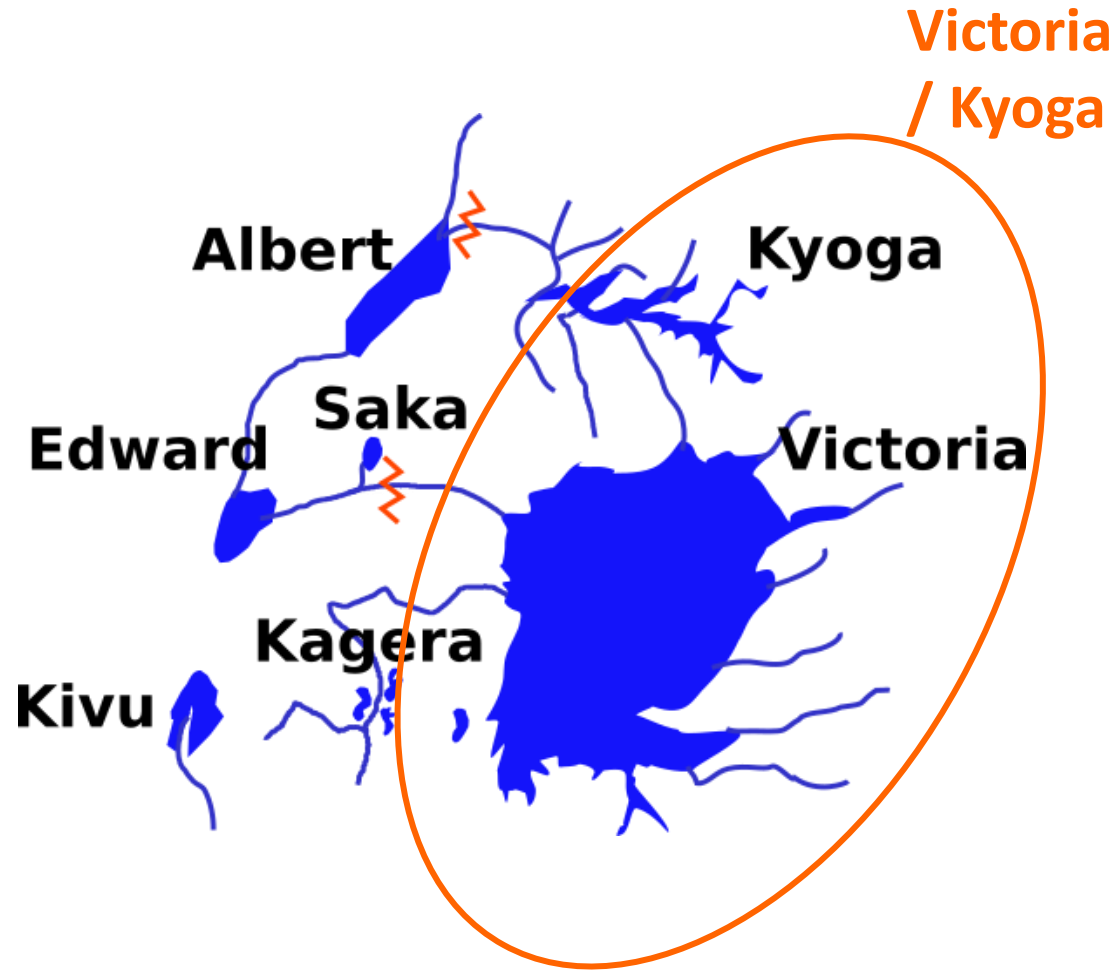


Did the different ecological groups colonise Lake Victoria independently?



PCA on whole-genomes separates Lake Victoria/Kyoga cichlids from others

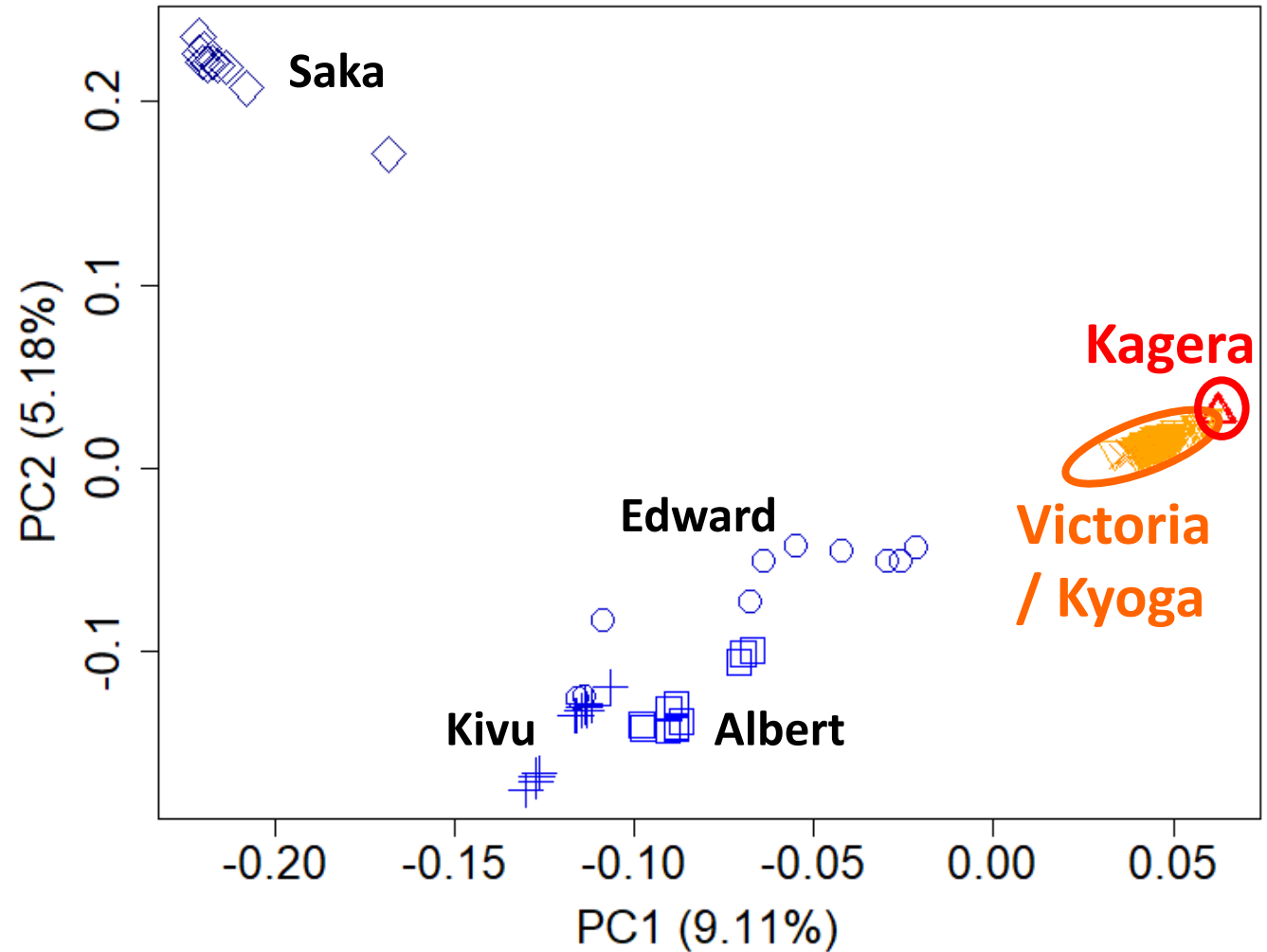
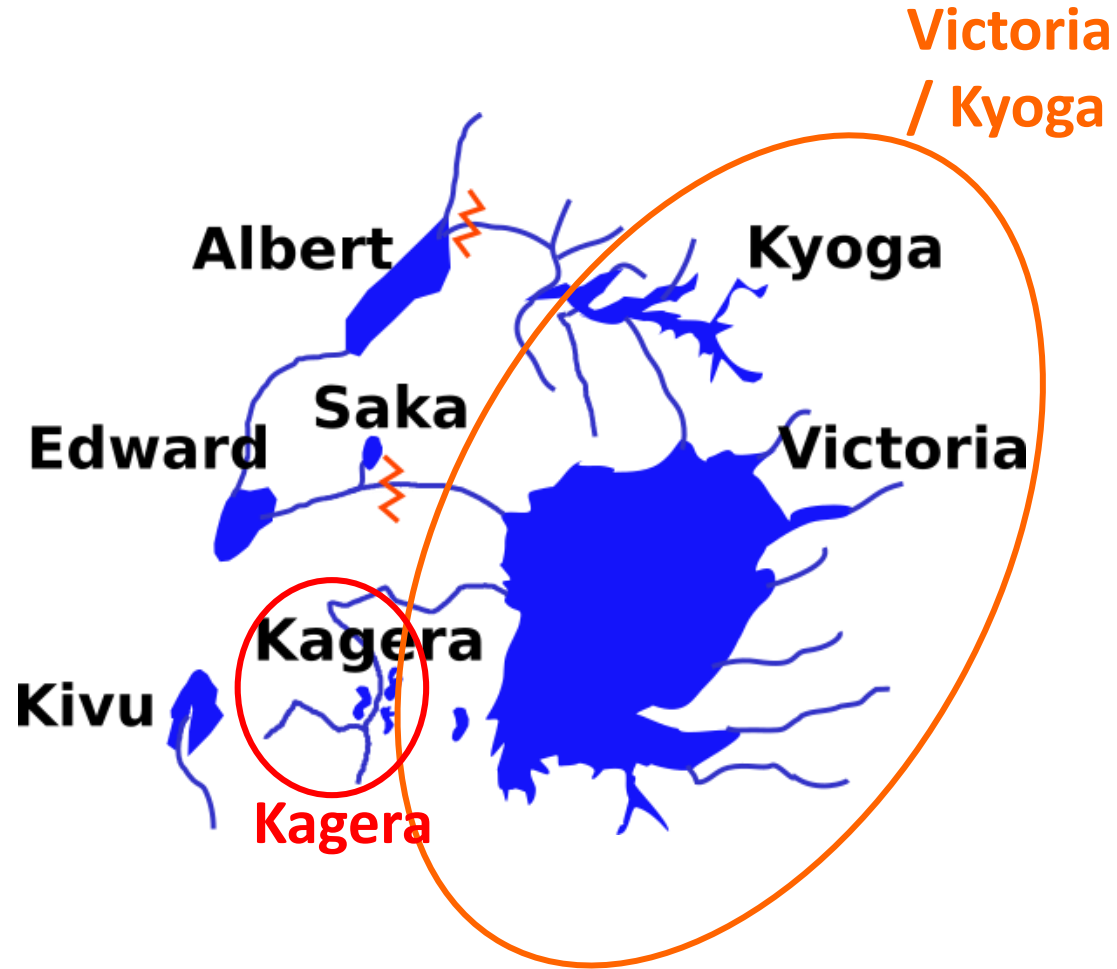
(152 genomes, 1.6M LD-pruned SNPs)



➤ migration barrier

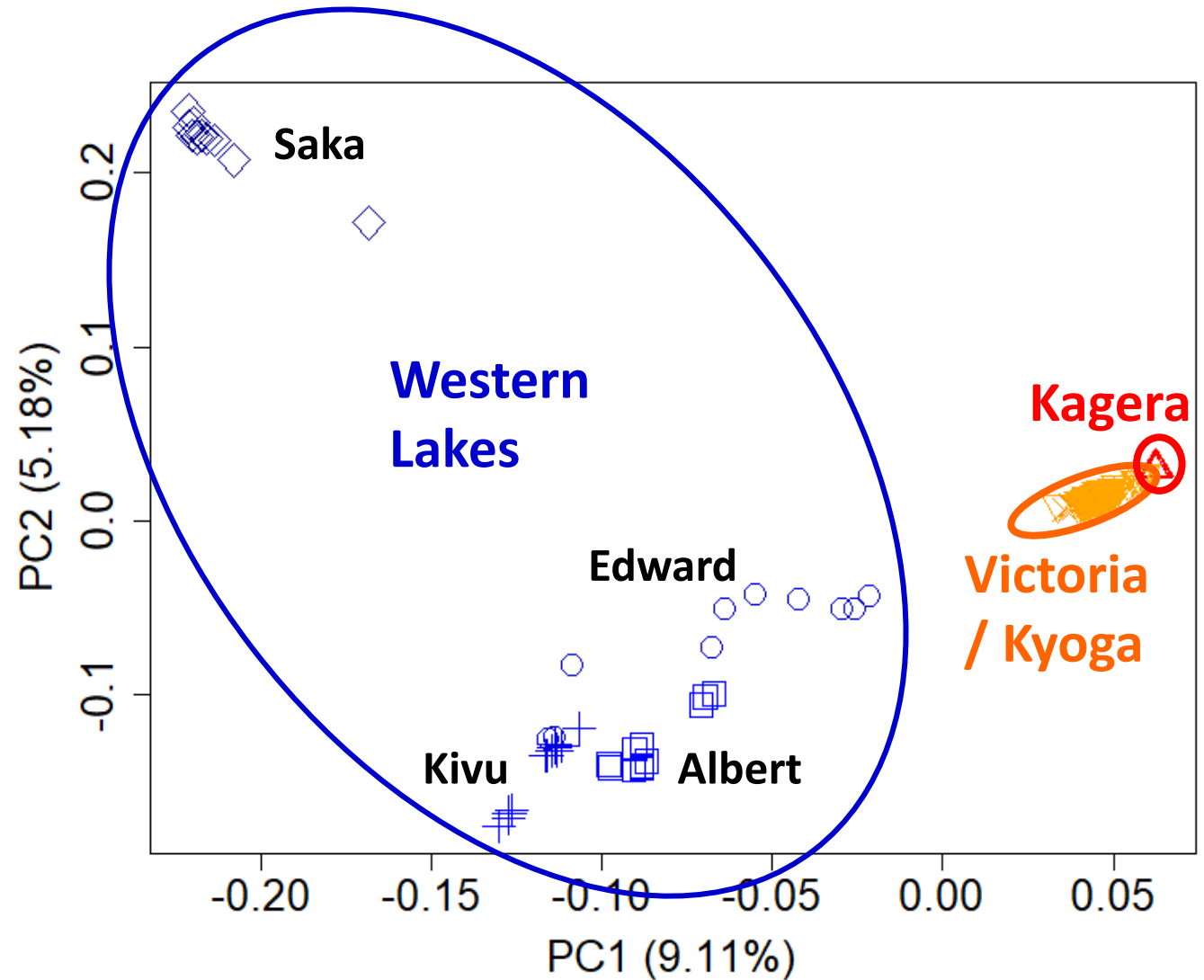
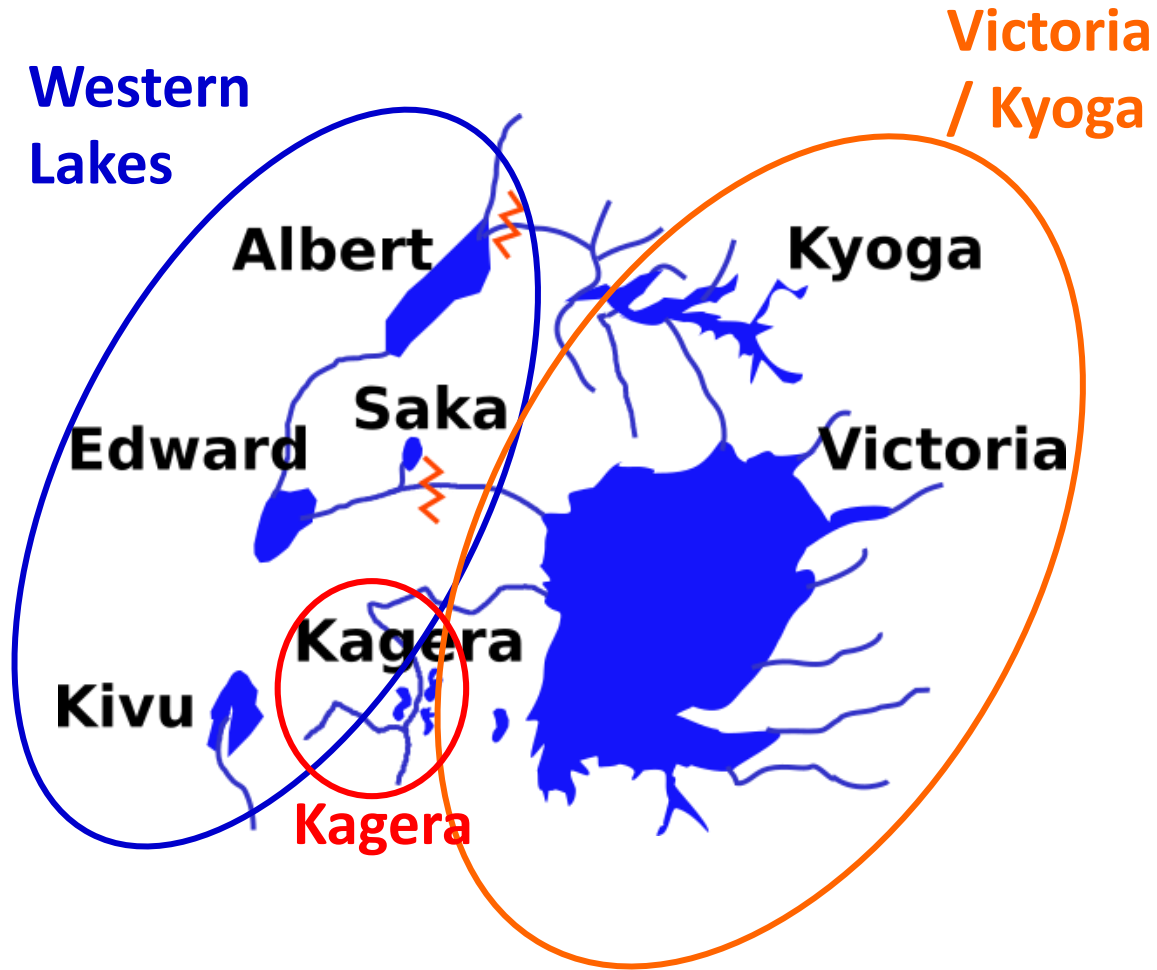
PCA on whole-genomes separates Lake Victoria/Kyoga cichlids from others

(152 genomes, 1.6M LD-pruned SNPs)



PCA on whole-genomes separates Lake Victoria/Kyoga cichlids from others

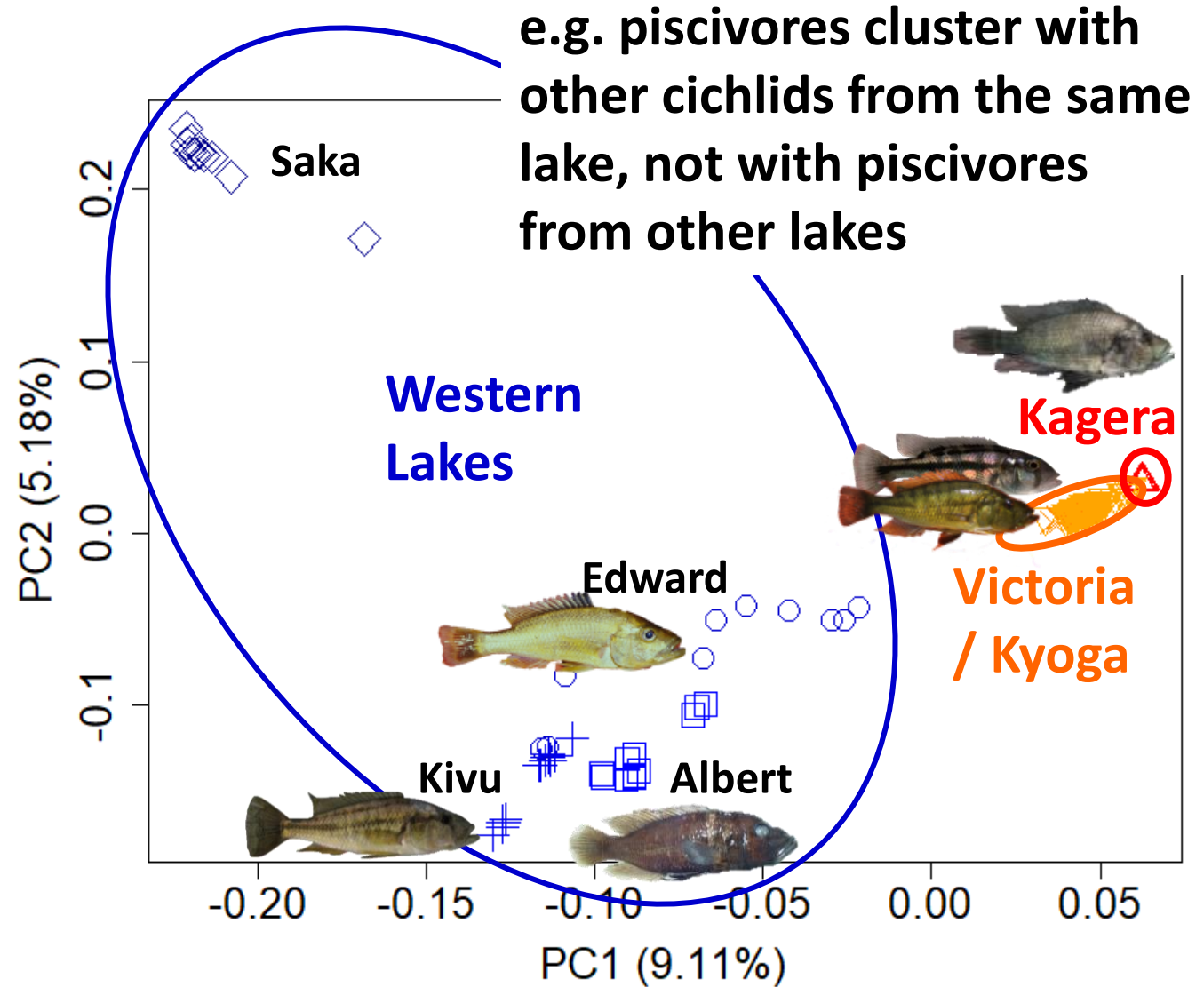
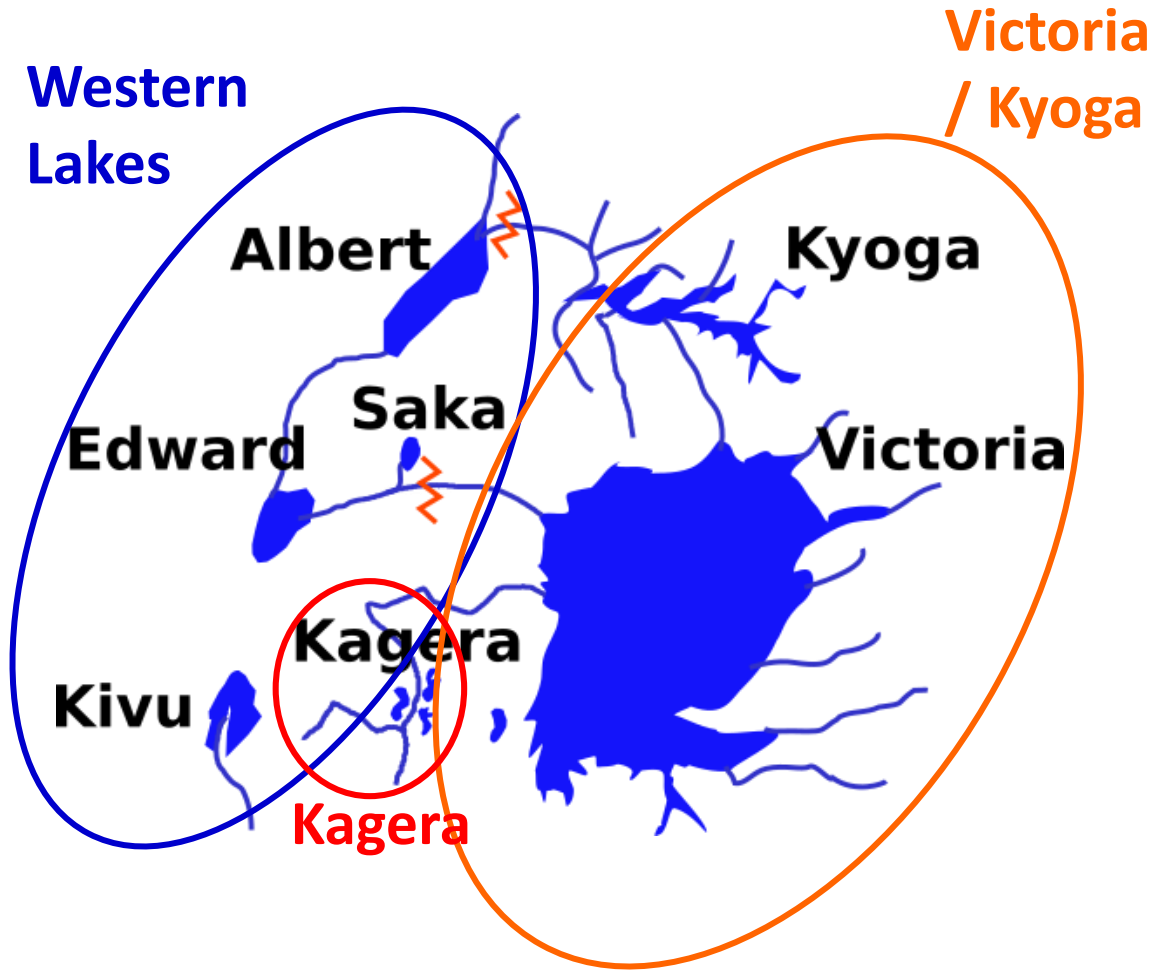
(152 genomes, 1.6M LD-pruned SNPs)



➤ migration barrier

PCA on whole-genomes separates Lake Victoria/Kyoga cichlids from others

(152 genomes, 1.6M LD-pruned SNPs)

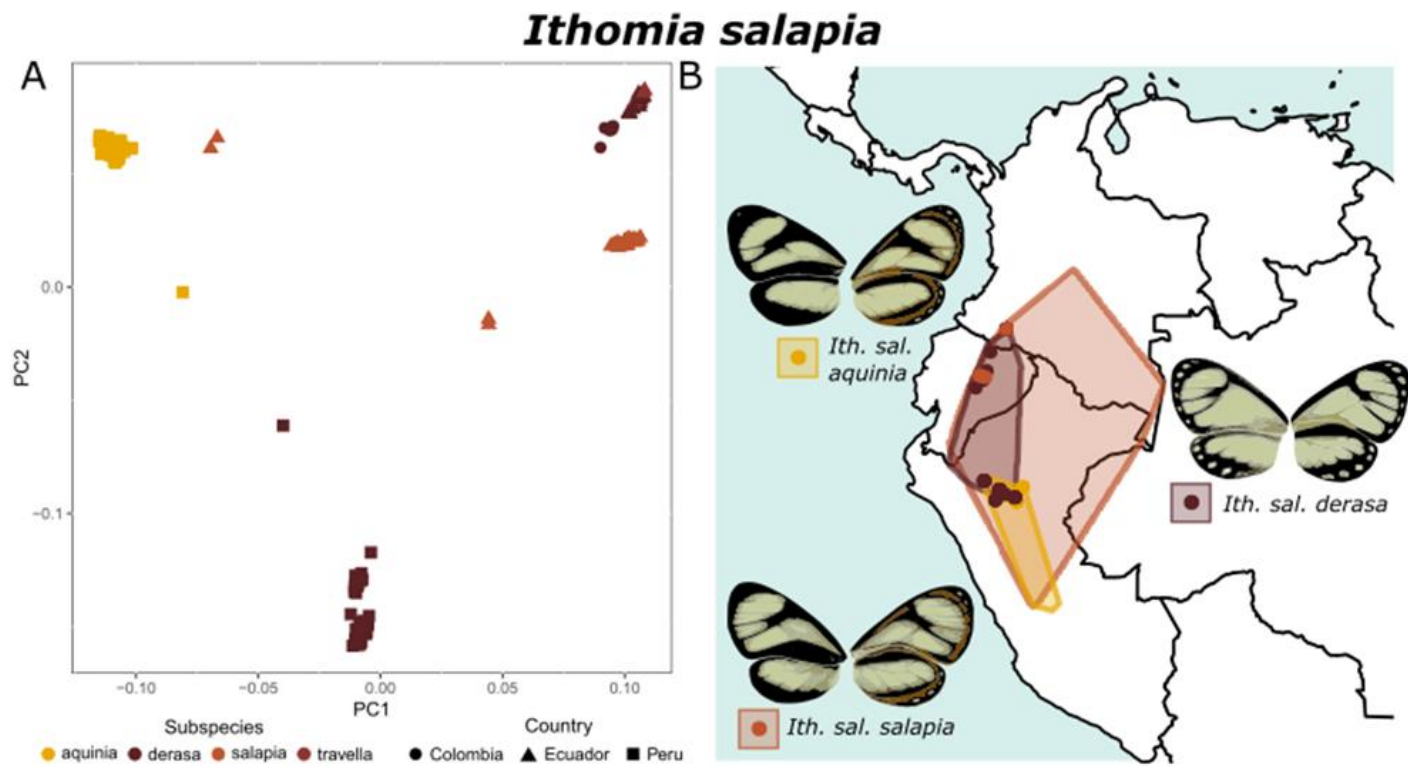
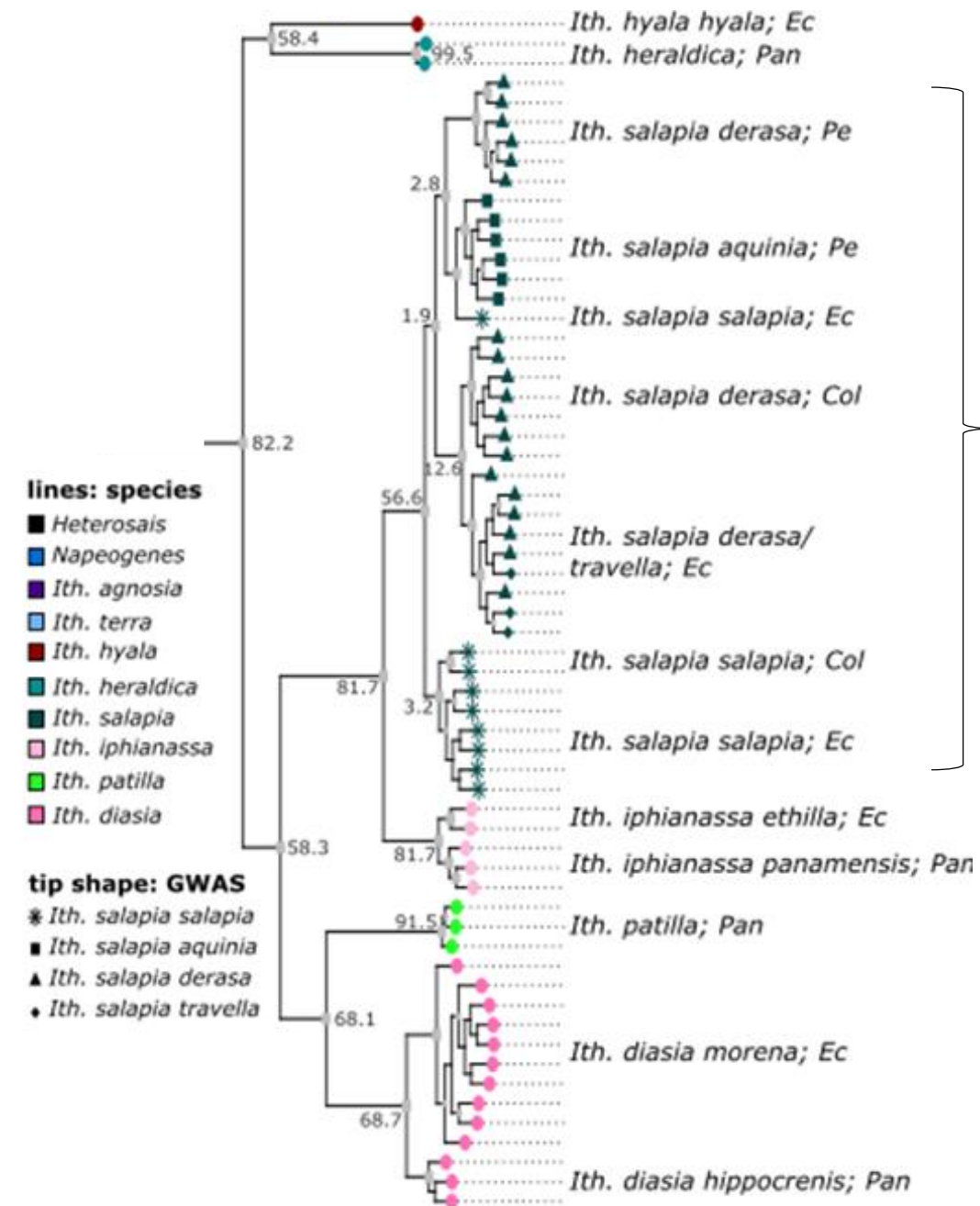


➤ migration barrier

Case studies to discuss

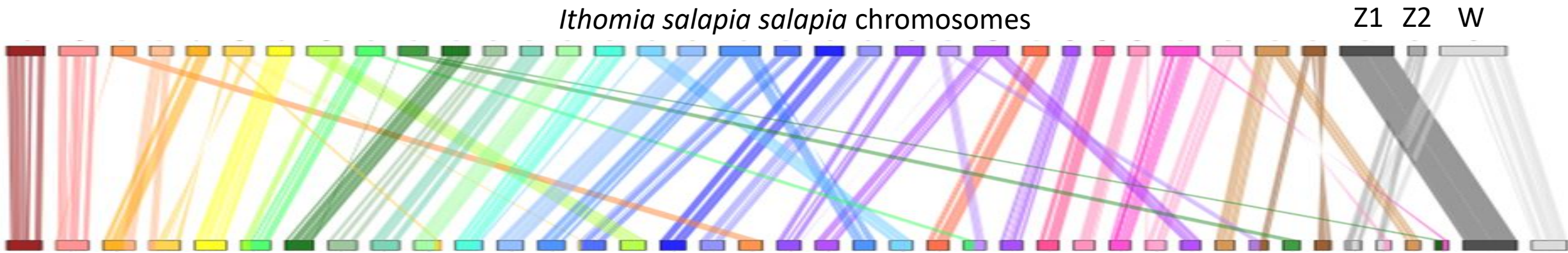
Three examples with very different optimal sampling designs

- Identifying barriers to gene flow and the genetic basis of relevant traits in a hybrid zone
- Inferring if a fish species community evolved in a lake or if the different species independently colonised the lake
- **Inferring if chromosomal rearrangements contribute to speciation in a butterfly species complex**

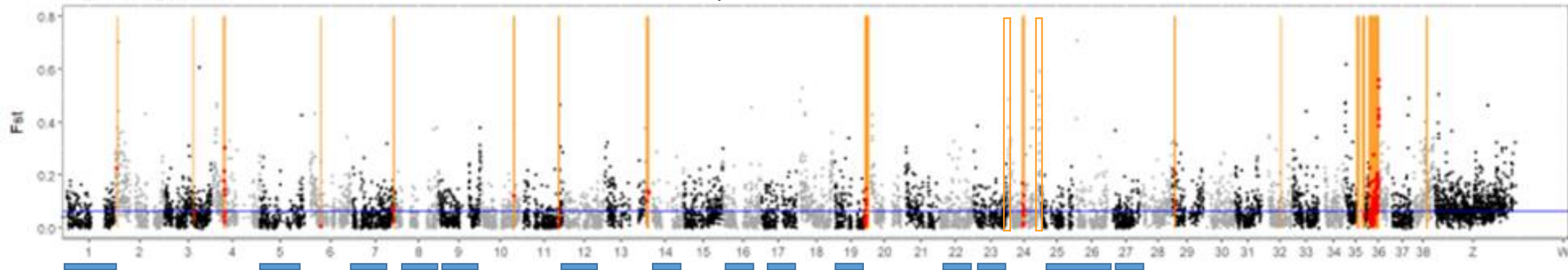


Ithomia salapia «subspecies» show massive chromosomal rearrangements, likely contributing to speciation

Ithomia salapia salapia chromosomes



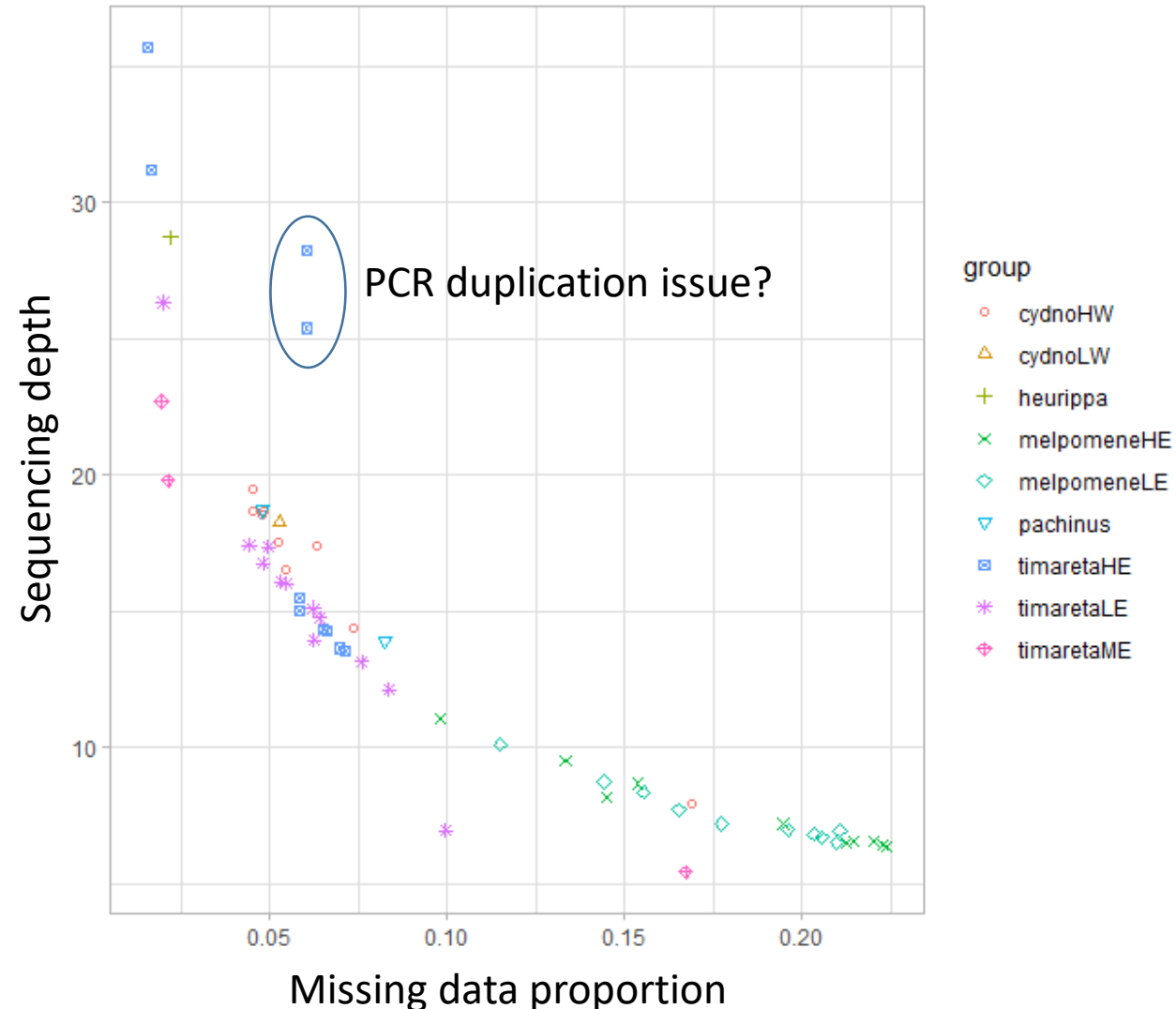
Ithomia salapia derasa chromosomes



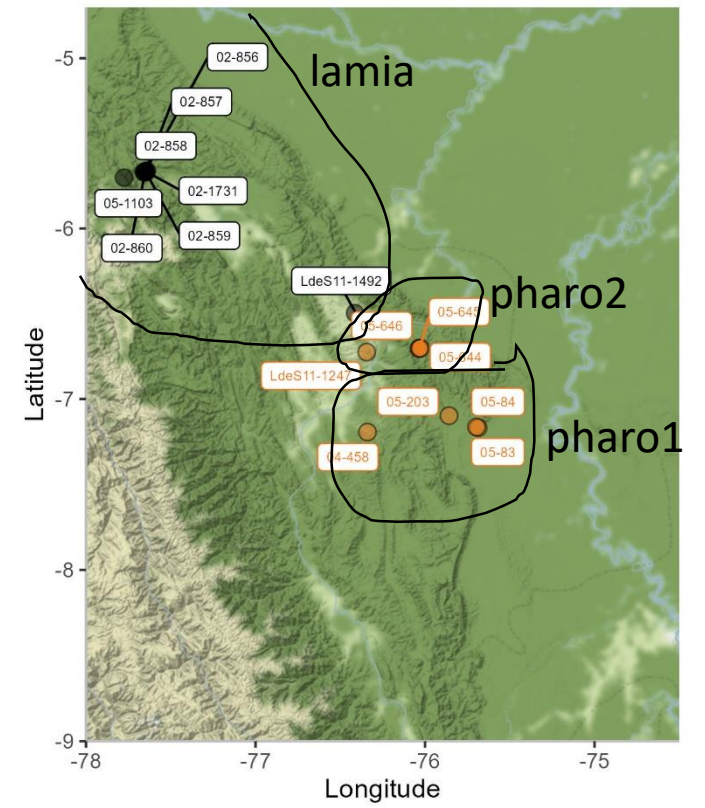
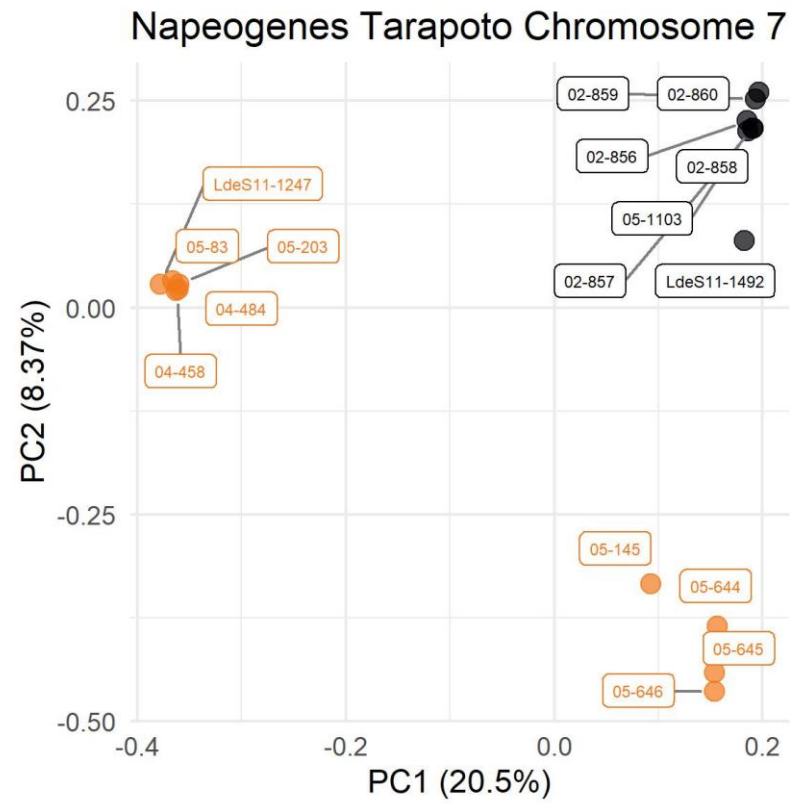
Once you get your data, make sure you trust it

Get a feeling for your data, do quality checks at all stages and plot your data

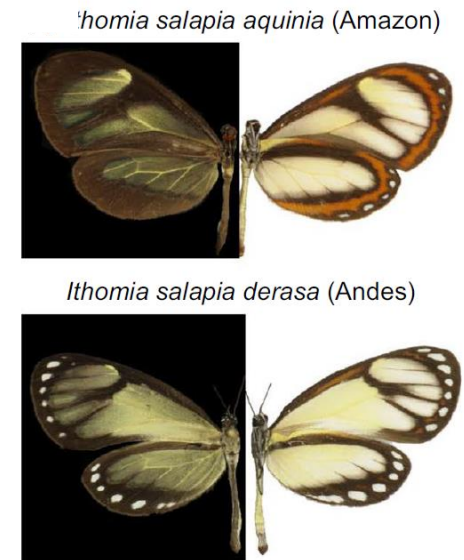
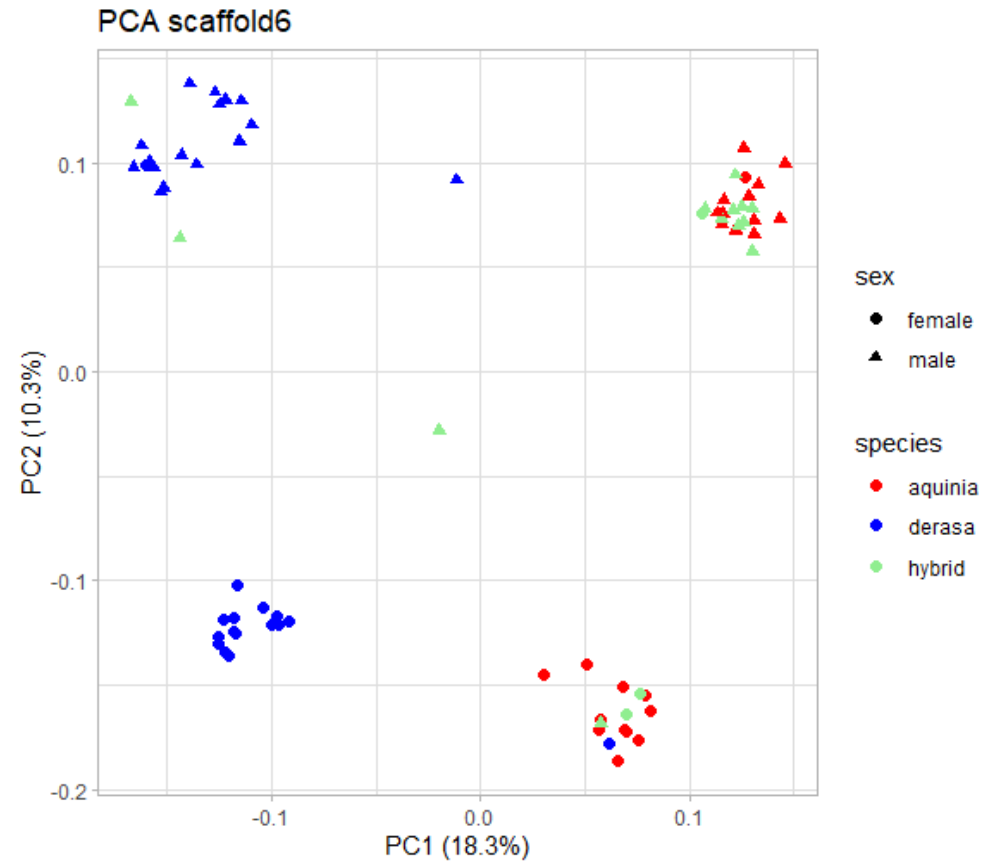
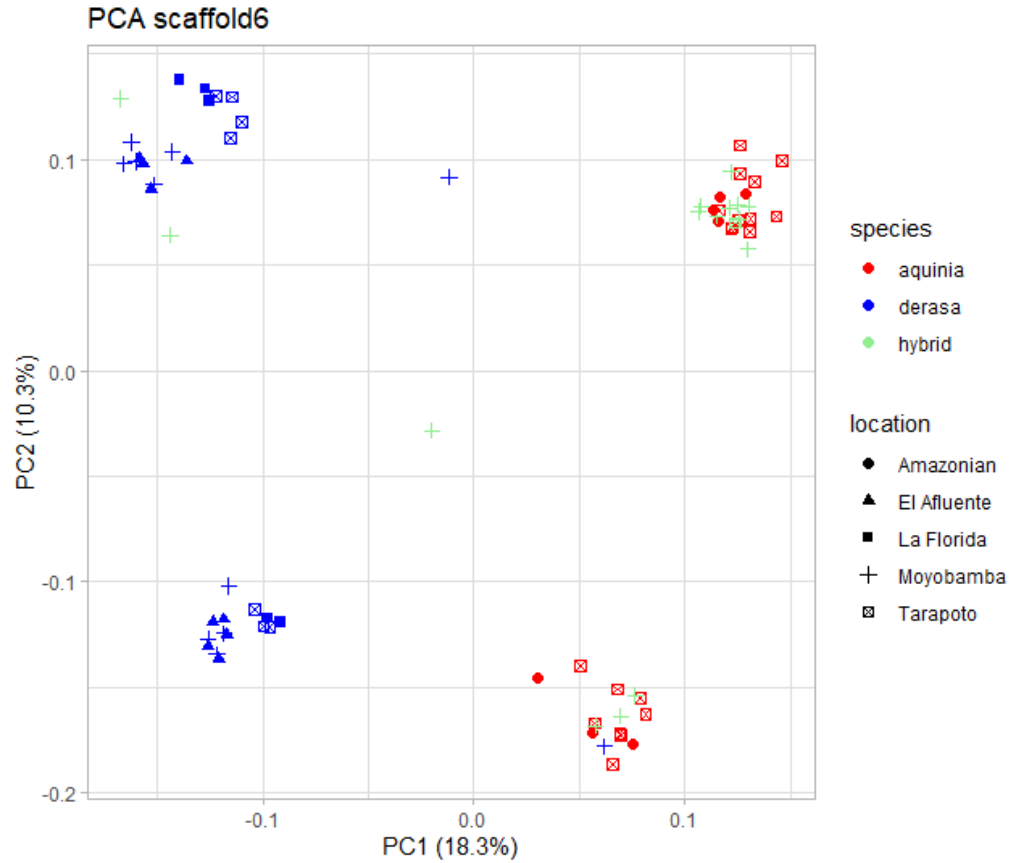
- Do any individuals show low mapping proportion?
- Do any individuals show high proportion of missing data or low sequencing depth?



Once you get your data, make sure you trust it



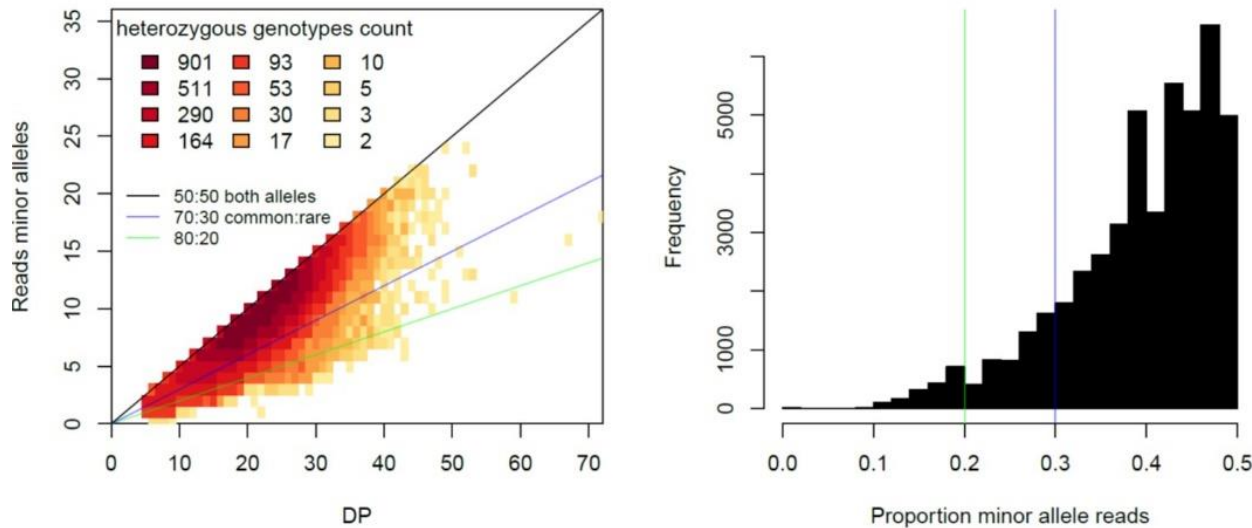
Once you get your data, make sure you trust it



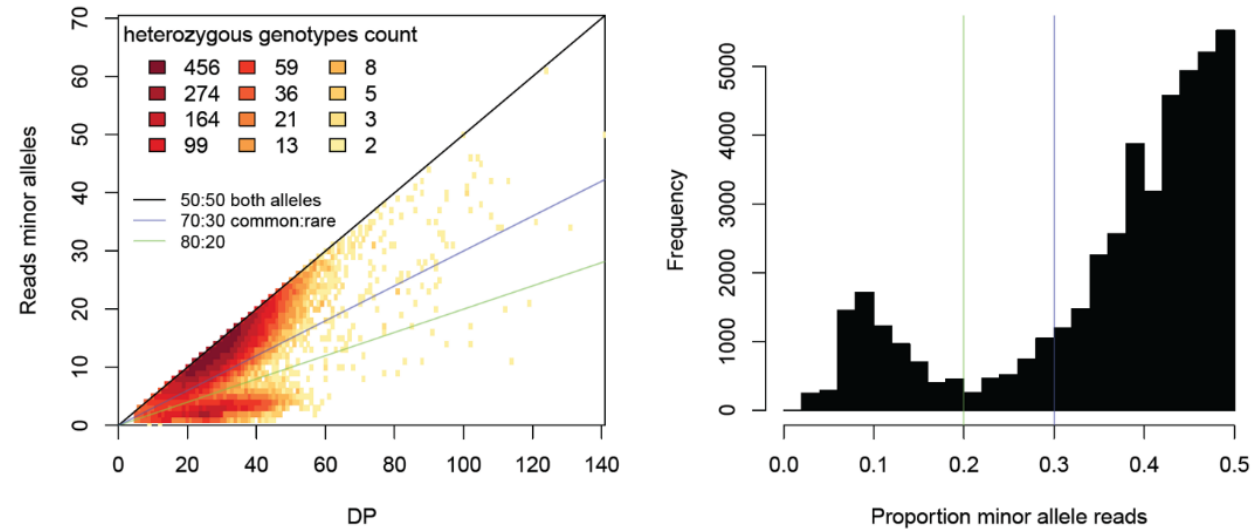
Once you get your data, make sure you trust it

Contamination can cause issues, including erroneous inference of introgression

well-sequenced individual

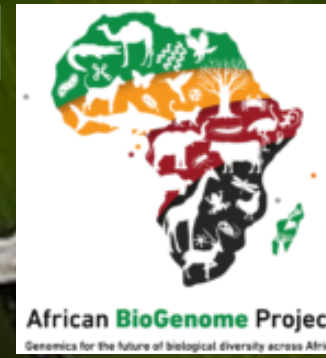


Individual with likely contamination





Project Psyche



ATLASea
Atlas des génomes marins



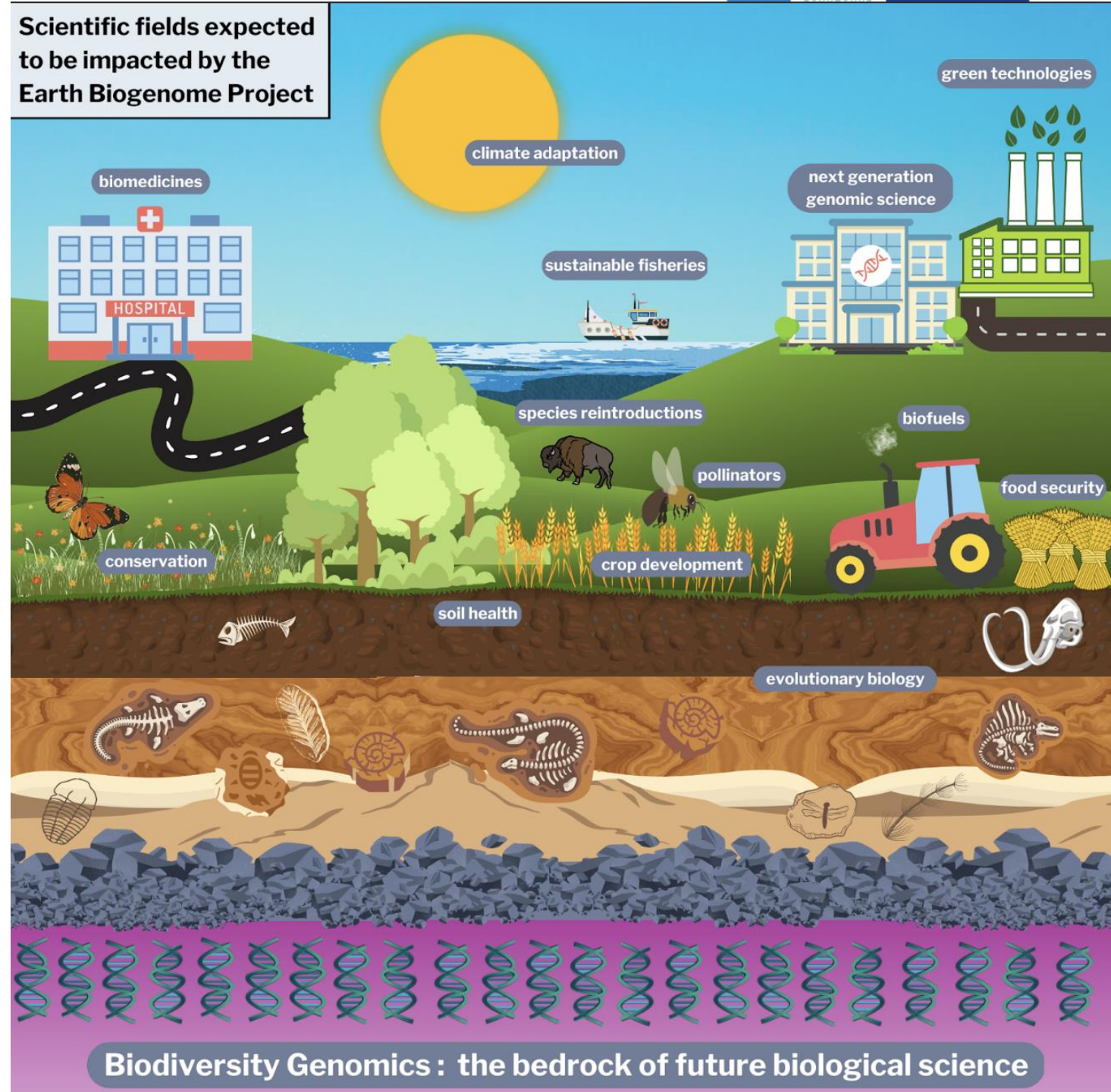
CREATING A NEW FOUNDATION FOR BIOLOGY

Sequencing Life for the
Future of Life



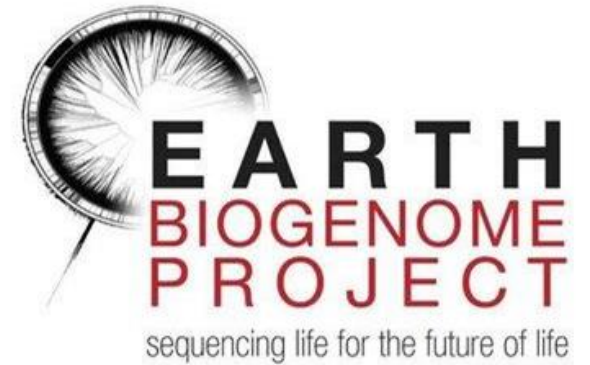
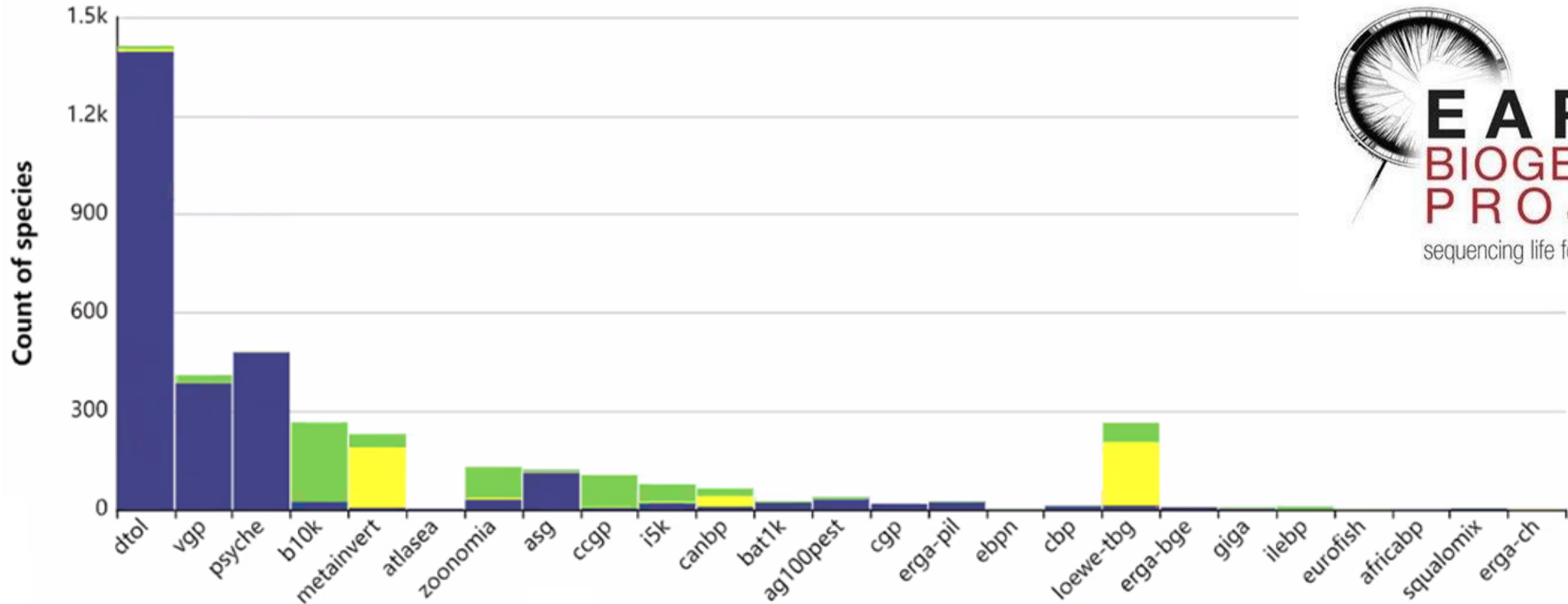
Reference genomes have global impact in a wide range of areas important for planetary health:

- Biomedicines
- Climate adaptation
- Biomonitoring
- Green technologies
- Next generation genomic science
- Sustainable fisheries
- Species reintroductions
- Conservation
- Pollinators
- Biofuels
- Food security
- Crop development
- Soil health
- Evolutionary biology

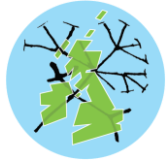


Large reference genome projects

contig scaffold chromosome complete genome



Major Projects at the Tree of Life Programme at the Wellcome Sanger Institute, Cambridge, UK



★ Darwin Tree of Life Project

- 70,000 species from Britain and Ireland [funded for ~8,000 species to 2026 currently]



★ Aquatic Symbiosis Genomics

- 1,000 species (500 symbiotic systems) from marine and freshwater



★ Project Psyche – Lepidoptera genomes for Europe

- Sequencing all 11,000 moths and butterflies of Europe



★ Vertebrate Genomes Project

- VGP Phase 1 (ordinal - 260 species) and Phase 2 (family) goals




★ European Reference Genome Atlas

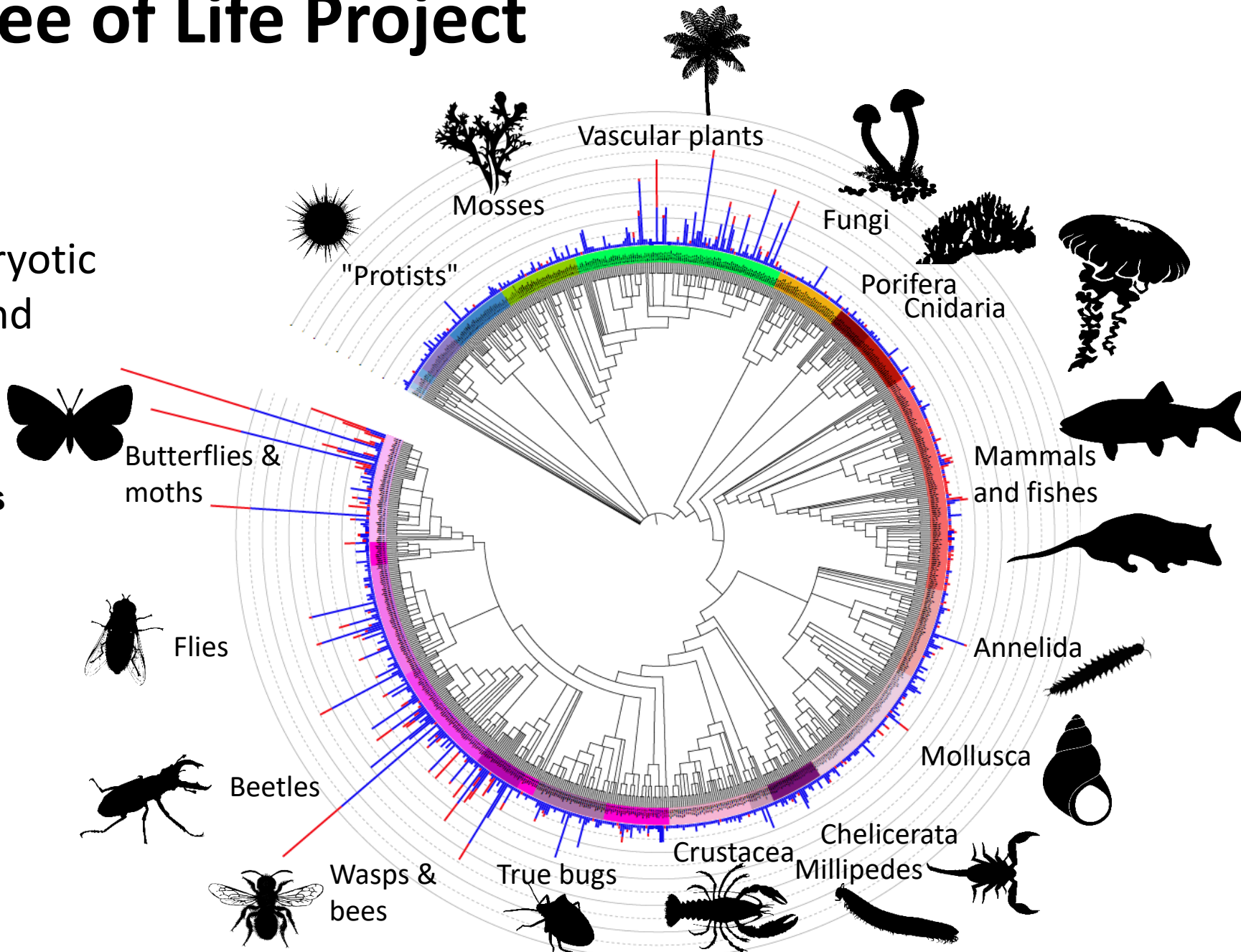
- Sequencing the genomes of all species in the European continent (~500 species)



Darwin Tree of Life Project

sequence all ~70,000 eukaryotic species in Britain and Ireland

 >5000 species in progress
 >1500 genomes complete





Project Psyche

Generate and explore reference genomes for all **~11,000 species** of **butterflies and moths (Lepidoptera)** of **Europe**.

Named after Psyche -
Greek goddess of the soul,
often depicted with butterfly
wings



Credit: Sidon, Sarcophagus relief of Psyche. Livius.org

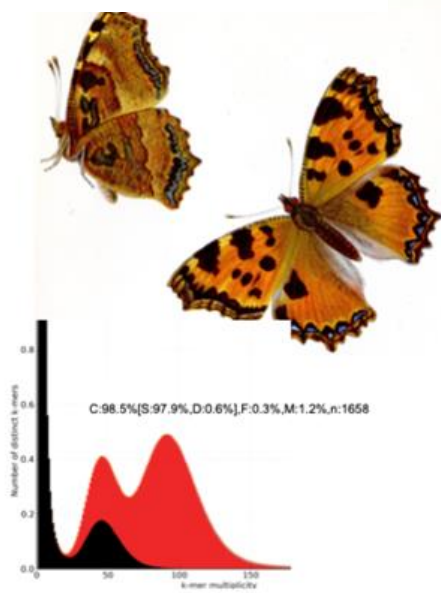


Why Lepidoptera?

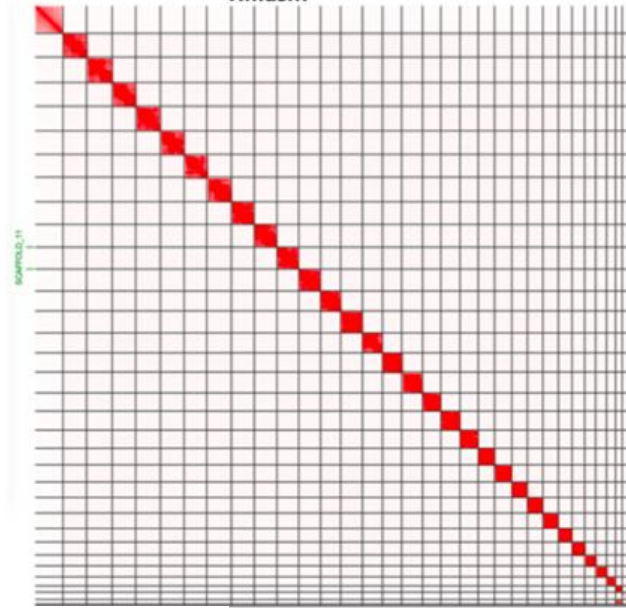


Lepidoptera genomes tend to assemble very easily

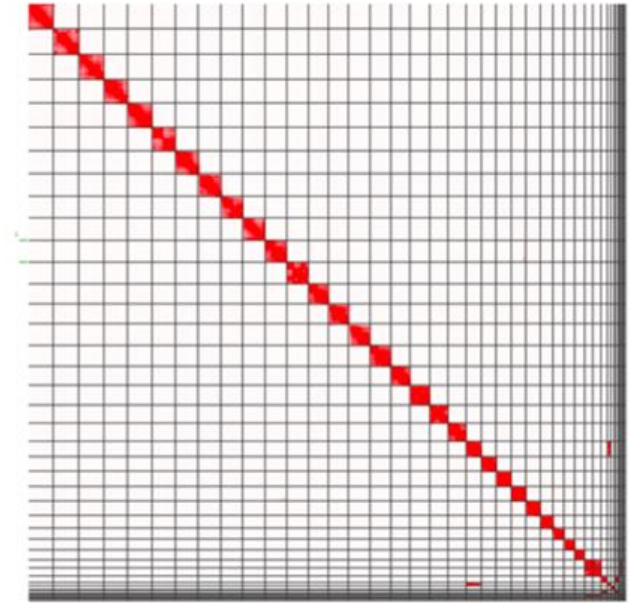
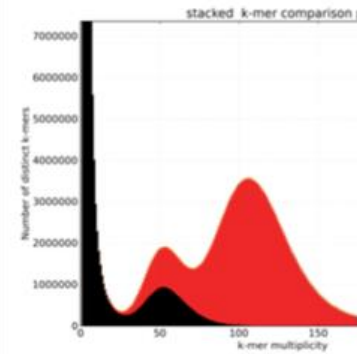
Nymphalis polychloros



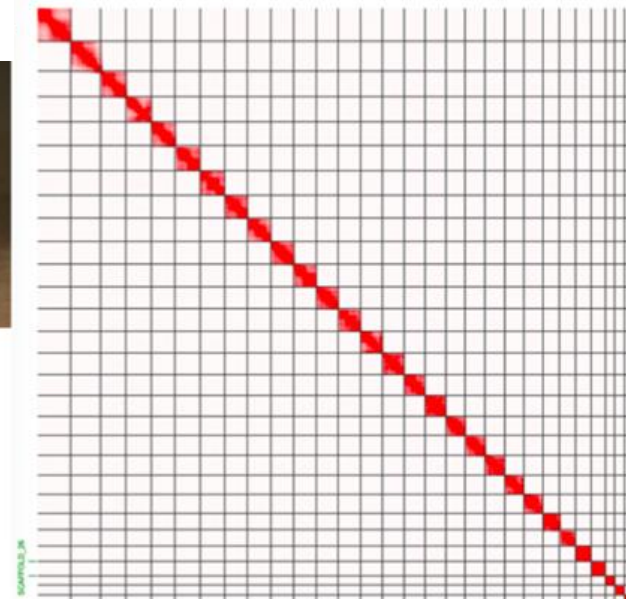
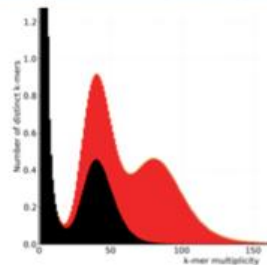
Hifiasm



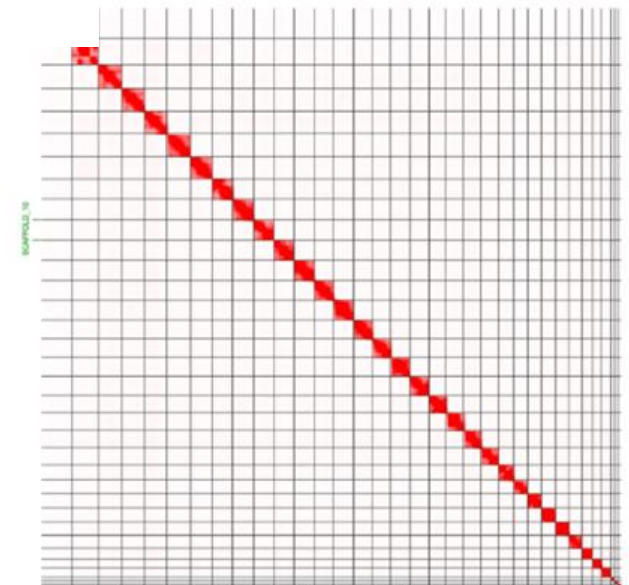
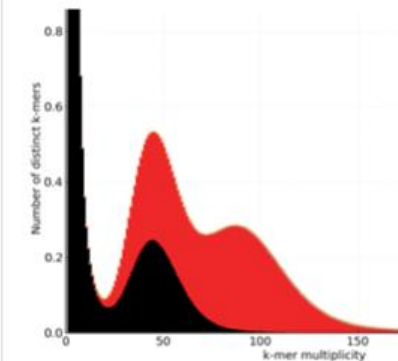
Pheosia tremula



Deilephila porcellus



Colias croceus



Why sequence all species of Lepidoptera in Europe?

To answer wide-reaching questions:

- Chromosome evolution
- Evolution of opsin genes
- Genetic basis of migration
- Population genetics
- Genetic diversity predictors
- & many more...

nature ecology & evolution

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature ecology & evolution](#) > [articles](#) > article

Article | [Open access](#) | Published: 21 February 2024

Comparative genomics reveals the dynamics of chromosome evolution in Lepidoptera

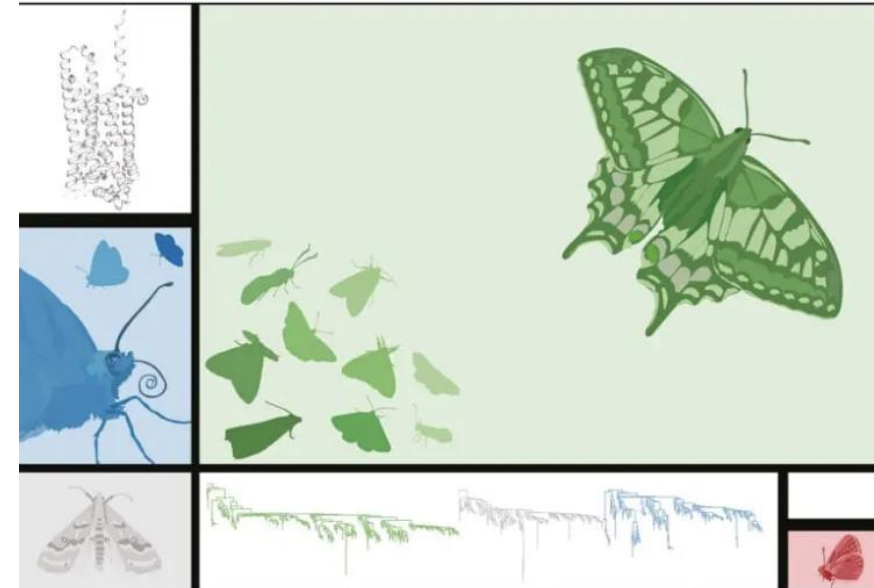
[Charlotte J. Wright](#) ✉, [Lewis Stevens](#), [Alexander Mackintosh](#), [Mara Lawniczak](#) & [Mark Blaxter](#) ✉

[Nature Ecology & Evolution](#) **8**, 777–790 (2024) | [Cite this article](#)

16k Accesses | 13 Citations | 255 Altmetric | [Metrics](#)

MOLECULAR BIOLOGY AND EVOLUTION

academic.oup.com/mbe



Opsin Gene Duplication in Lepidoptera: Retrotransposition, Sex Linkage, and Gene Expression

[Peter O Mulhair](#) ✉, [Liam Crowley](#), [Douglas H Boyes](#), [Owen T Lewis](#), [Peter W H Holland](#) [Author Notes](#)

Molecular Biology and Evolution, Volume 40, Issue 11, November 2023,

Psyche community



Charlotte Wright



Mark Blaxter

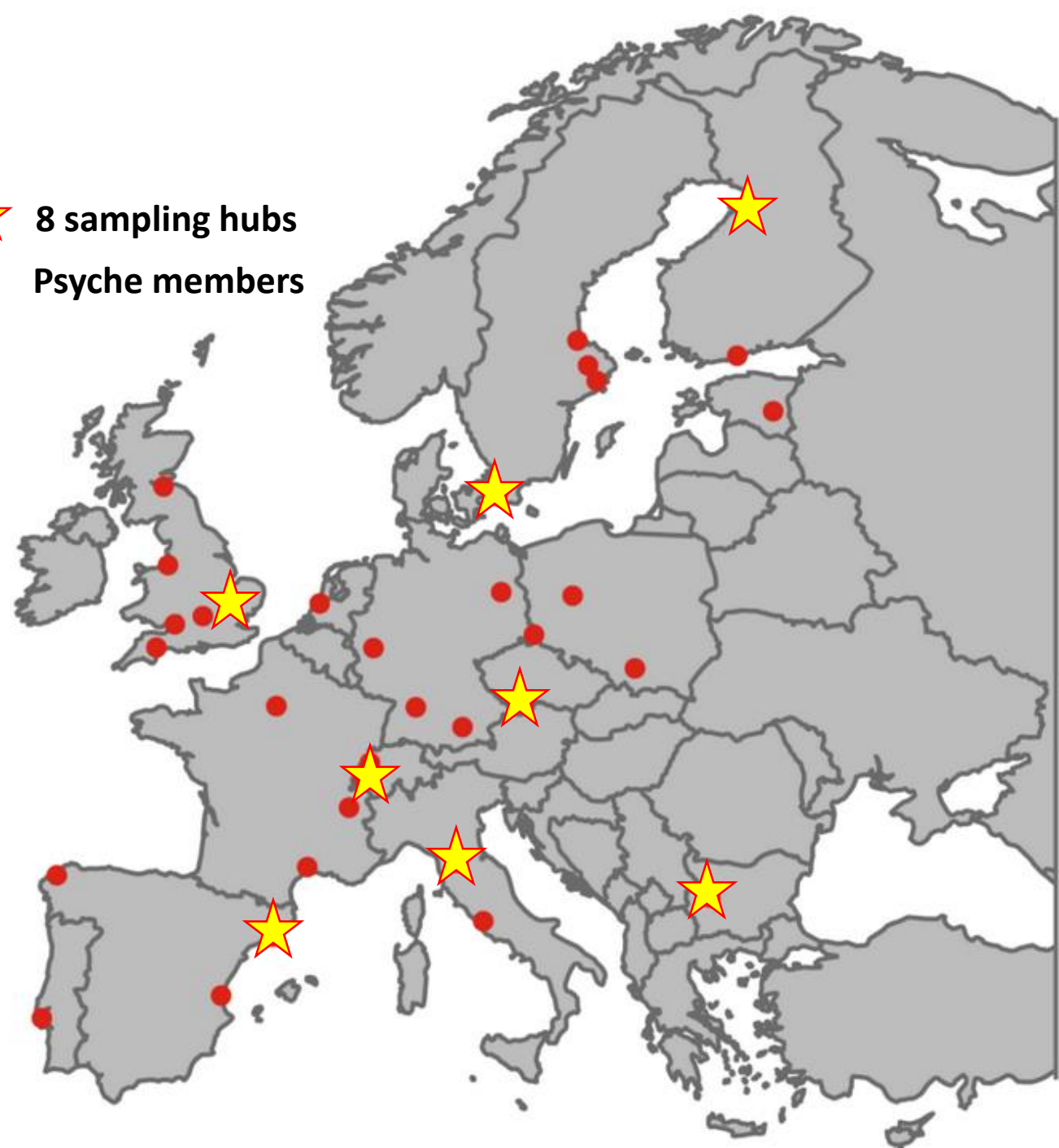


Joana Meier

Join us if you are interested
in butterflies and moths



- ★ 8 sampling hubs
- Psyche members





Phases of Project Psyche

Generate and explore reference genomes for all **~11,000 species** of **butterflies and moths (Lepidoptera)** of **Europe**.

Phase 1: Generate 2,000 genomes

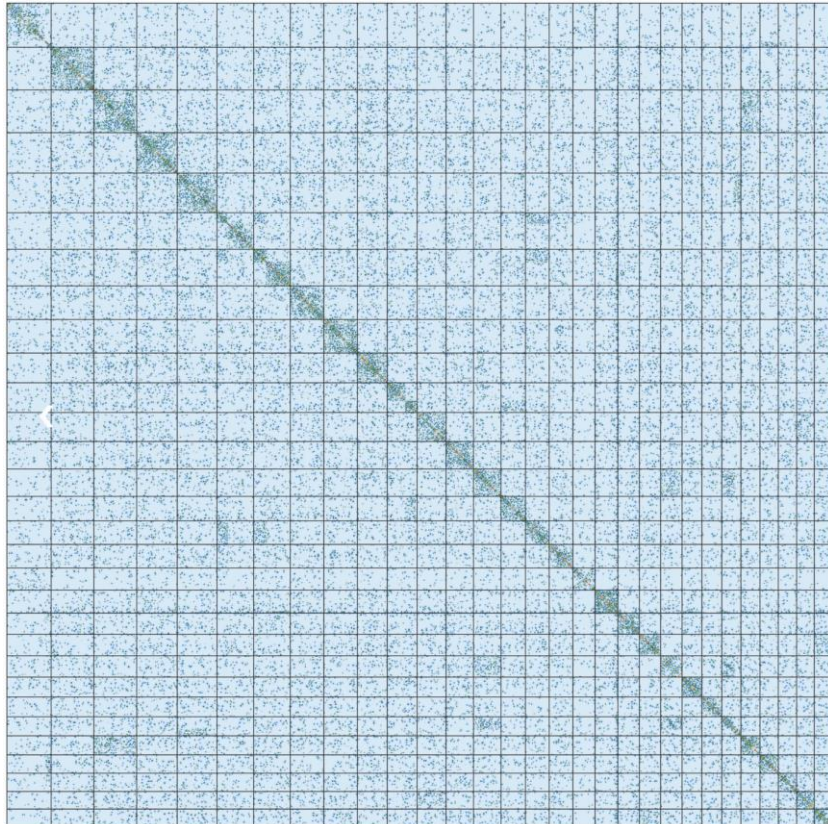
- Focus on taxonomic diversity
- Samples from sampling hubs
- Sequencing at the Wellcome Sanger Institute
- Writing a whitepaper to celebrate the first 1000 genomes
- Writing a set of 3-5 publications on the first 1000 genomes

Plans for later phases:

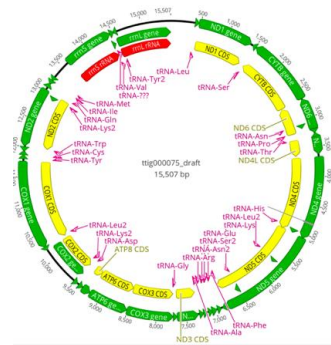
- Establishing more sampling hubs
- Establishing more sequencing hubs

Genome assembly of the moth and all its co-bionts

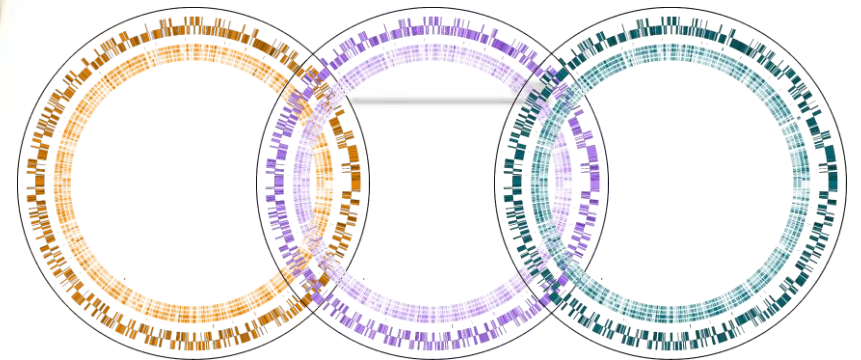
***Phalera bucephala*:**
one moth, five genomes



31 nuclear chromosomes



mitochondrial genome



Wolbachia A, B1, B2

Openly published genomes and Genome Notes: Collector(s) as first author(s)

<https://www.darwintreeoflife.org/genomes/genome-notes>

Wellcome Open Research


Wellcome Open Research 2023, 8:278 Last updated: 18 SEP 2023



DATA NOTE

The genome sequence of the Dark Spectacle, *Abrostola*

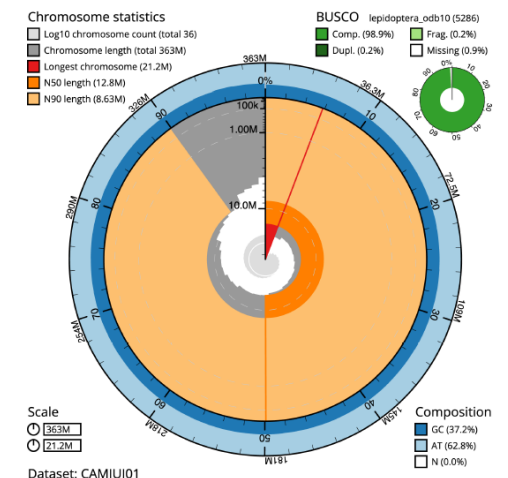
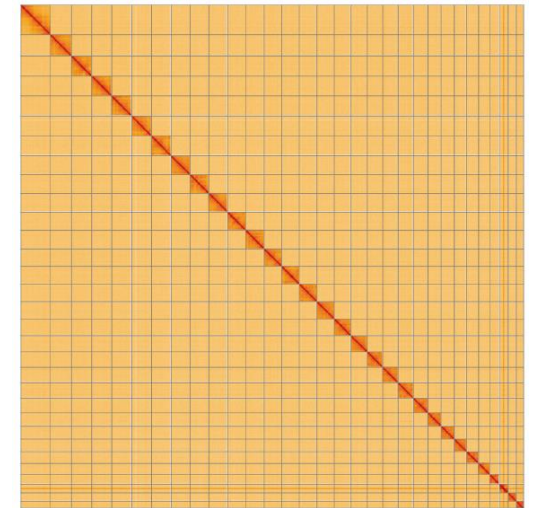
triplasia (Linnaeus, 1758) [version 1; peer review: 2 approved]

Douglas Boyes¹⁺, Owen T. Lewis ²,

University of Oxford and Wytham Woods Genome Acquisition Lab,
Darwin Tree of Life Barcoding collective,
Wellcome Sanger Institute Tree of Life programme,
Wellcome Sanger Institute Scientific Operations: DNA Pipelines collective,
Tree of Life Core Informatics collective, Darwin Tree of Life Consortium

¹UK Centre for Ecology & Hydrology, Wallingford, England, UK

²University of Oxford, Oxford, England, UK



Psyche enables LepEU: population genomics across Lepidoptera

LepEU: European Lepidopteran
Population Genomics Consortia

<https://lepeu.github.io/>

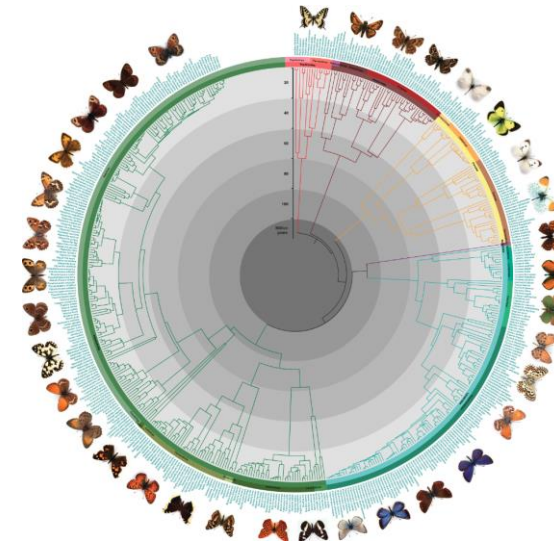
*Quantifying regional genetic diversity, inbreeding,
and local adaptation of leps to inform biodiversity
management in Europe*



Project Psyche
Reference genomes of
all European Lepidoptera



Patrícia Beldade
Lisbon University, Portugal



Chris Wheat
Stockholm University,
Sweden

COST action:

Lep10KGenomes: Utilizing 10,000 genomes of European Lepidoptera

- Organise workshops and training courses
- Research exchanges for early career researchers
- Work with external stakeholders to address societal needs using genomic information
- Devise best practices
 - For reference genomes - **Psyche**
 - For population genomics - **LepEU**



Niklas Wahlberg
Lund University,
Sweden



Funded by the Horizon 2020 Framework Programme
of the European Union

Find out more here:

<https://tinyurl.com/Lep10k-COST>

Join Project Psyche here:

<https://tinyurl.com/projectpsyche>

**How to generate a
reference genome?**

Check if a reference genome exists



Genomes on a Tree ([GoaT, https://goat.genomehubs.org/](https://goat.genomehubs.org/))

Genomes on a Tree (GoaT)

Genomes on a Tree (GoaT): A versatile, scalable search engine for genomic and sequencing project metadata across the eukaryotic tree of life. Challis *et al.* 2023. Wellcome Open Res 2023, **8**:24 [doi:10.12688/wellcomeopenres.18658.1](https://doi.org/10.12688/wellcomeopenres.18658.1)

GoaT needs been built using [GenomeHubs](#) to help coordinate efforts across the [Earth Biogenome Project](#) (EBP) Network at all stages from planning through sequencing and assembly to publication. [read more...](#)

Search GoaT

Type to search GoaT taxon index (e.g. Canidae)

include descendants include estimates empty columns result columns query builder

Off On Off

First, check if a reference genome already exists

- Genomes on a Tree ([GoaT, https://goat.genomehubs.org/](https://goat.genomehubs.org/))
- If yes, get it e.g. from NCBI (www.ncbi.nlm.nih.gov/home/genomes)

https://www.ncbi.nlm.nih.gov/datasets/genome/GCA_917862395.2

National Library of Medicine
National Center for Biotechnology Information

Search NCBI Log In

NCBI Datasets Taxonomy **Genome** Gene Command-line tools Documentation

Genome assembly iHelSar1.2 reference

Download datasets URL **FTP** Actions

Submitted GenBank assembly	GCA_917862395.2
Taxon	Heliconius sara
WGS project	CAKJTVO2
Assembly type	haploid
Submitter	WELLCOME SANGER INSTITUTE
Date	Apr 6, 2023

View the legacy Assembly page

BLAST the reference genome

https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/917/862/395/GCA_917862395.2_iHelSar1.2/

Index of /genomes/all/GCA/917/862/395/GCA_917862395.2_iHelSar1.2

Name	Last modified	Size
Parent Directory	-	-
GCA_917862395.2_iHelSar1.2_assembly_report.txt	2023-04-18 17:05	34K
GCA_917862395.2_iHelSar1.2_assembly_stats.txt	2023-04-18 17:05	31K
GCA_917862395.2_iHelSar1.2_feature_count.txt	2024-02-15 01:58	168
GCA_917862395.2_iHelSar1.2_genomic.fna.gz	2023-04-18 17:05	107M
GCA_917862395.2_iHelSar1.2_genomic.gbff.gz	2023-04-18 17:05	130M
GCA_917862395.2_iHelSar1.2_genomic.gaps.txt.gz	2023-04-18 17:05	1.4K
GCA_917862395.2_iHelSar1.2_wgsmaster.gbff.gz	2023-04-18 17:05	1.4K
README.txt	2024-04-11 16:11	54K
annotation_hashes.txt	2024-02-17 10:02	410
assembly_status.txt	2024-07-15 13:53	14
md5checksums.txt	2024-02-17 10:02	625
uncompressed_checksums.txt	2024-06-14 01:21	167

[HHS Vulnerability Disclosure](#)

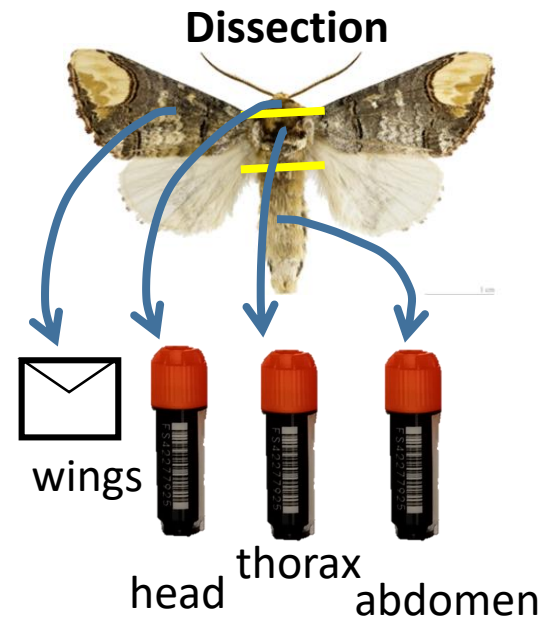
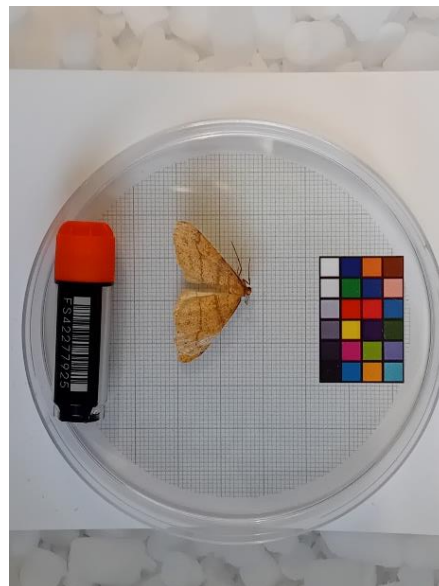
```
wget https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/917/862/395/GCA_917862395.2_iHelSar1.2/GCA_917862395.2_iHelSar1.2_genomic.fna.gz
```


1. Getting specimens with well-preserved DNA



Dry shipper (-150°C)
(with liquid nitrogen in the walls)

- Ideally:
 - Fresh tissue (or blood for birds or fishes)
 - Flash-frozen (in liquid nitrogen or cryogenic dry shipper) and kept at max -70°C
- ensure you have all metadata for the specimen (GPS coordinates, photo of the specimen, etc)



2. Sequencing the specimen

For Project Psyche
we use these tissues:



thorax

PacBio sequencing (long reads) to assemble contigs

ATGTGTCATGGGACATATGTGTCATGGGACATGAGAGAGAGA ← contig = consensus sequence of aligned reads
GAGAGAGAATGTGTCATGGGACAT
GGACATATGTGTCATGGGACATGA
TATGTGTCATGGGACATGAGAGAGAGA

head

HiC sequencing to scaffold contigs together

GTGTCATG ATGTGTCATG ATGAGAGAG GAGAGGGA
ATGTGTCATGGGACATATGTGTCATGGGACATGAGAGAGAGAGAGA GAGAGAGGGACATAAGAGTGTGTCATG
ATGTGTCATGGGACATATGTGTCATGGGACATGAGAGAGAGAGAGANNNGAGAGAGGGACATAAGAGTGTGTCATG ← scaffold = sequence with gaps

abdomen

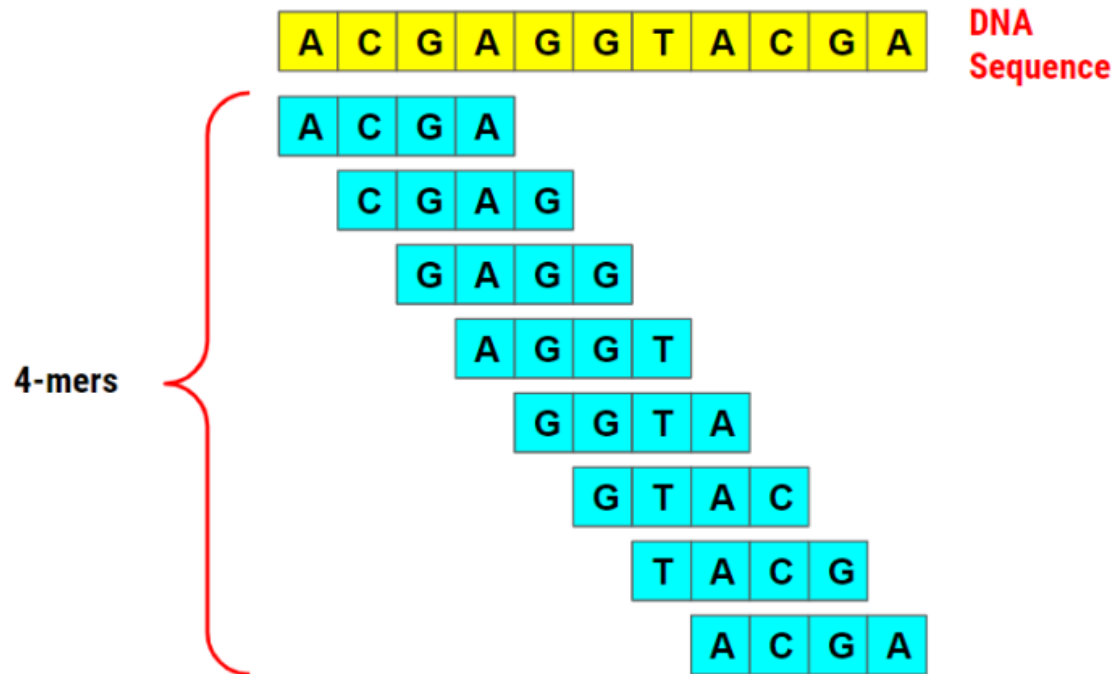
RNA sequencing for genome annotation

GTGAGTCATGGGAGAGGGACATCAAAAAAAAAAAAAAAAAA
AGAGGTCATGTGAGTACATCGATCGGAAAAAAAAAAAAAAAAA

For T2T (telomere-to-telomere) genomes,
add ONT ultra-long reads (>150 kb)

3. First statistics with kmer analyses

Estimating the genome size, depth of coverage, ploidy, levels of heterozygosity & duplication, identifying contamination

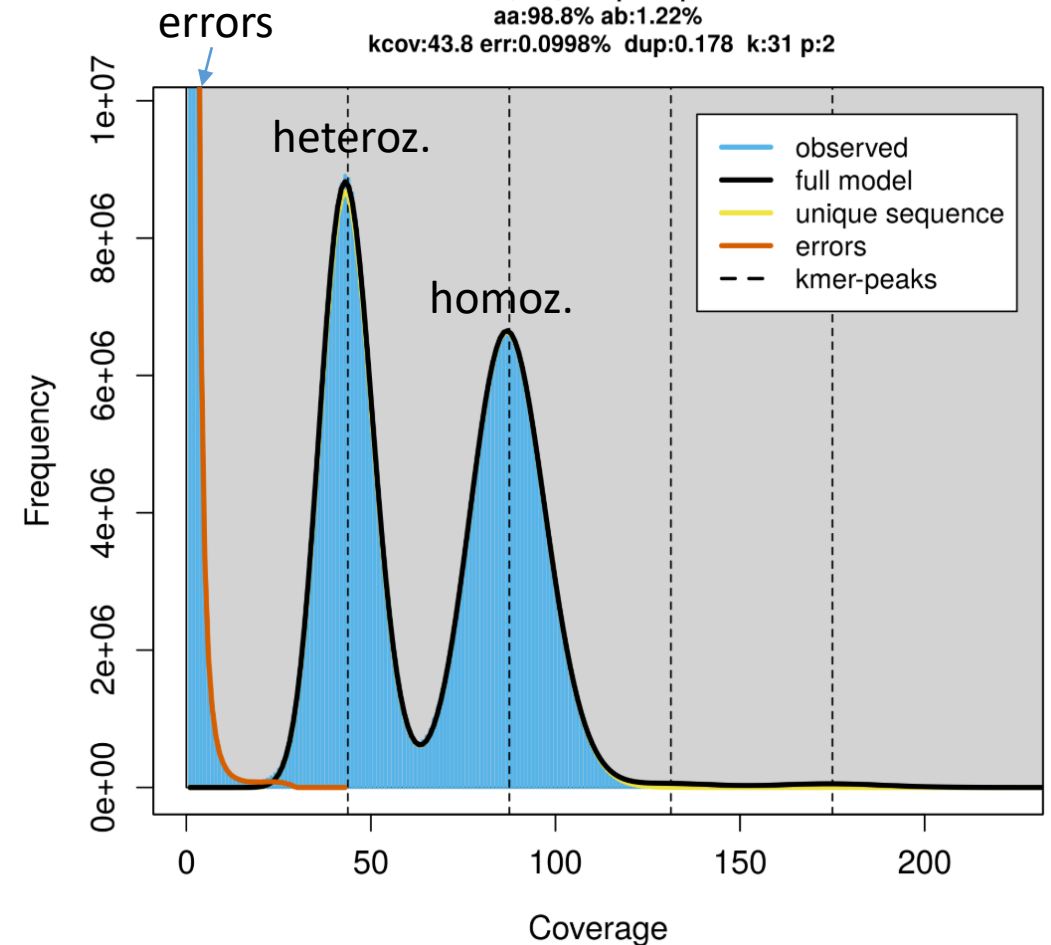


Genomescope

Histogram of frequency of k-mers

GenomeScope Profile

len:281,671,112bp uniq:87.5%
aa:98.8% ab:1.22%
kcov:43.8 err:0.0998% dup:0.178 k:31 p:2



4. Assembling the genome

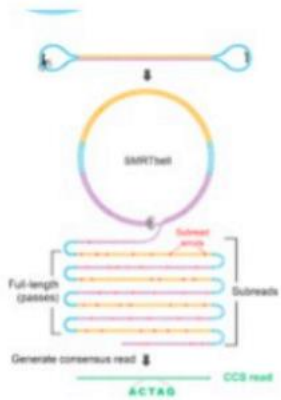
DToL Current Pipeline



For mitochondria genome assembly

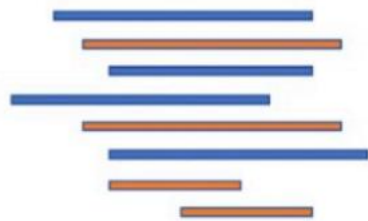
- Sequencing technologies: PacBio HiFi + HiC (Arima or Qiagen)

Slide from Marcela Uliano-Silva



Assembly

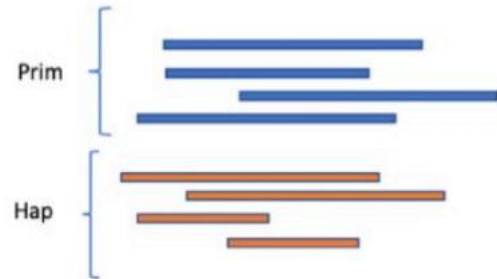
Hicanu
or Hifiasm



2 - asmstats,
BUSCO, merqury

Haplotype separation

Purge dups



3 - asmstats,
BUSCO, merqury

Scaffolding

Yahs scaffolding
(Arima or Qiagen
HiC)



4 - asmstats,
BUSCO,
merqury, HiC
heatmap

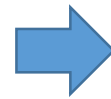
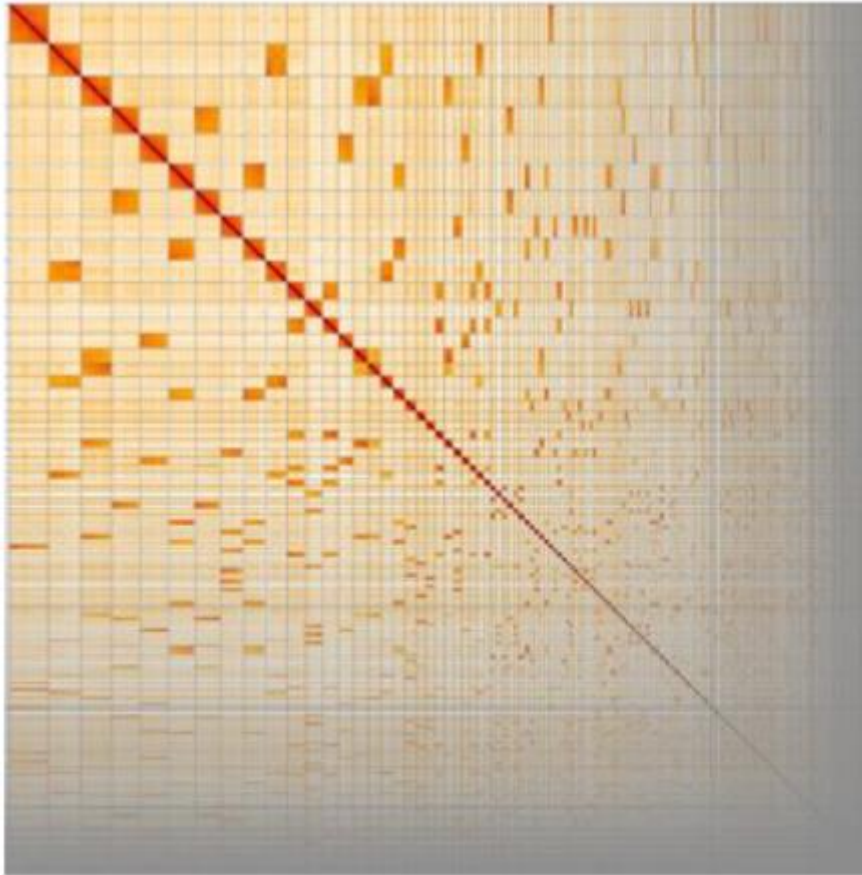
Curated assembly

5 - asmstats,
BUSCO,
merqury, HiC
heatmap

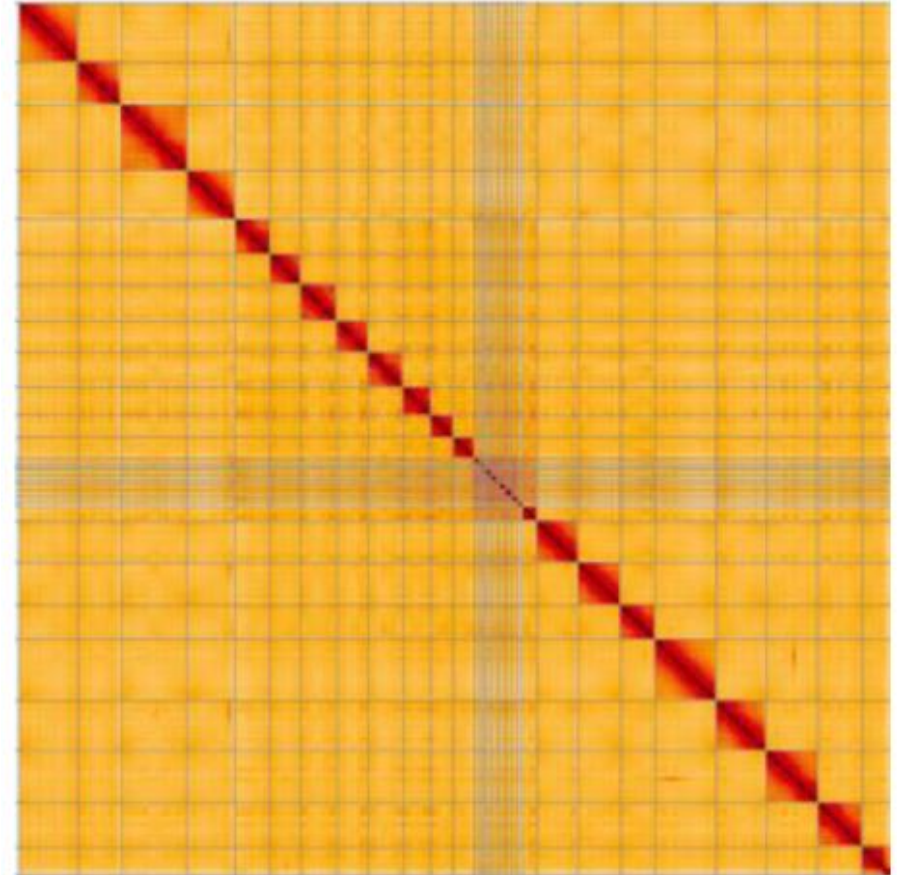
1- Kmer
Jellyfish/
GenomeScope,
asmstats,
smudgeplot

Manual curation to correct errors by assemblers

HiC map before manual curation



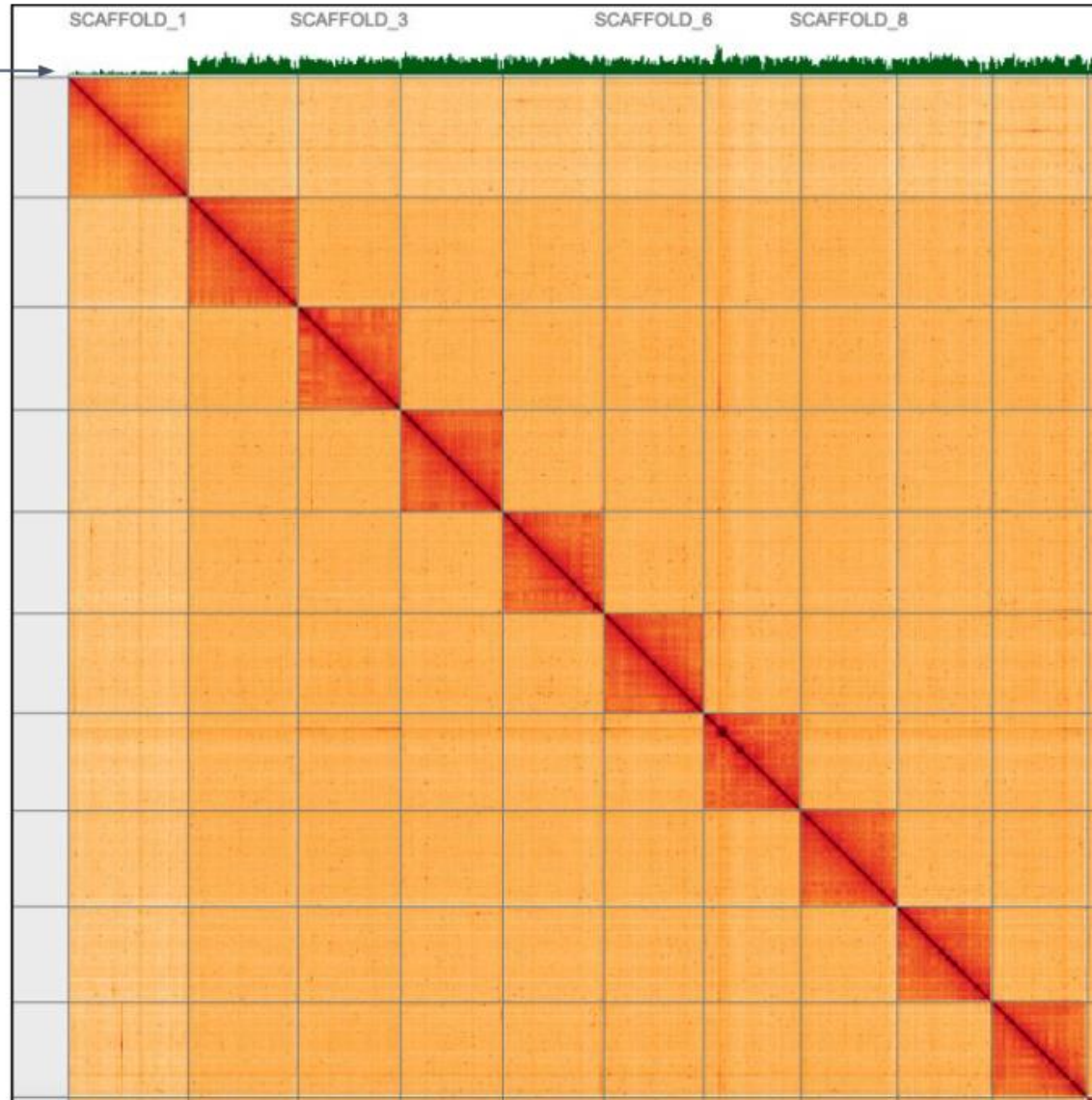
After manual curation with PretextView



https://gitlab.com/wtsi-grit/rapid-curation/-/blob/main/Interpreting_HiC_Maps_guide.pdf

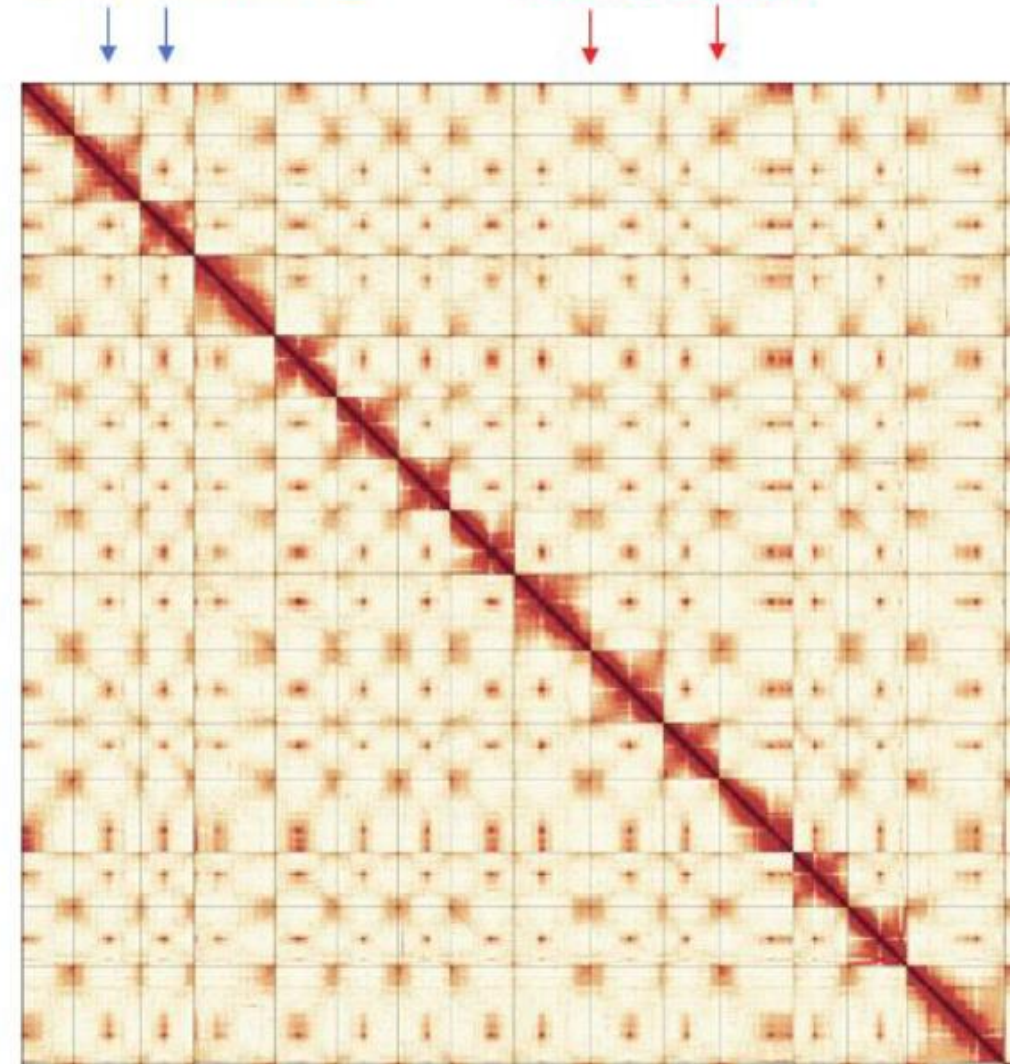
HiC maps will also show interesting features

Identifying
sex
Chromosomes
(half coverage)



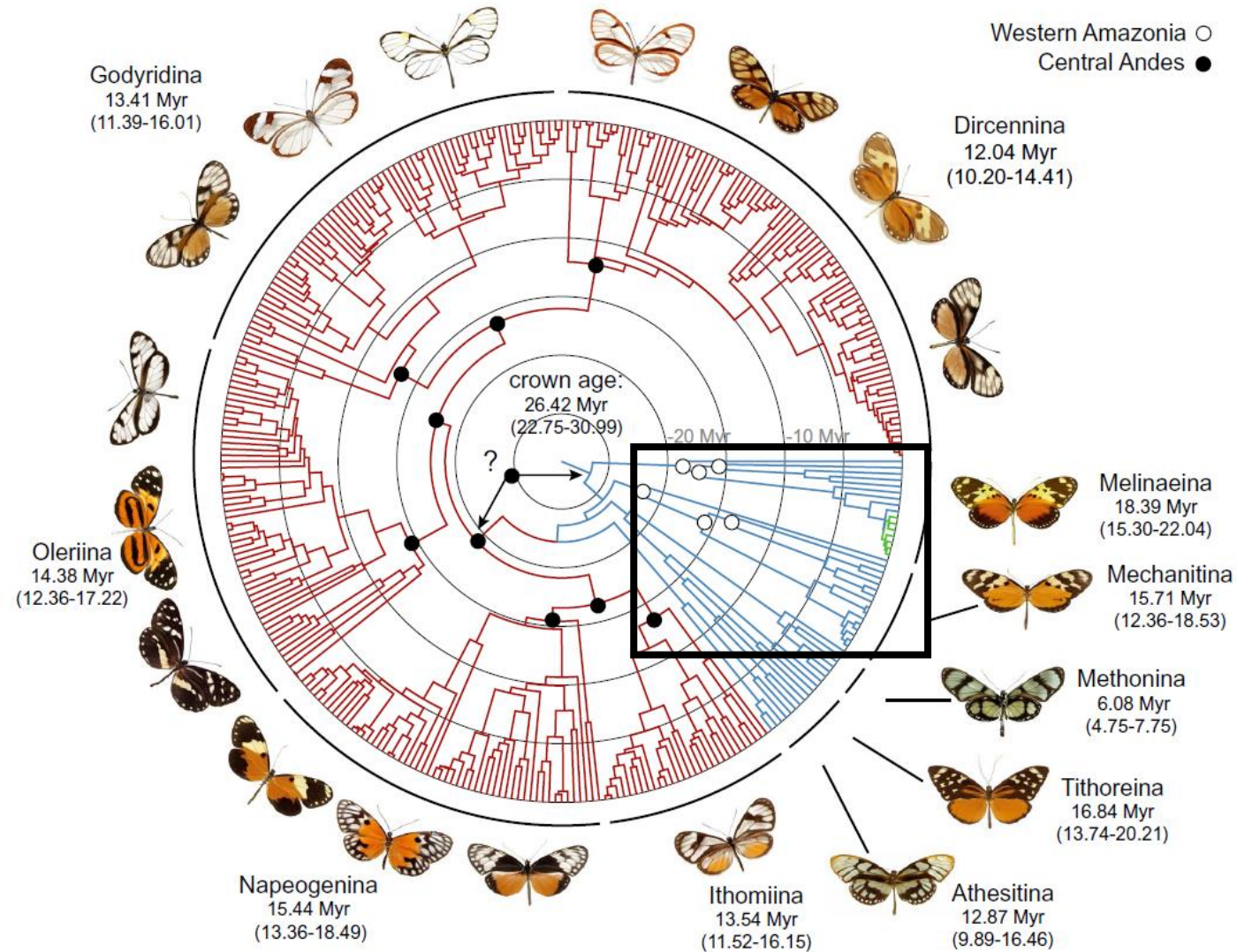
Centromeres, eg

Telomeres, eg

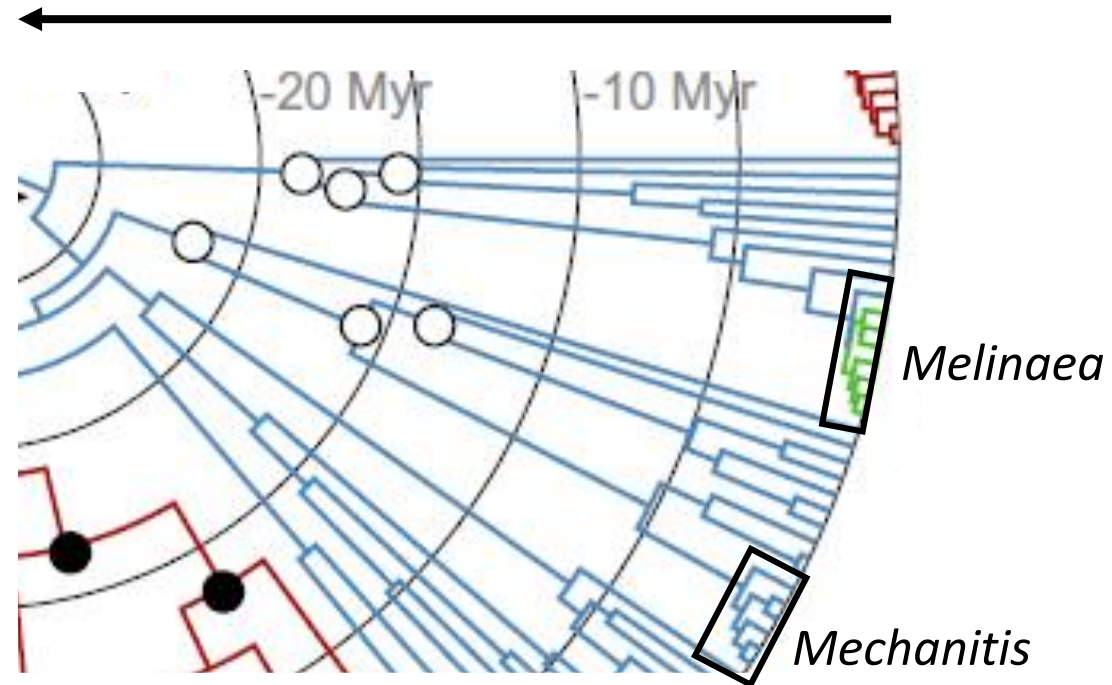


**An example of using reference
genomes from my own work**

Comparing rapidly and slowly speciating Ithomiini genera



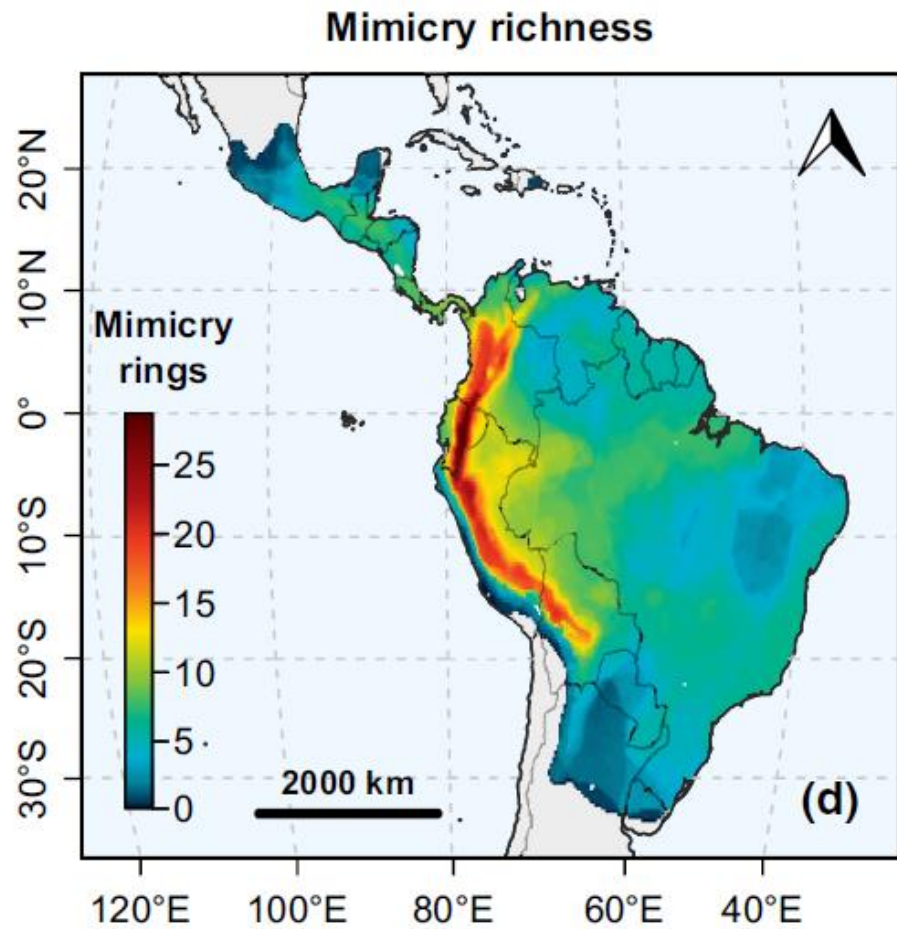
Time since last common ancestor



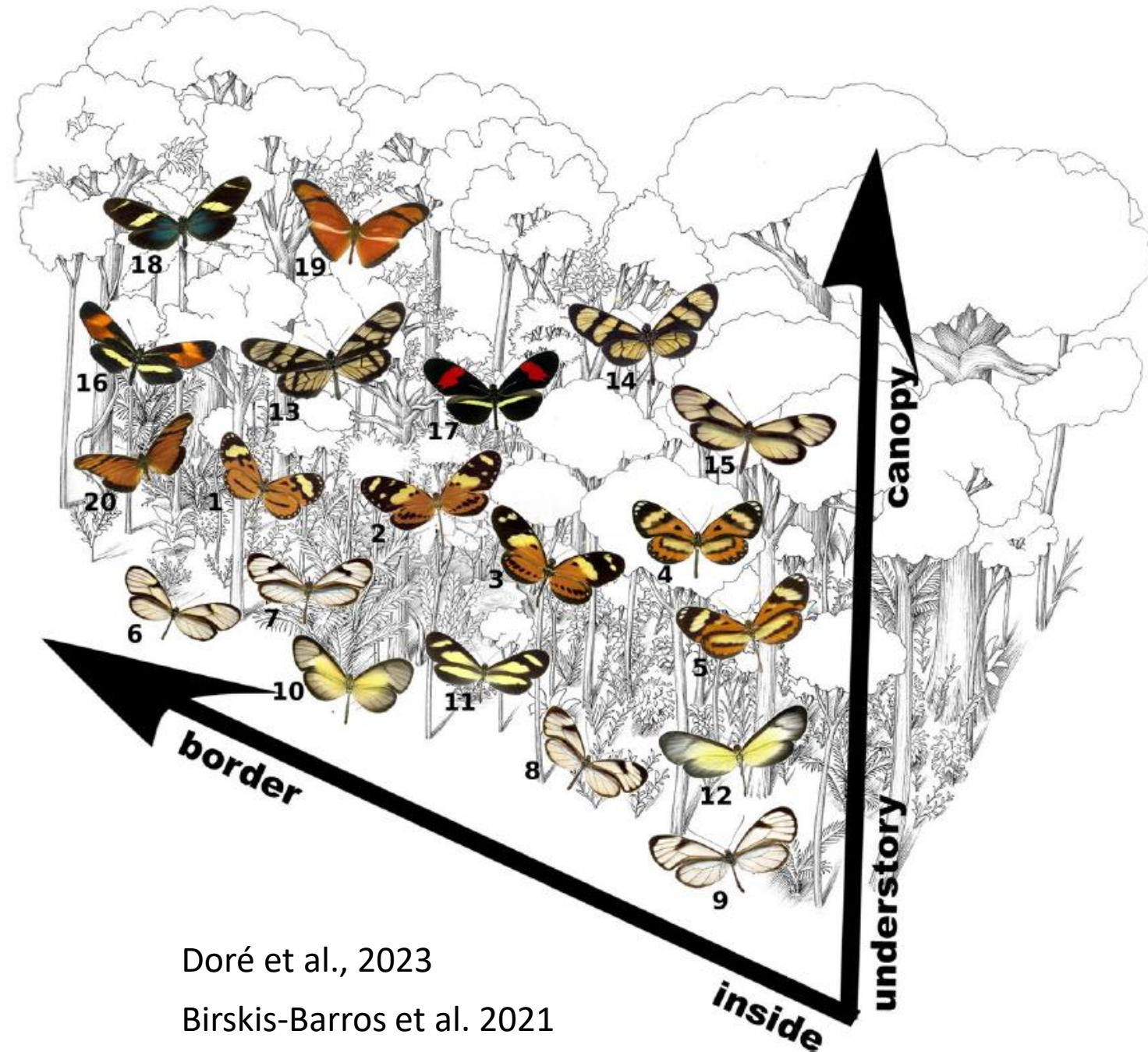
Müllerian mimicry rings of ithomiini and others



Co-mimics converge in microhabitat use and associated traits



Doré et al., 2021

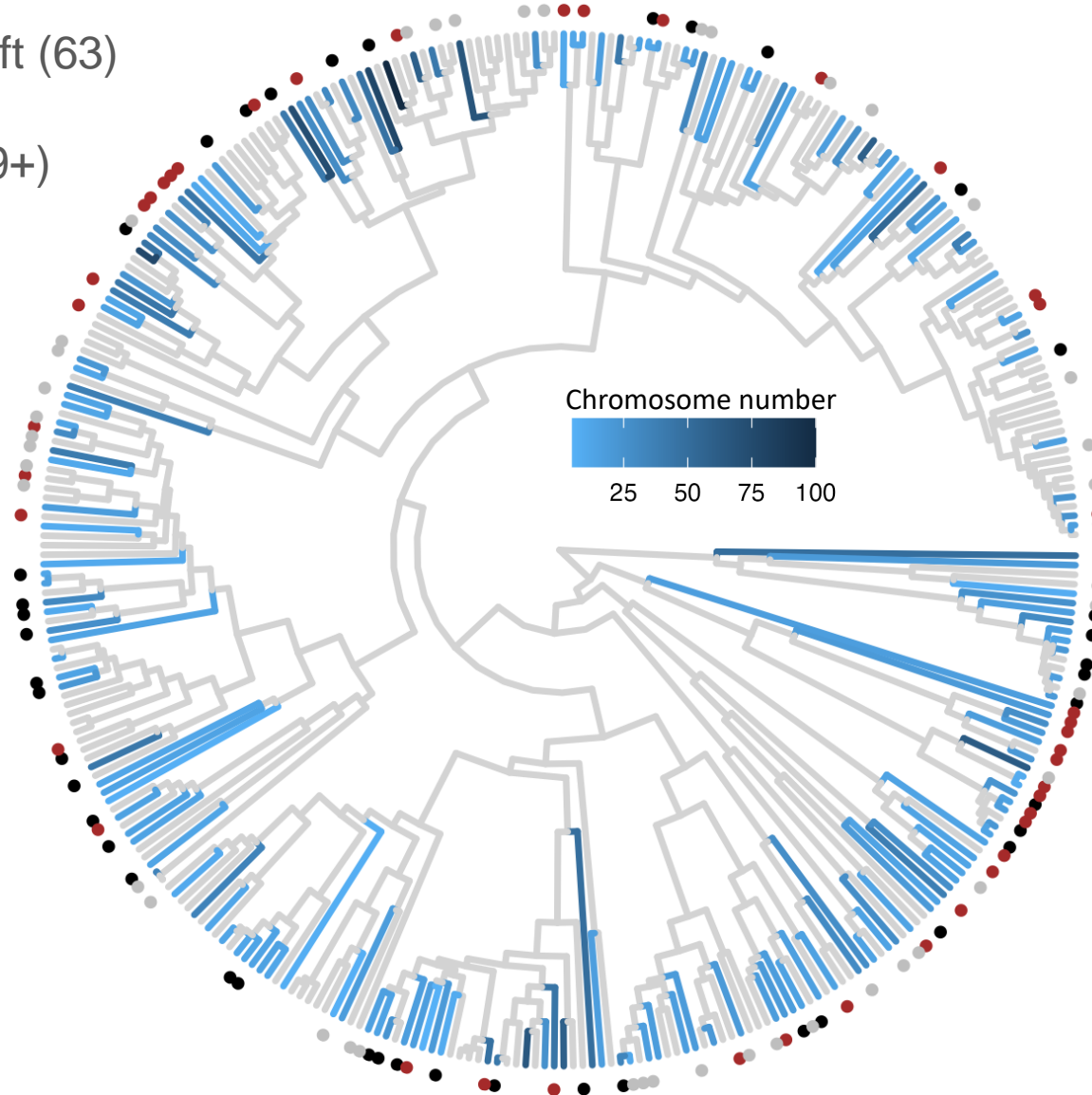


Doré et al., 2023

Birskis-Barros et al. 2021

Sequencing 200 ithomiini genomes

- Curated/Draft (63)
- Lab (47)
- Sampled (49+)



**Patricio
Salazar**



**Karin
Näsvall**

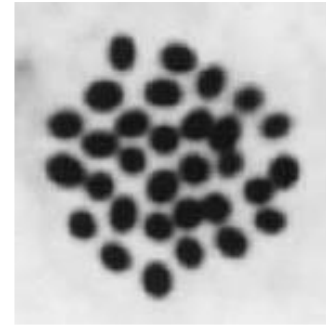
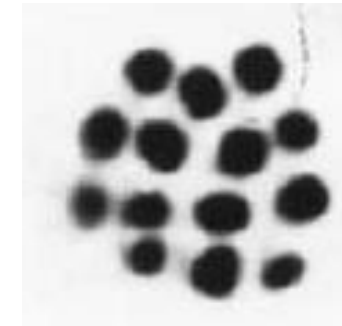
Do chromosomal rearrangements contribute to rapid diversification?

most butterflies have 30-31 chromosomes

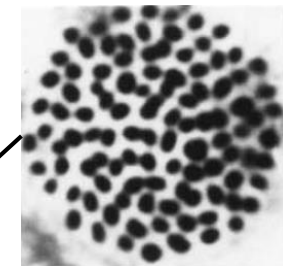
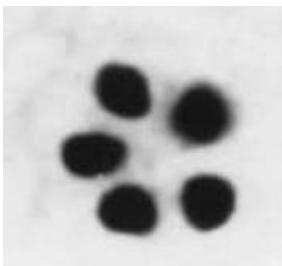
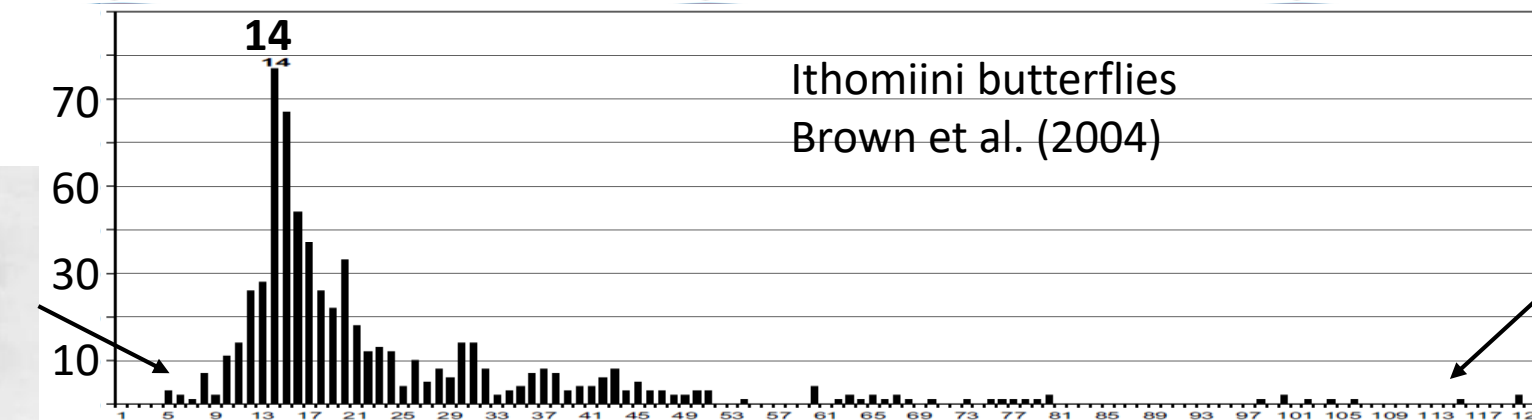
De Vos et al. (2020): 2400 Lepidoptera (including 171 ithomiini)

Melinaea species:

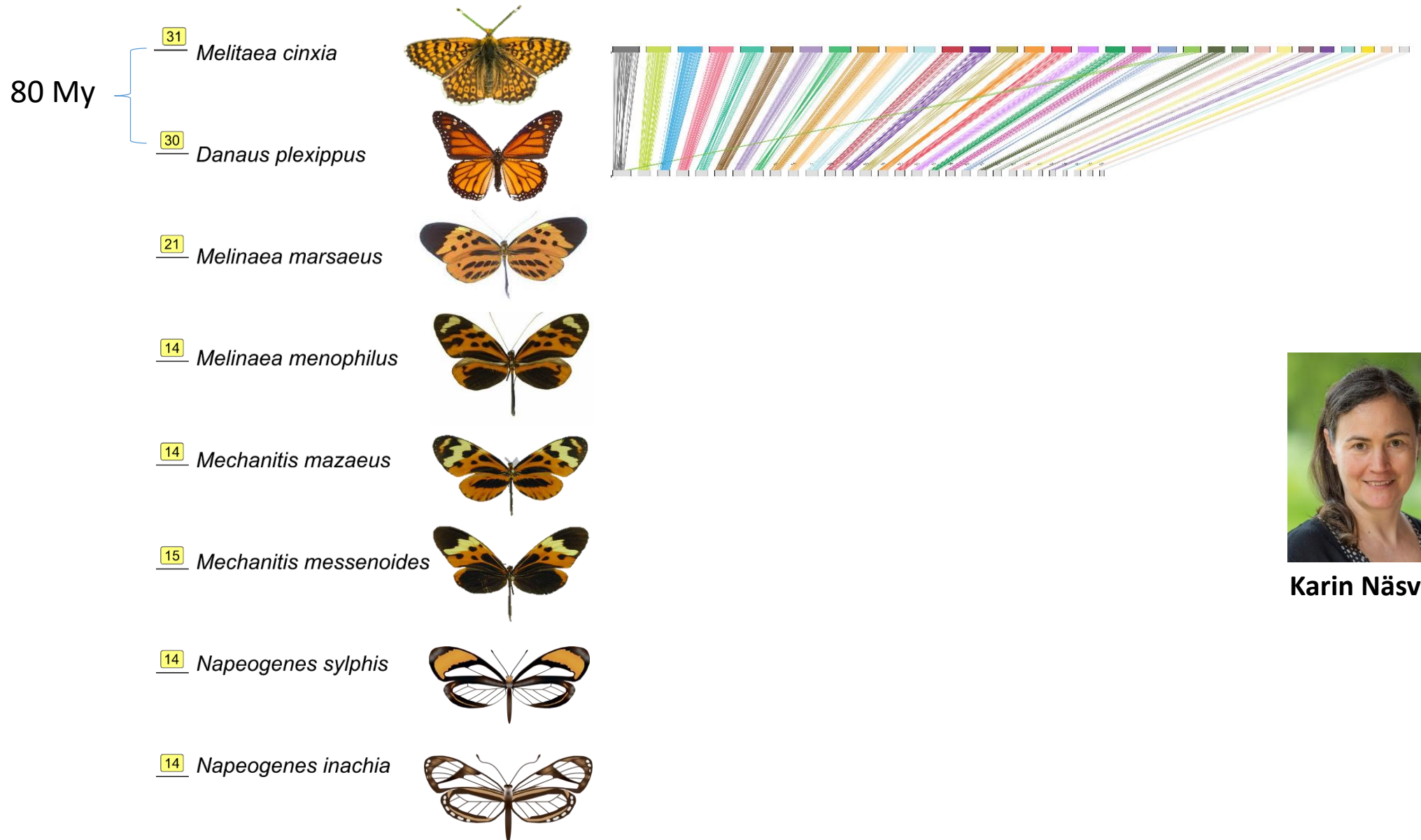
13-30 chromosomes



Keith Brown et al., 2004, Hereditas

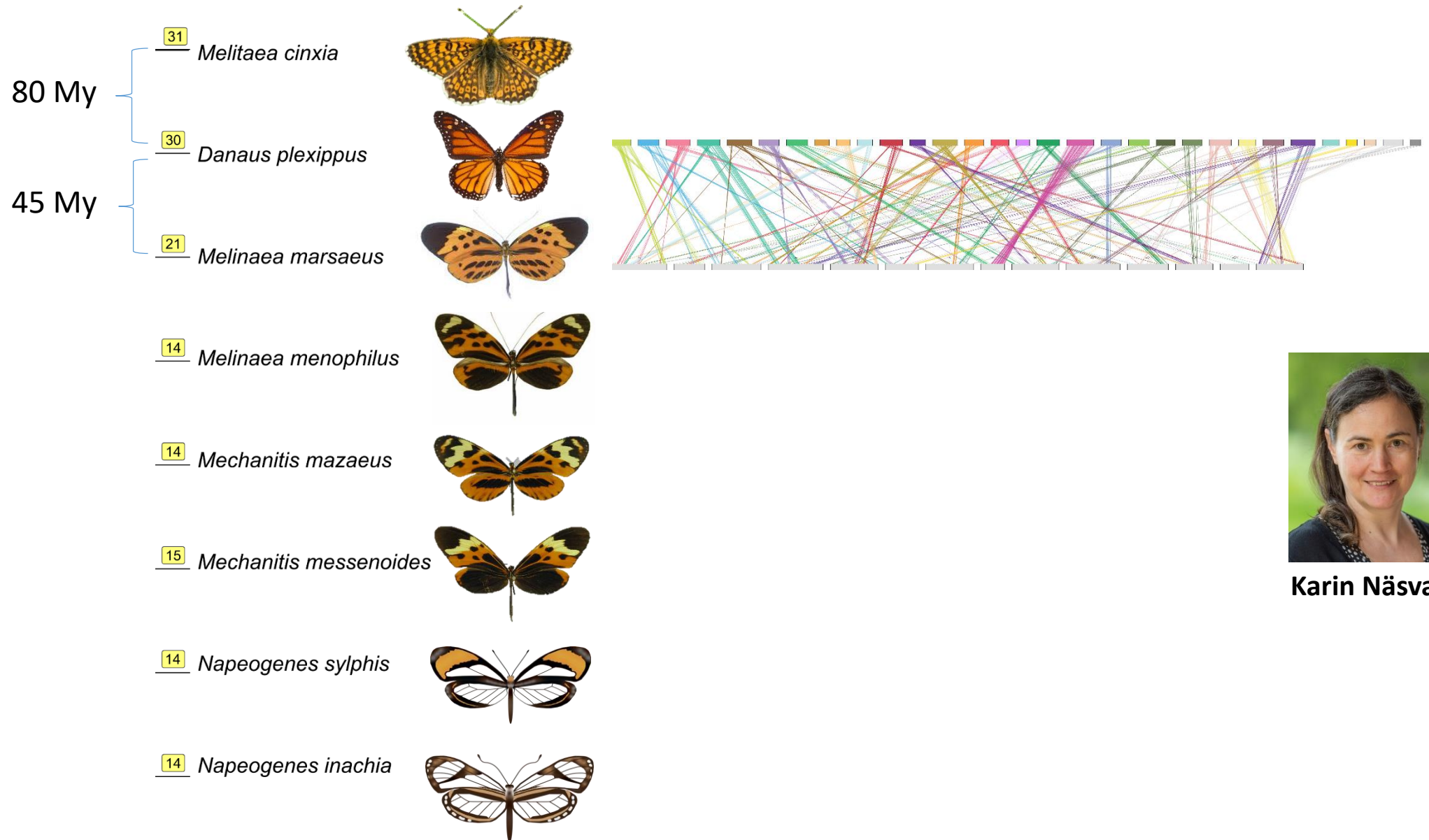


Large-scale chromosomal rearrangements



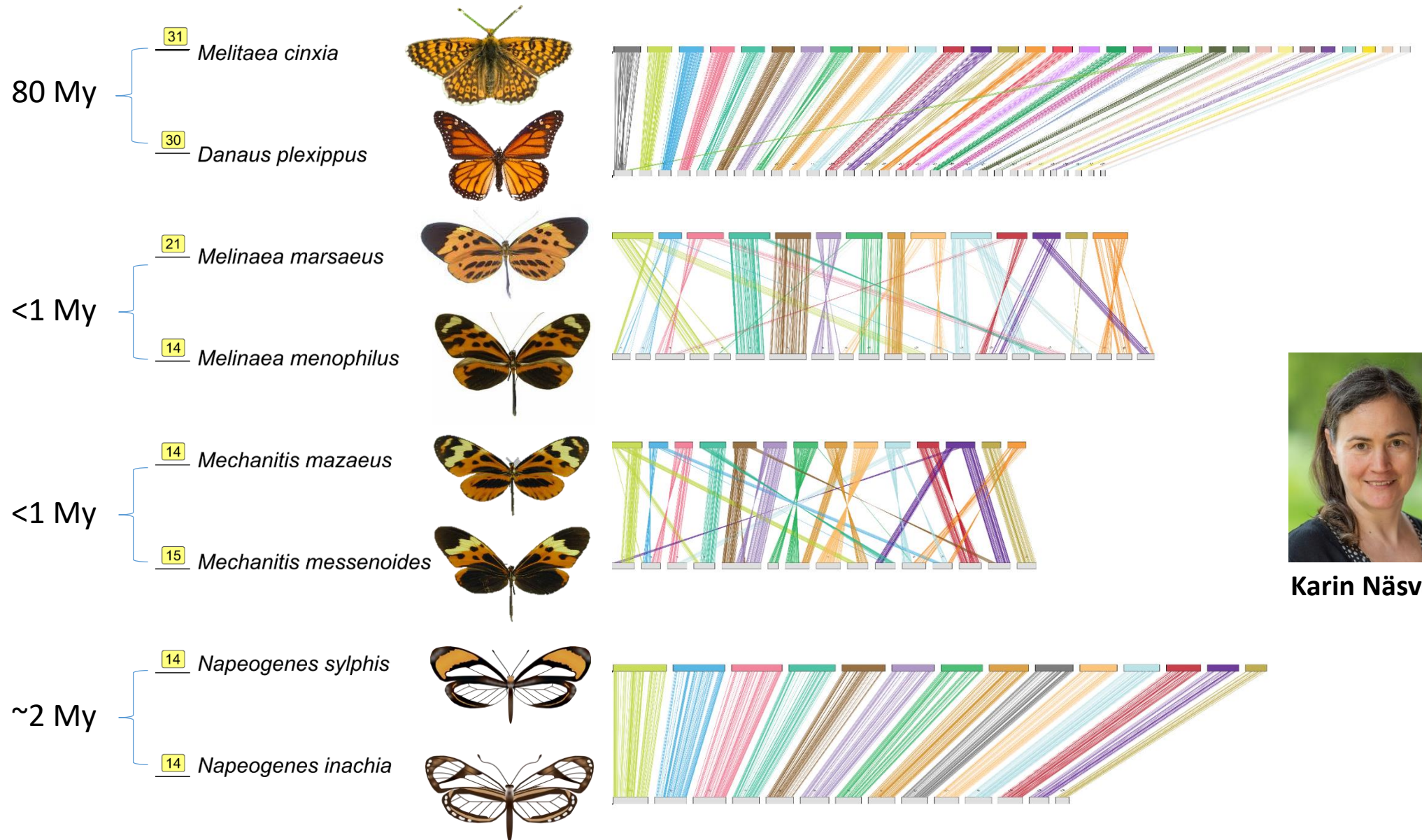
Karin Näsvall

Large-scale chromosomal rearrangements



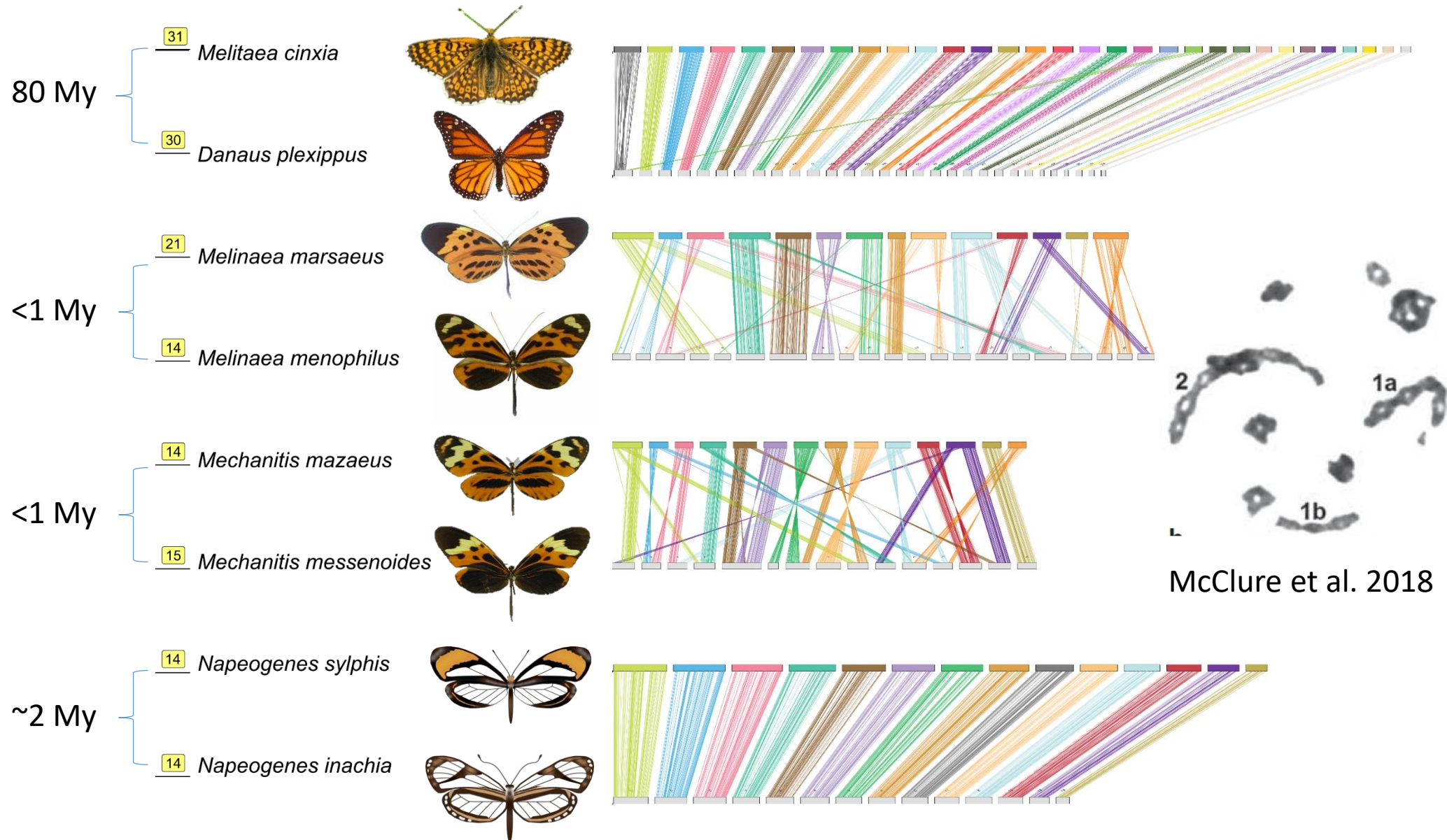
Karin Näsvall

Large-scale chromosomal rearrangements

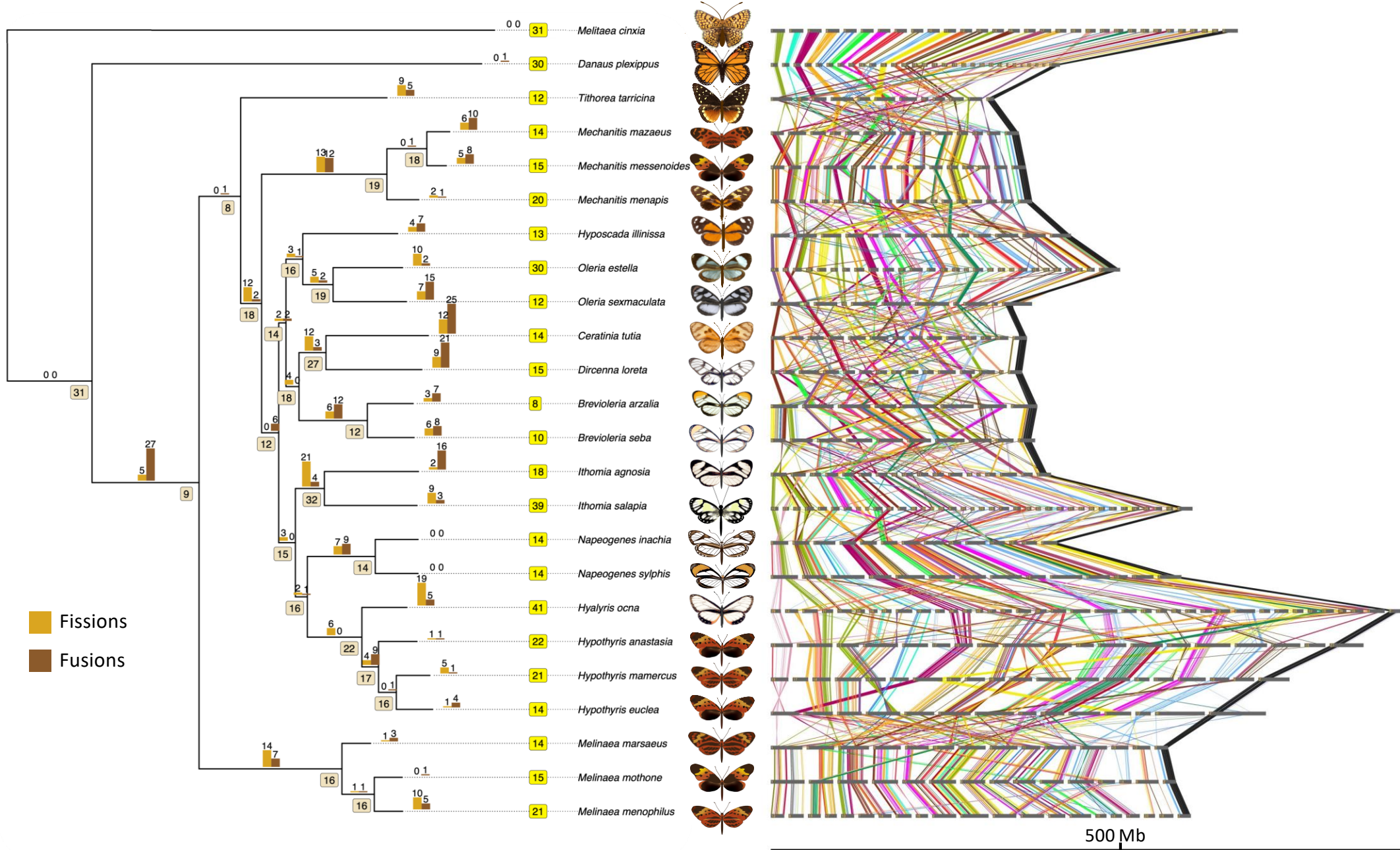


Karin Näsvali

Large-scale chromosomal rearrangements

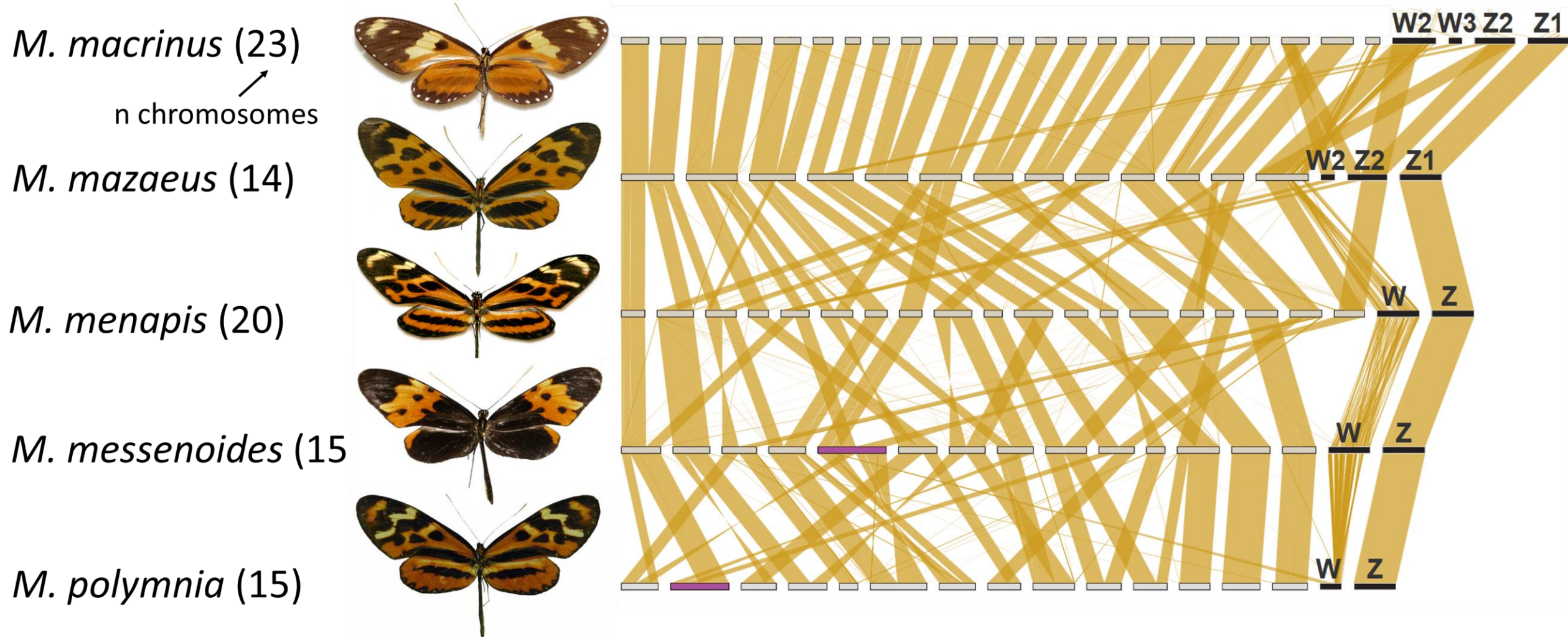


Large-scale chromosomal rearrangements

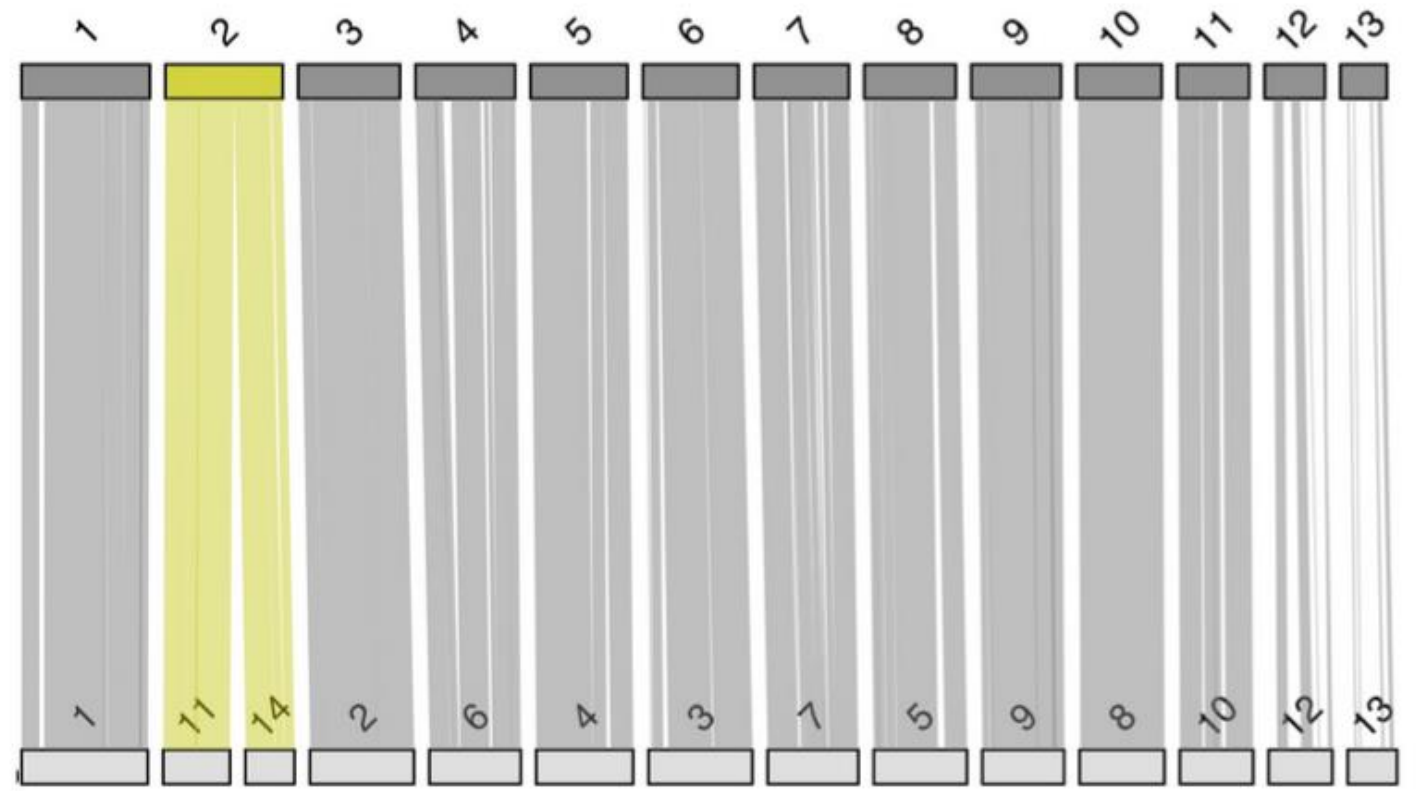


Karin Näsvall

Species of recent radiations (<1 Mya) show highly rearranged chromosomes

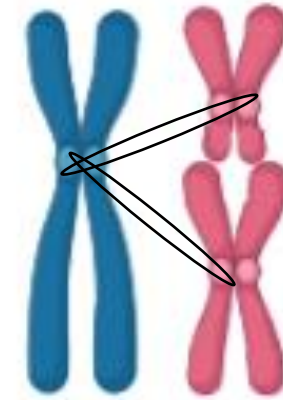
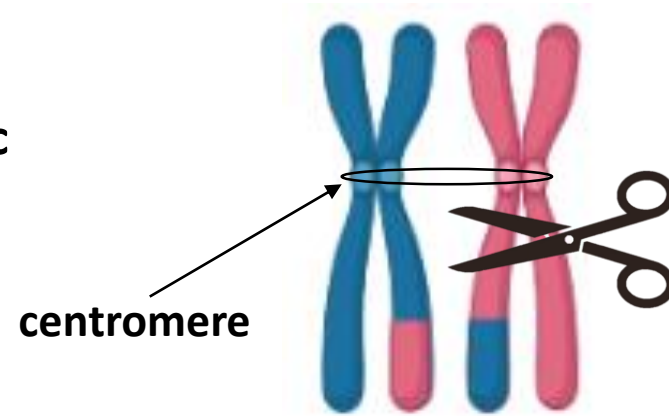


Some individuals are heterozygous for fission-fusion rearrangements

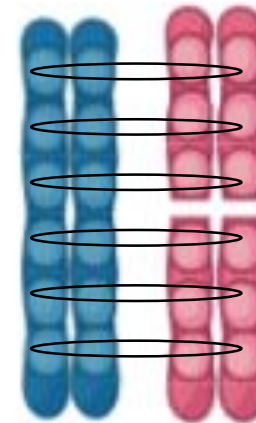
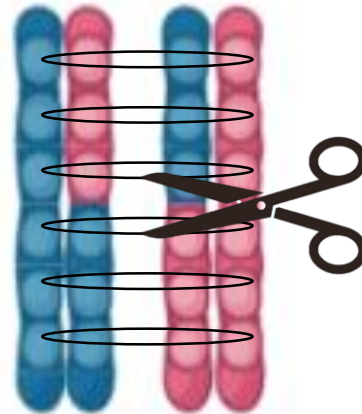


Butterflies have holocentric chromosomes

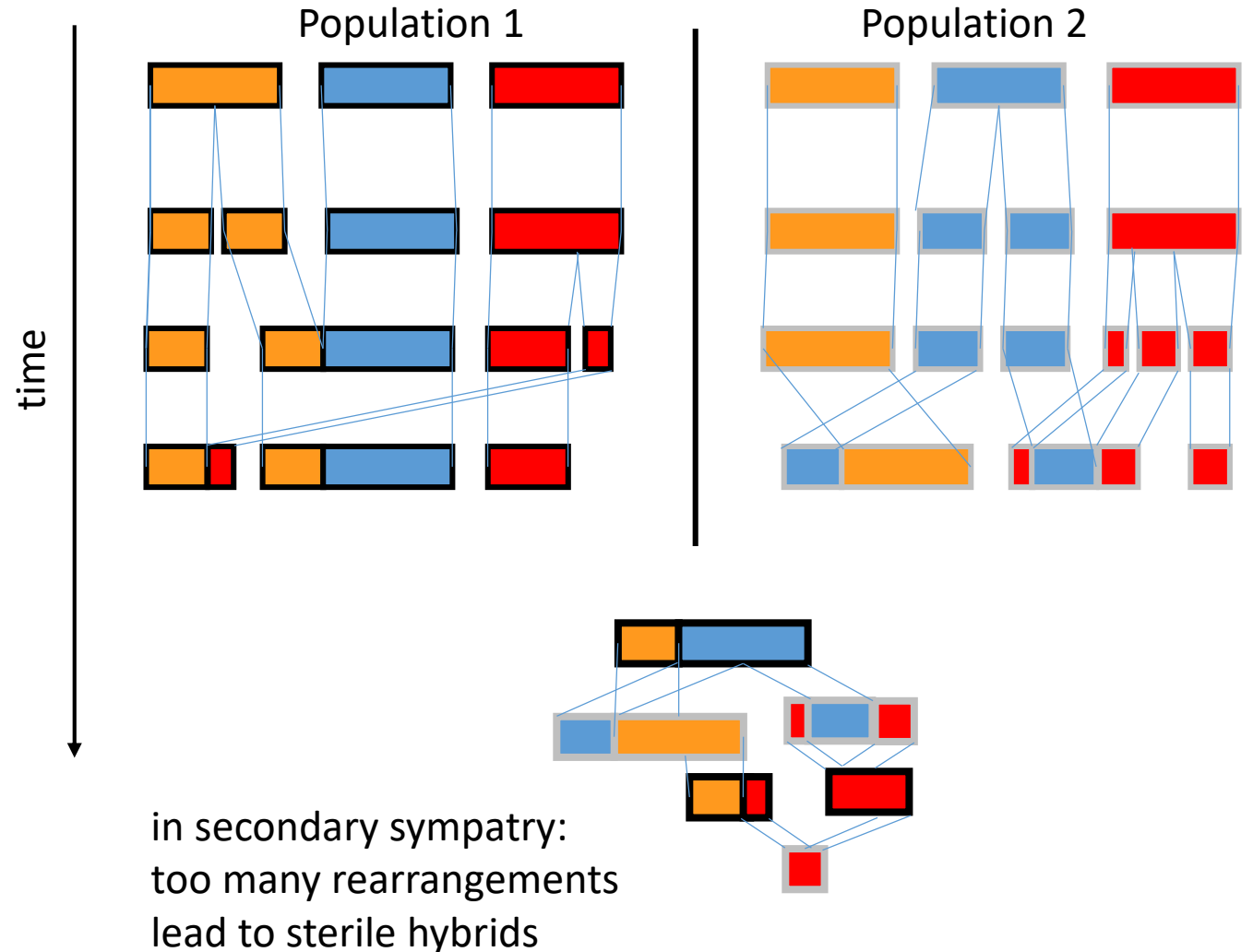
Monocentric
(like us)



Holocentric
(all butterflies)



Hypothesis: Chromosomal rearrangements during periods of geographic isolation drive speciation







Many thanks to my wonderful team and key collaborator

@Sanger



Arif Maulana
PhD Student



Dr Karin Näsvall
Postdoctoral Fellow



Dr Nicol Rueda
Postdoctoral Fellow



Dr Patricio A.
Salazar-Carrión
Postdoctoral Fellow



Jonah Walker
PhD Student



Eva van der
Heijden
PhD student

@IKIAM
en Ecuador



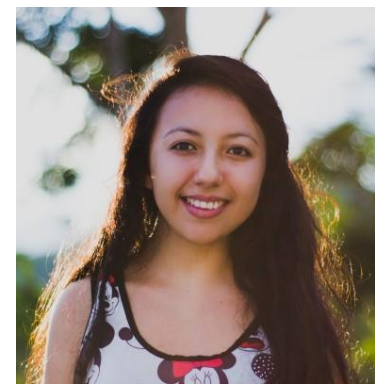
Franz Chandi



Kimberly
Gavilanes



Alex Arias



María José
Sánchez



Prof Caroline Bacquet