

Bioinformatics for Next-Generation Sequencing
ICIPE, Nairobi, November 2014

Evolutionary Genomics 1: Comparative and Phylo-genomics

Simon Martin

Department of Zoology
University of Cambridge



Evolutionary Genomics

- Lecture 1 (this one): Comparative and phylogenomics
 - Genome organisation and synteny
 - Gene family evolution
 - Phylogenomics
- Lecture 2: Population genomics
 - Diversity and population structure
 - Adaptive and neutral evolution
 - Demographic history and hybridisation
 - Speciation

Genome organisation and synteny

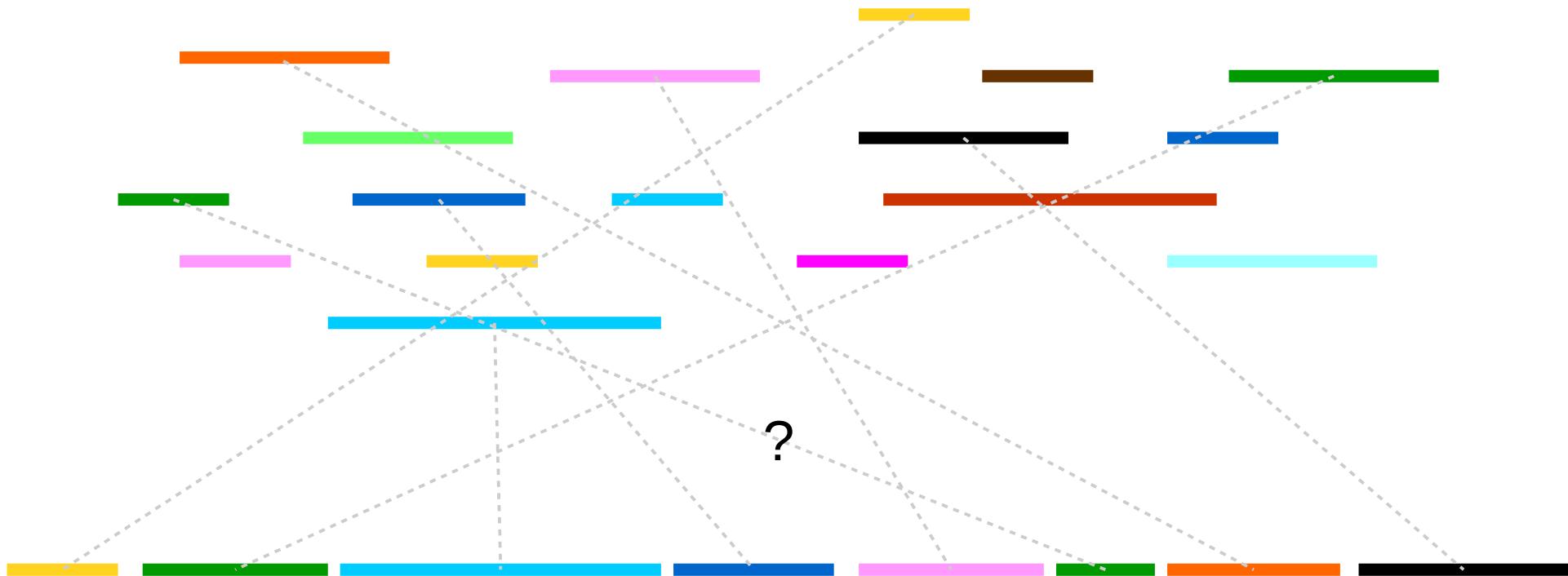
- How has the size and organisation of the genome changed between species?
- Macrosynteny
 - Chromosomal fusions
 - Splits
 - segmental duplications
- Microsynteny
 - Gene order
 - indels
 - inversions

Macrosynteny

- Requires “complete” genome assemblies
- Linkage map (or chromosome-sized scaffolds)
- Annotated genes for both species
- Homologous gene pairs

From contigs to chromosomes

- Next-gen approaches usually produce fairly small contigs.
 - Usually thousands to tens/hundreds of thousands
- So how do we assemble chromosomes?

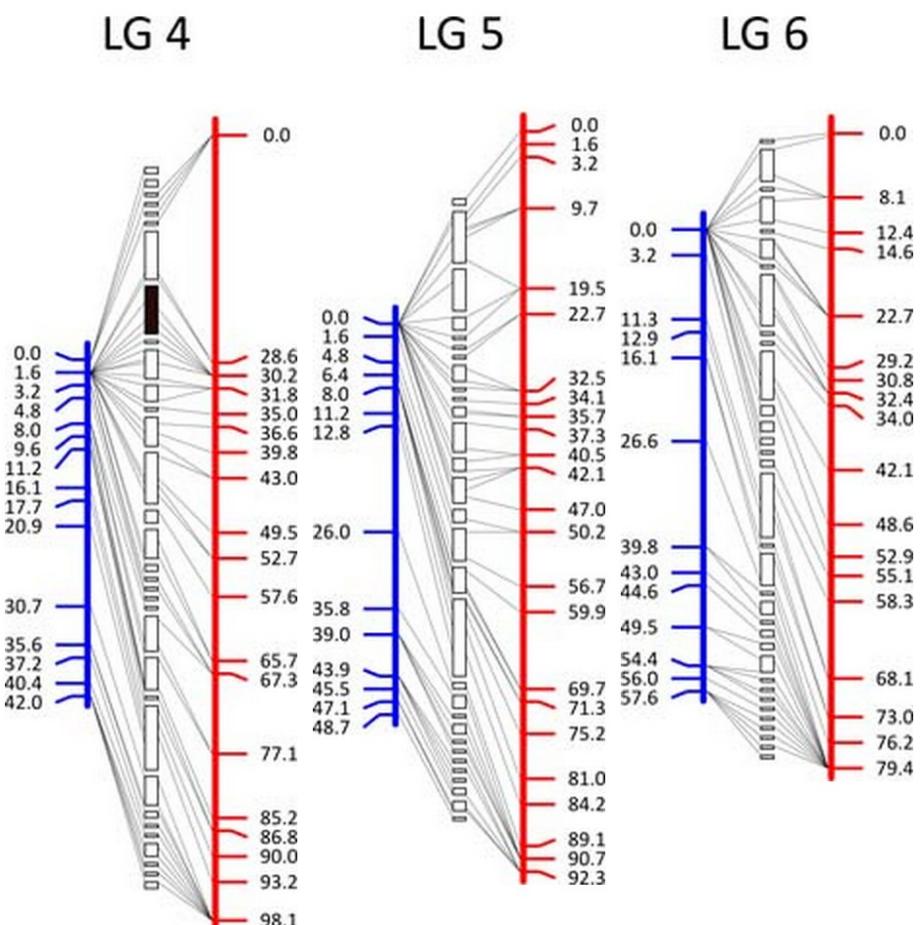


From scaffolds to chromosomes

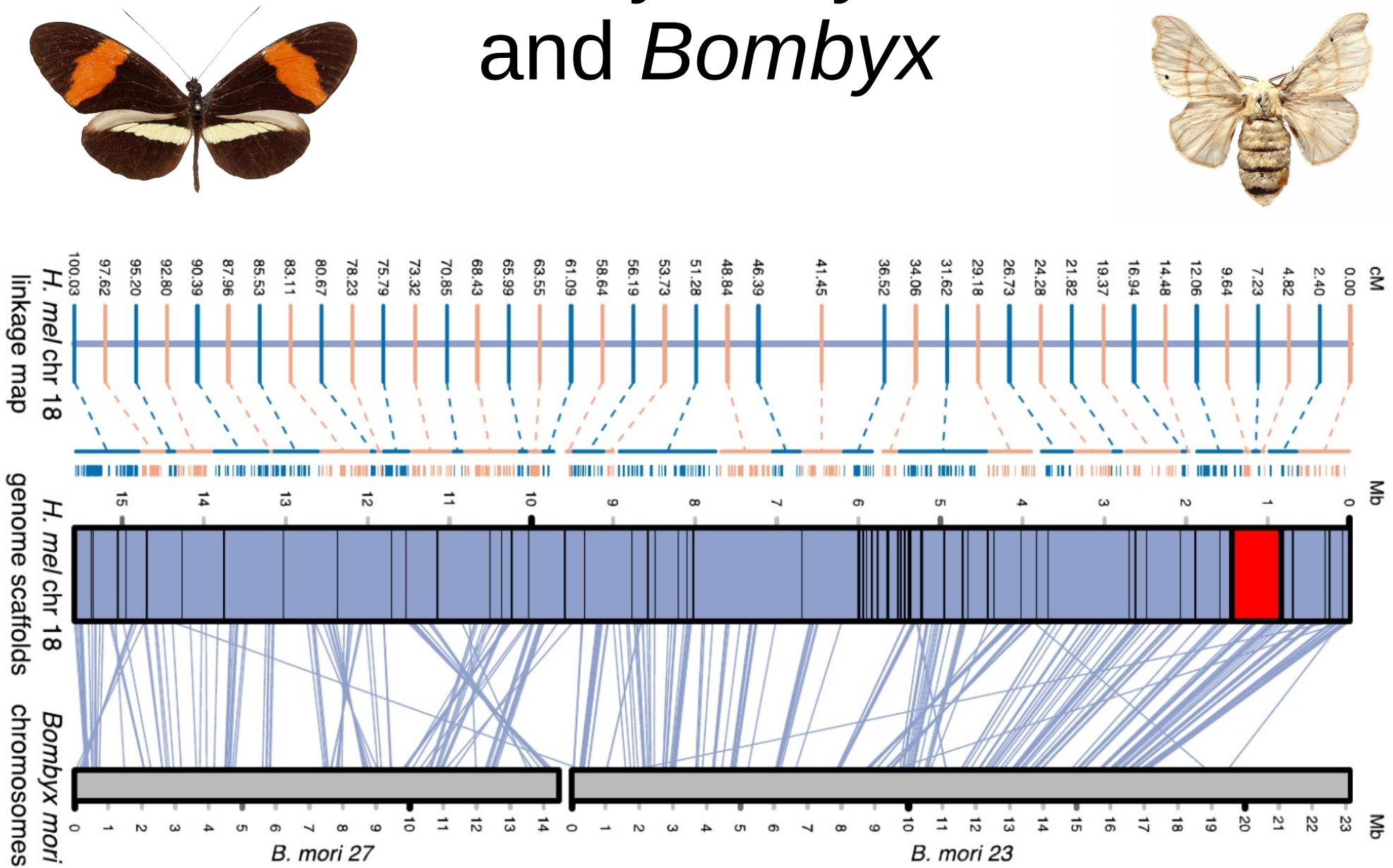
- Connect **genetic linkage maps** to scaffold sequences



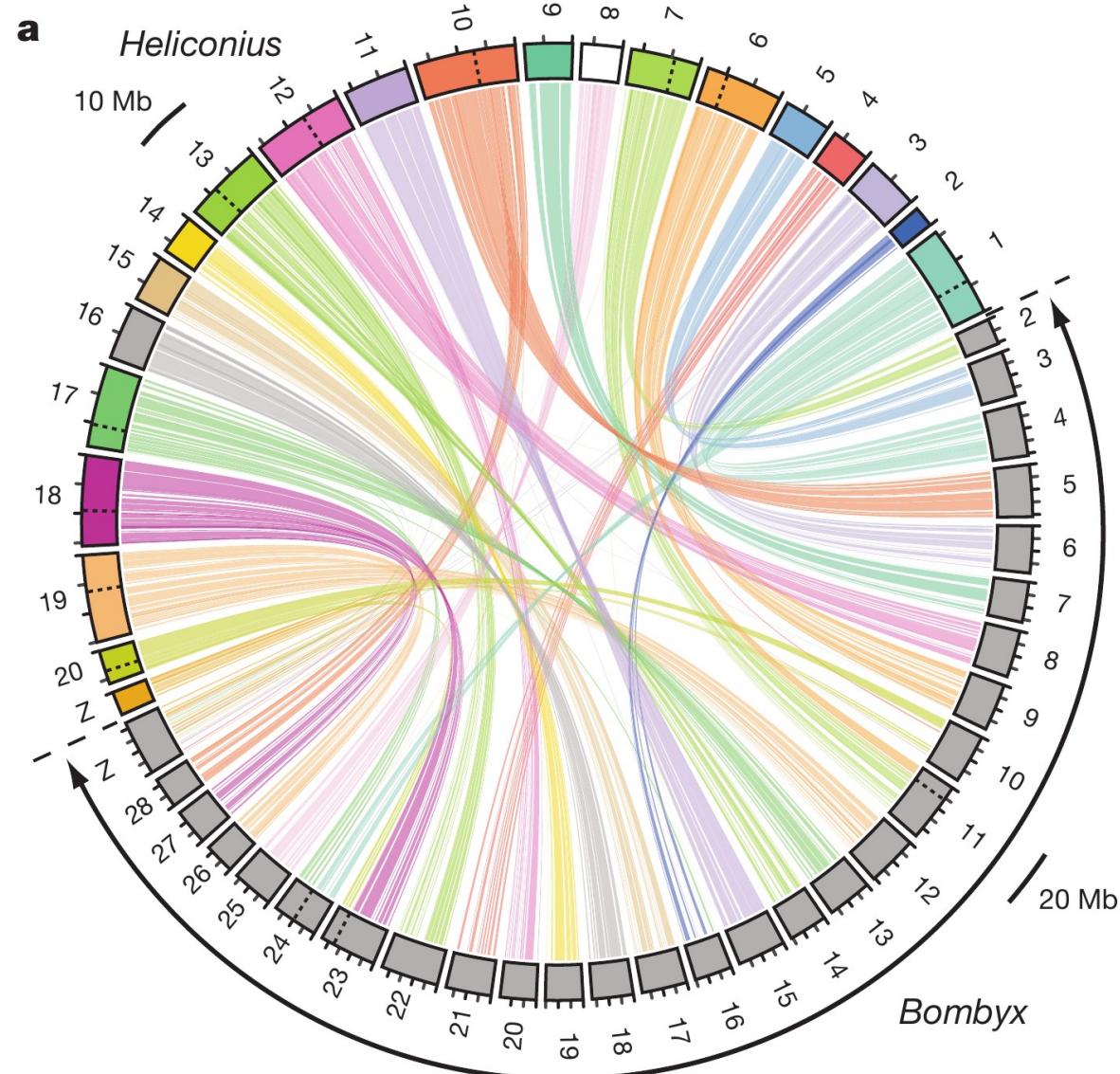
Marker	Scaffold Number in Fugu ver5	Map position (cM)		Primer sequences	
		Male Map	Female Map	Forward (5'-3')	Reverse (5'-3')
f799		357	0.0	- ggtttggcacattctgact	ttagcattaggccgacgag
f756		357	0.0	- gtaaaggctggatggactg	ggctggaaaaactctgact
f1414		588	-	0.0 catatggcctggaaacactc	gagggcgaggaagagattt
f1424		599	0.0	0.0 cgaaaggccacagacagatt	ttgcccgaataatctggta
f1468		655	-	0.0 cacgtctgcgttgtgtata	ctcgtataggctggtgat
f1486		676	0.0	0.0 aaaagcaccagatgtgacc	acttggccctggggactat
f209		52	1.6	28.6 ccccttttcctgtcatcatt	tctctgtgggtcatctgttt
f338		52	1.6	28.6 ggtgtggagccctcagagat	ttgcatgaattggttctgt
f560		53	1.6	- cggccctatgagtattgaa	gggagcagaaggtaatgtgaa
f351		53	-	30.2 gctgactgaaaactgcacca	atgaaaccatccaacccat
f551		53	-	30.2 tacttcatgcggcacttgg	gagtcgaagattgttcaca
f91		53	1.6	30.2 tgcatgacaaaaatgaaaggac	tcattttccggcttgttt
f1217		53	1.6	30.2 atctggcctgttgttcacc	gagctcagaaaaaccagtgtgg
f1594		822	1.6	30.2 aggaatgtggctcaagtgg	gccagtgtgcataaaaaaca
f1747		121	1.6	30.2 ccacctgtacaaacacacacag	agtatttgcttgttacccc
f76		121	1.6	- ggcccactttcacatct	tctaaaactggccagcagaa
f1746		121	1.6	30.2 agactgagggtgtatcggtt	cgtgtgcacaaagggtgtat
f32		121	1.6	- ttatcagctctcaaact	agggtcggtgtgggtaa
f1745		121	-	31.8 gaggcattcactggctttc	agggggaaaatttaggatgcac
f650		201	1.6	31.8 gccaacgcacttgaaaggt	aaaccagcagtttcccaga
f1322		201	1.6	31.8 ctcccaaaaaacccatct	ctttttccaaaaatgtatct



Chromosomal synteny in *Heliconius* and *Bombyx*



Chromosomal synteny in *Heliconius* and *Bombyx*

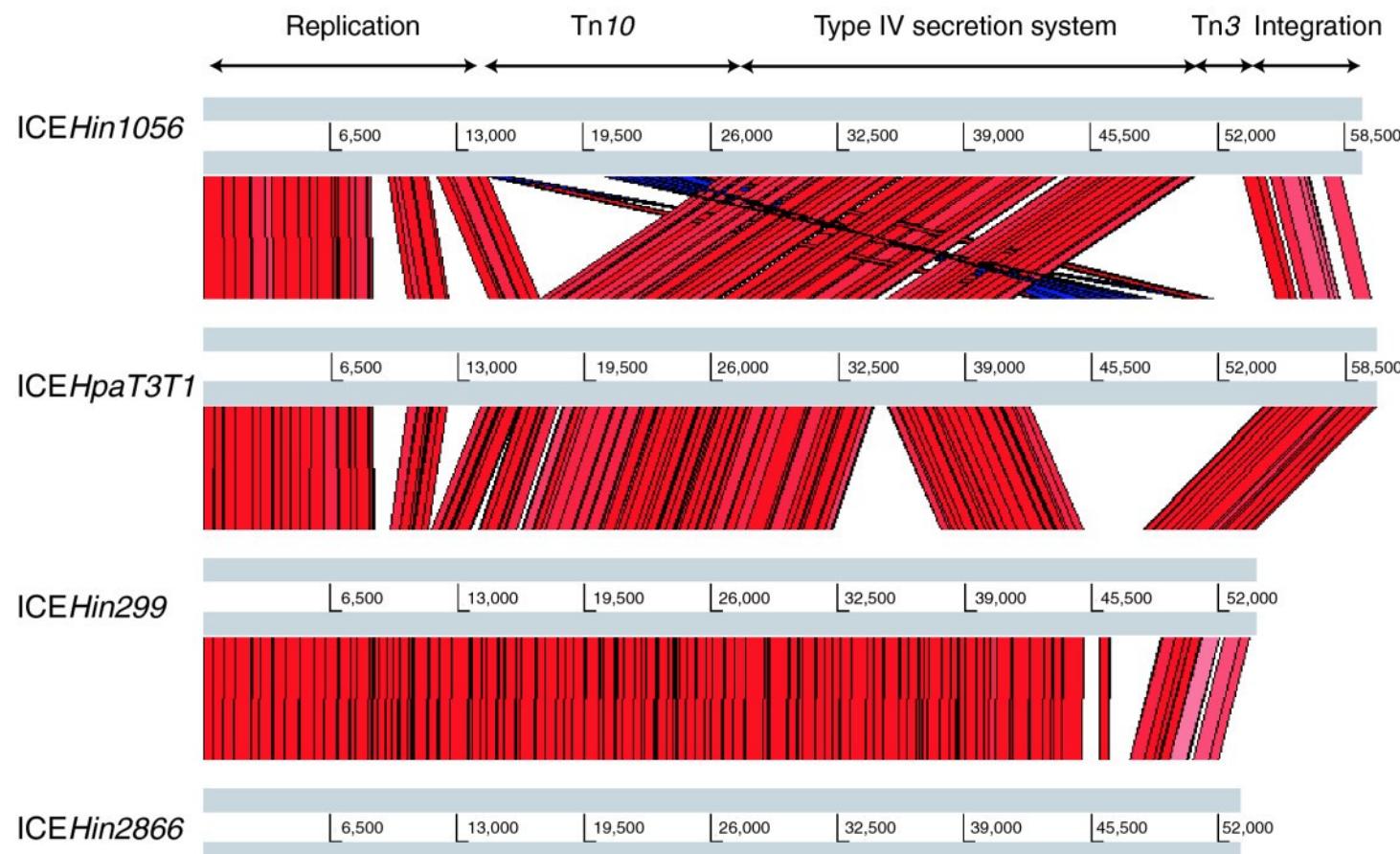


Some software for macrosynteny

- InParanoid – find orthologous gene pairs
- OrthoCluster – identify synteny blocks and chromosomal rearrangements
- CIRCOS – beautiful chromosome plots!

Microsynteny

ACT (Artemis Comparison tool)

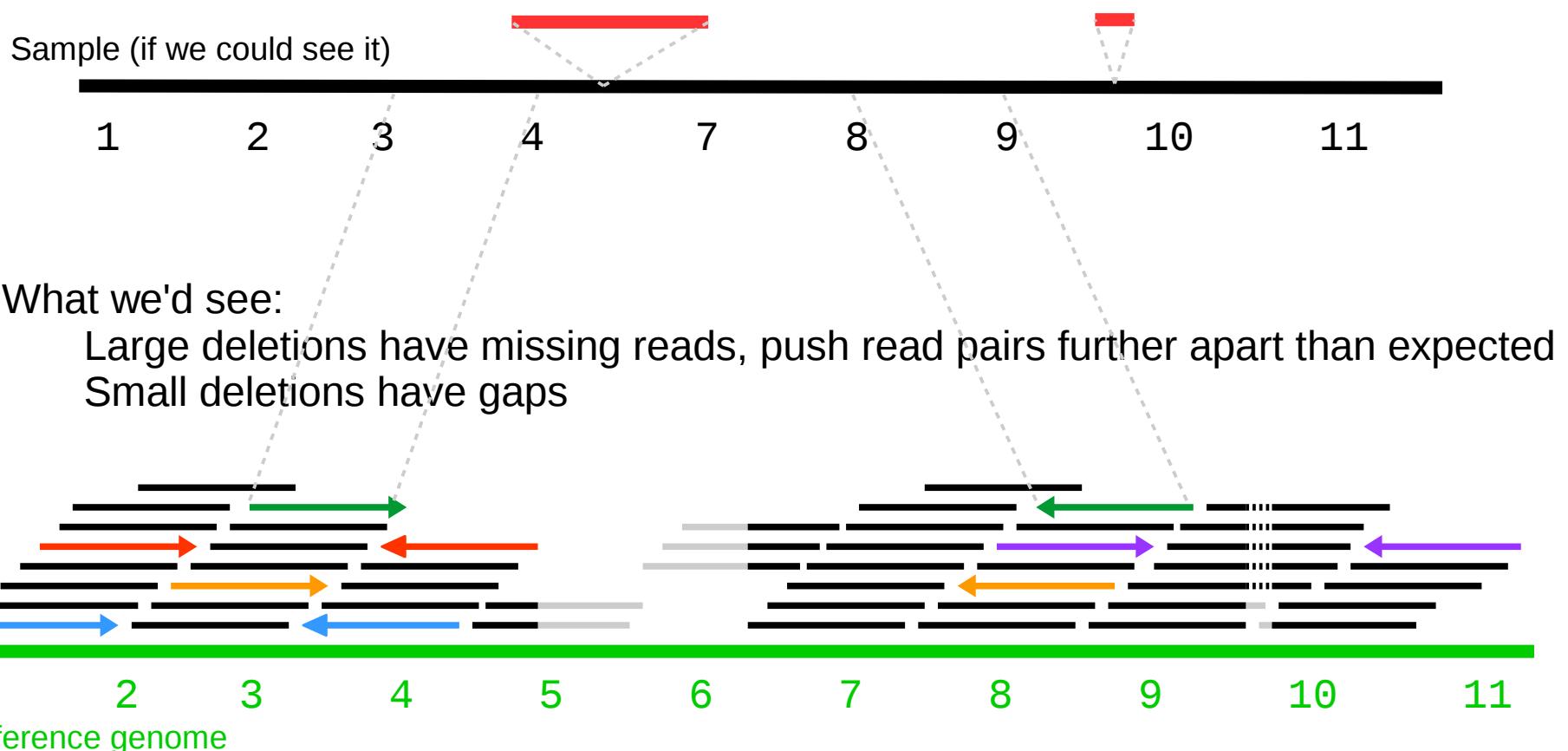


Can you study genome structure evolution with just one assembled genome?

- Yes, you can study “structural variants” (SVs) among individuals from the **same species or closely related species**.
- By resequencing single individuals and mapping the reads to the reference you can identify
 - Deletions
 - Insertions
 - Duplications
 - Inversions

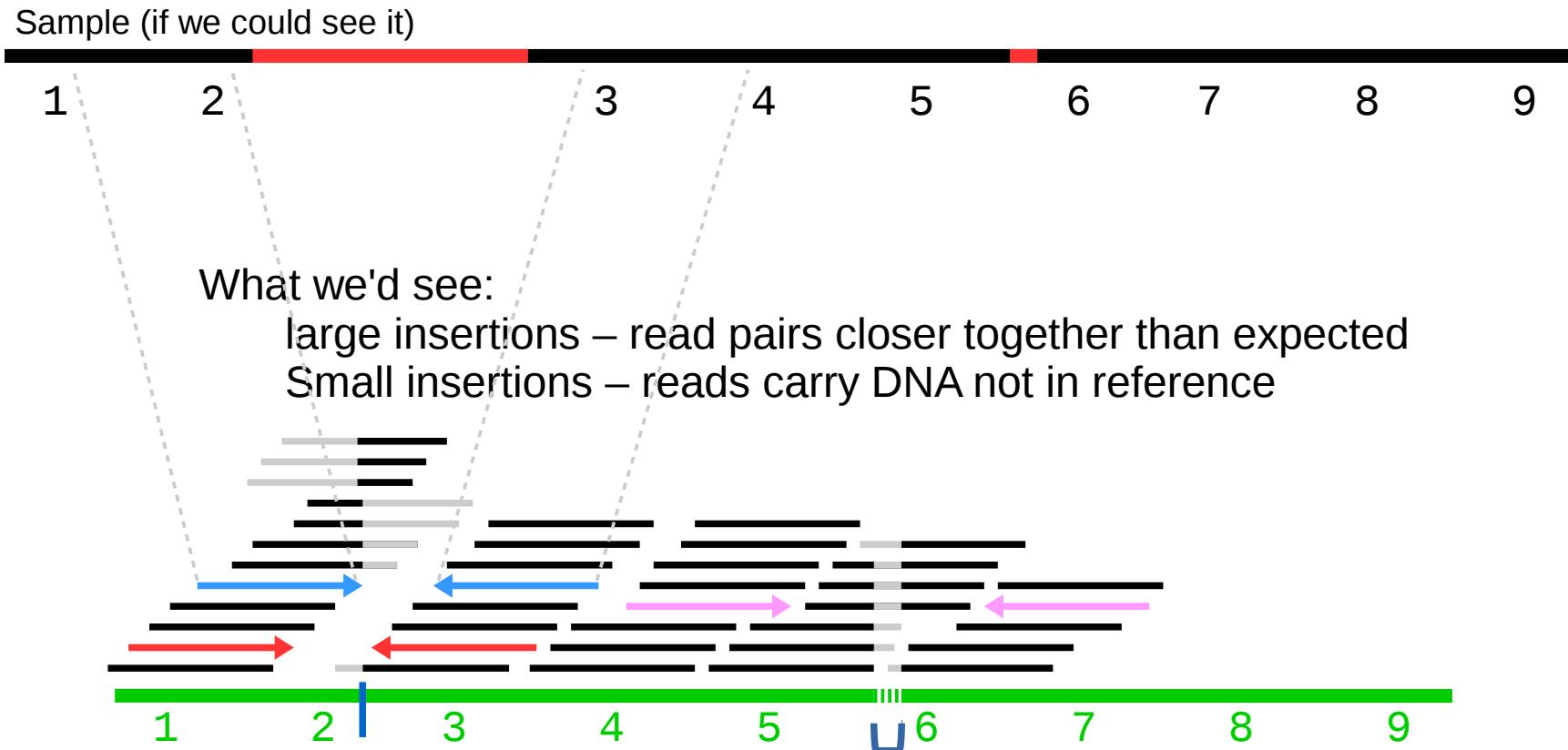
Using next gen data to study structural variation

Deletions



Using next gen data to study structural variation

Insertions



Using next gen data to study structural variation

Tandem duplications

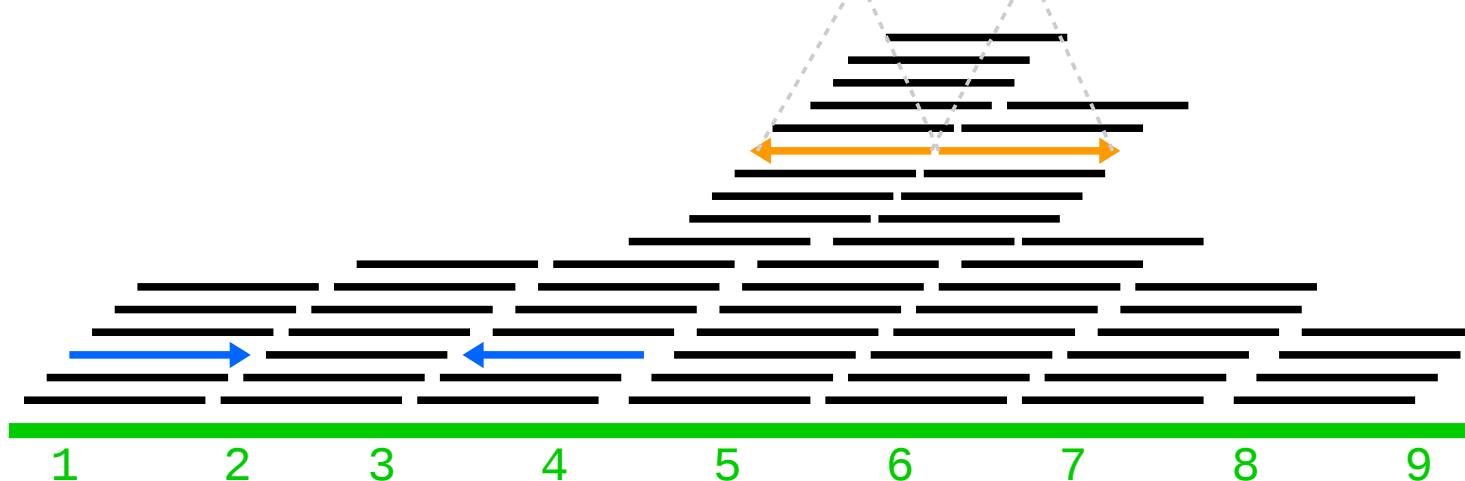
Sample (if we could see it)



What we'd see

Duplicated regions have increased read depth

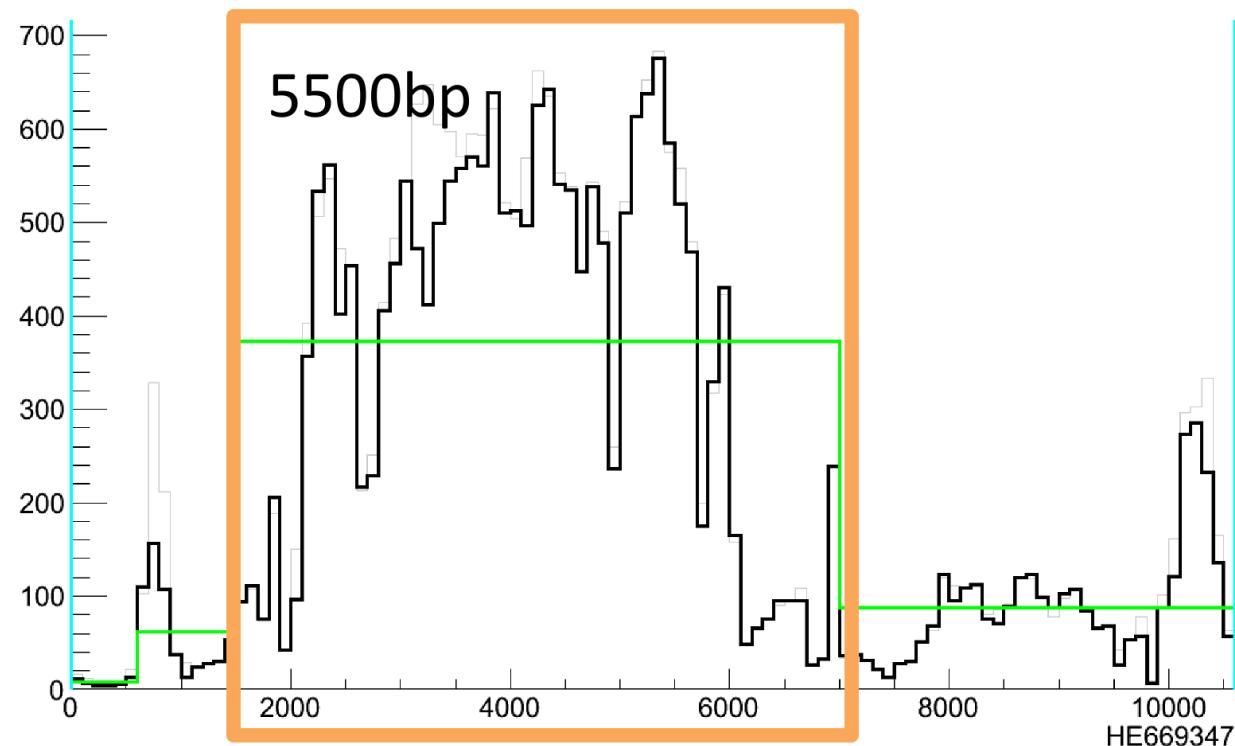
Read pairs have altered orientation



Using next gen data to study structural variation

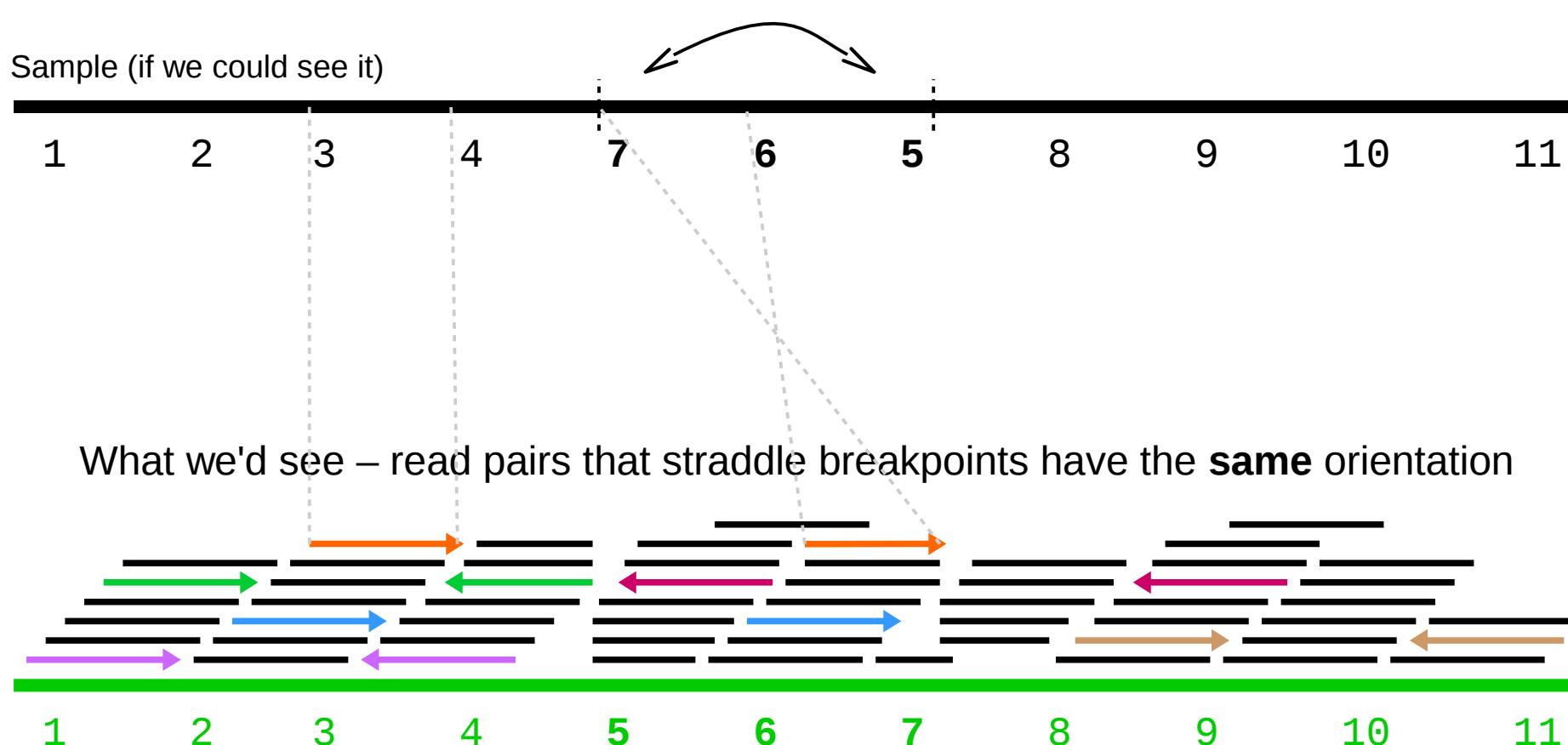
Duplications relative to reference

In reality, depth of coverage is quite variable, so only large duplications are detectable



Using next gen data to study structural variation

Inversions



Using next gen data to study structural variation

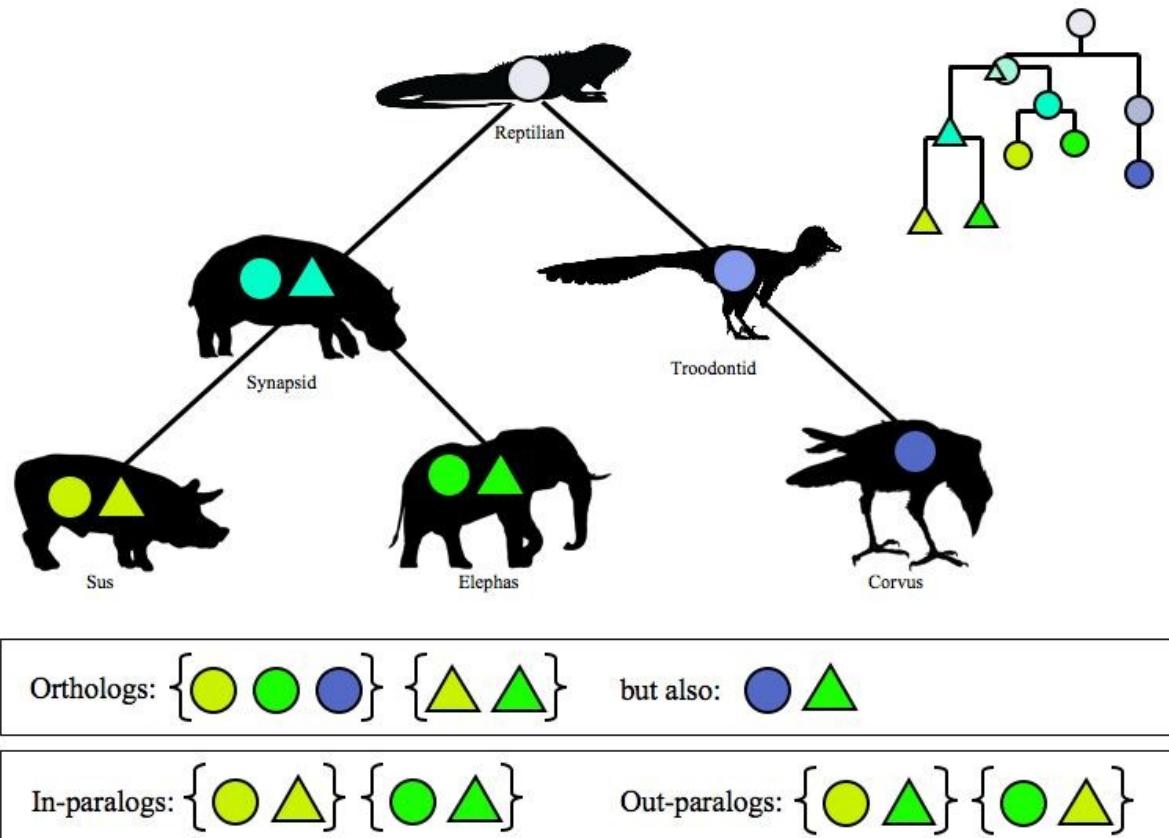
Some software

- **CNVnator** (Abyzov et al. 2012 Genome Research)
- **Pindel** (Ye et al. 2009 Bioinformatics)
- **Genome STRiP** (Handsaker et al. 2011 Nature Genetics)
- **DELLY** (Rausch et al. 2012 Bioinformatics)

Gene family evolution

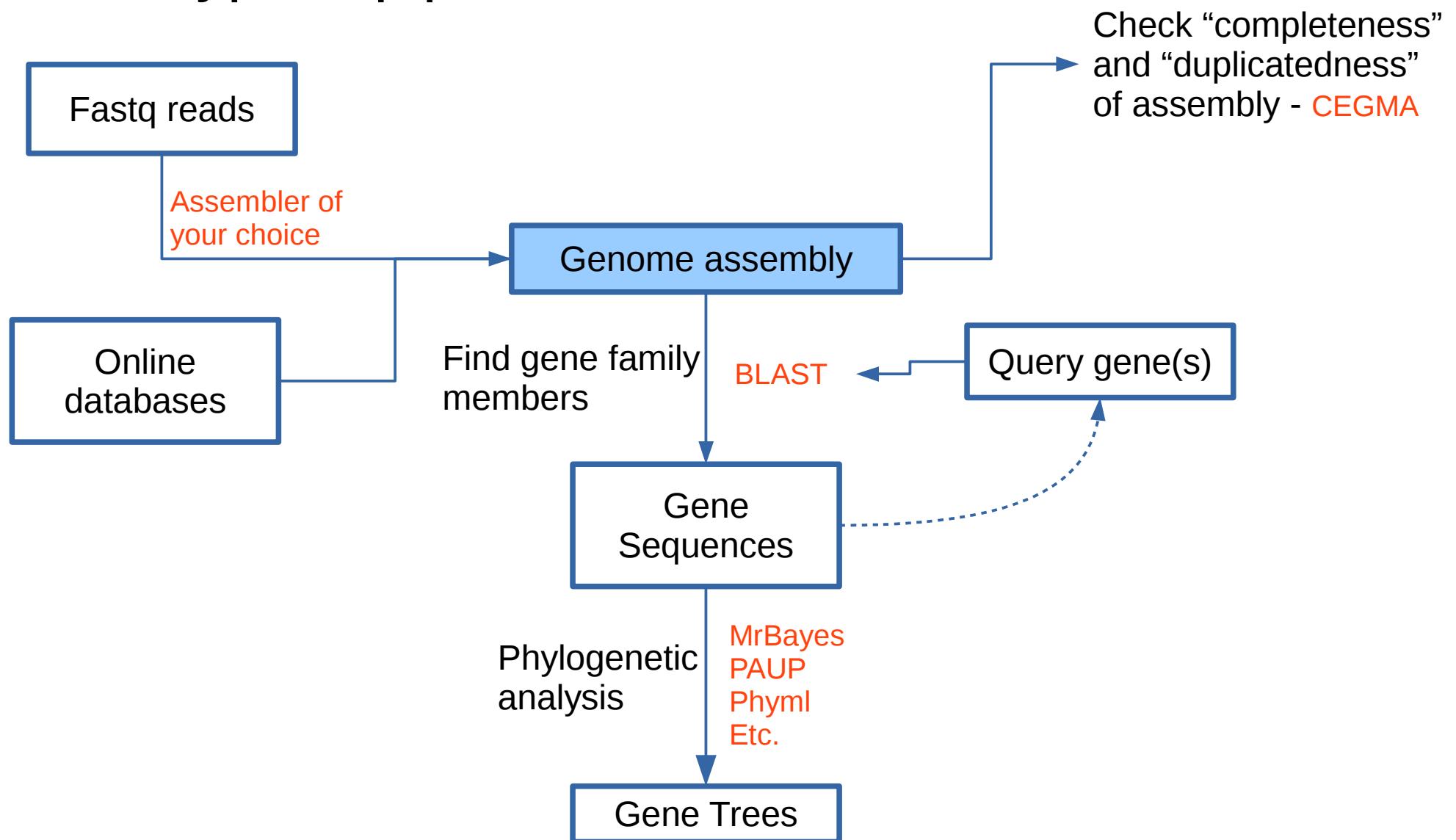
- Gene duplication leads to “families” of related genes
- E.G. There are ~80 Cytochrome P450 genes in the *Drosophila* genome
- We can study the relationship among homologous genes in one or more species

NOTE: **Orthologs** and **paralogs** are both types of **homologs**



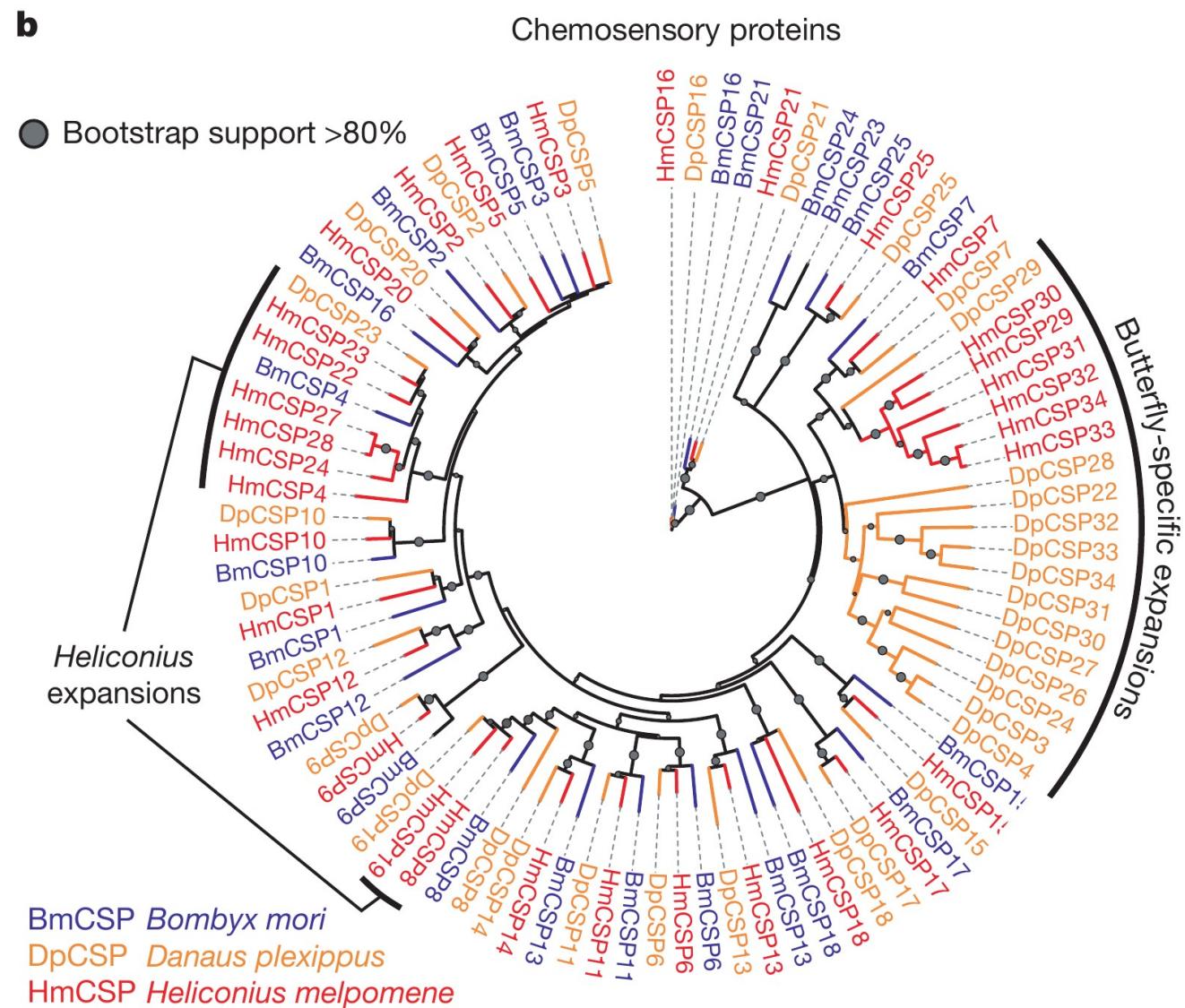
Gene family evolution

A typical pipeline

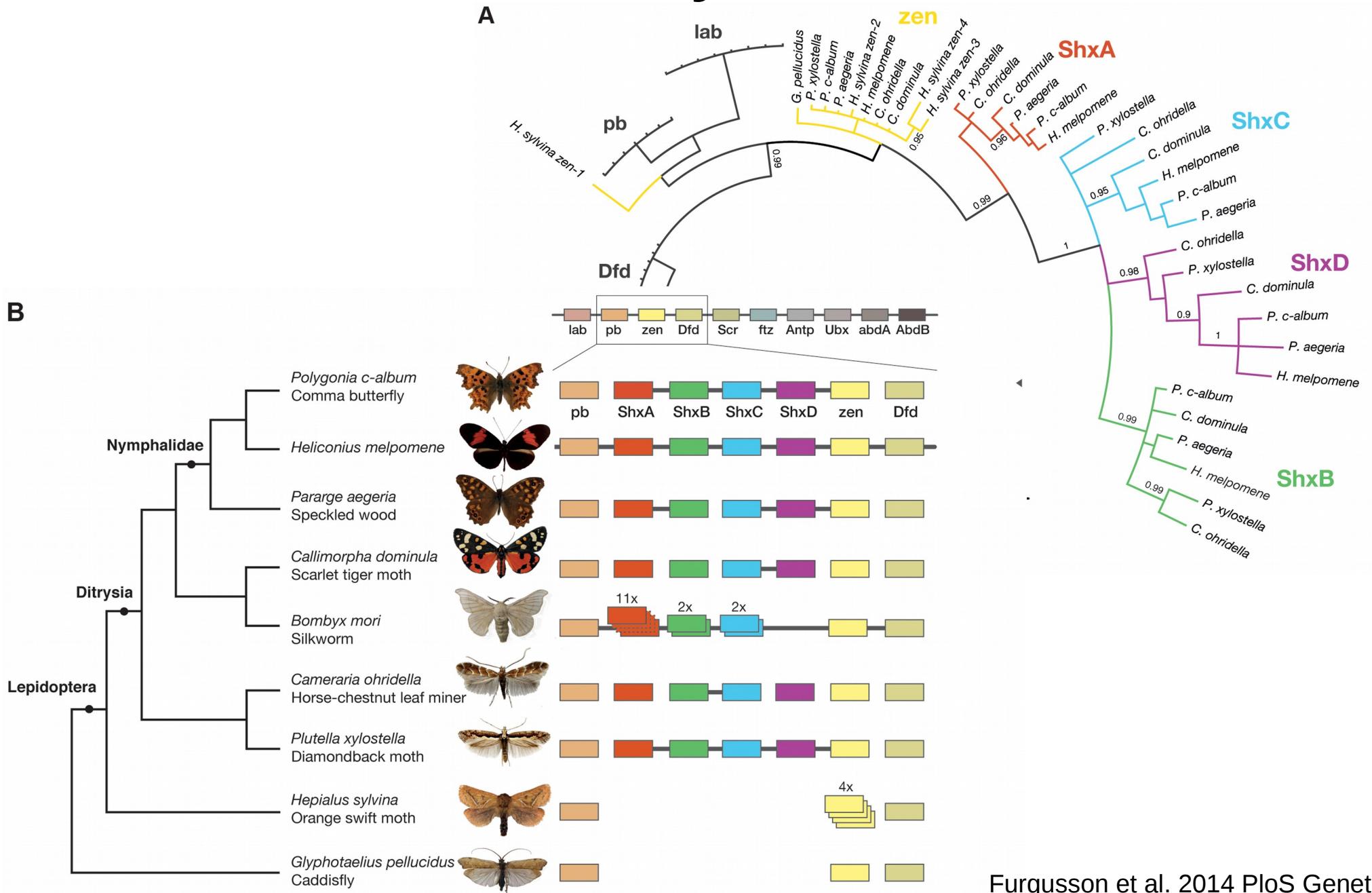


Gene family evolution

Chemosensory gene family expansions in butterflies.

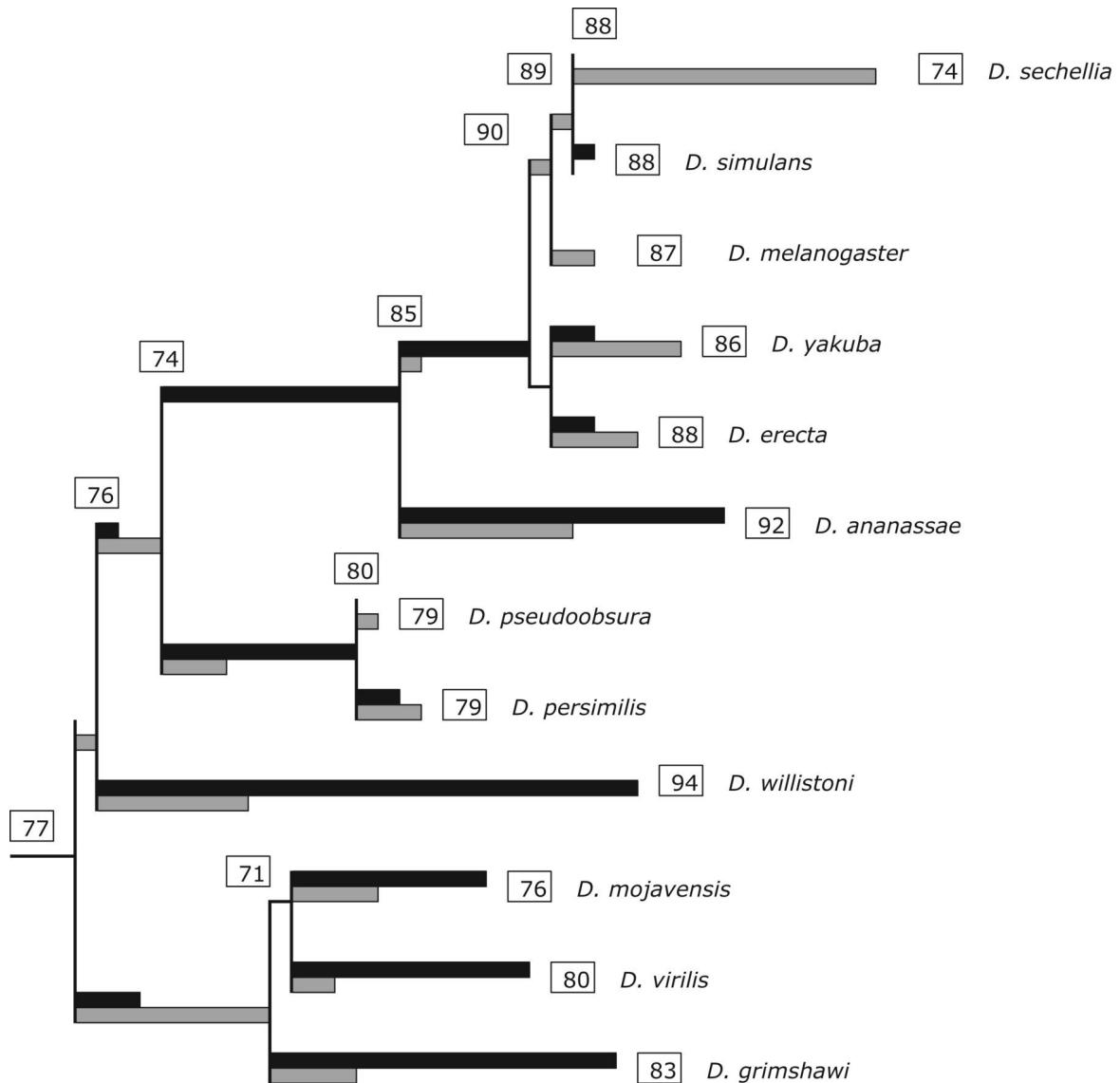


Gene family evolution



Gene family evolution

P450 gene gain and loss
across the *Drosophila*
phylogeny.

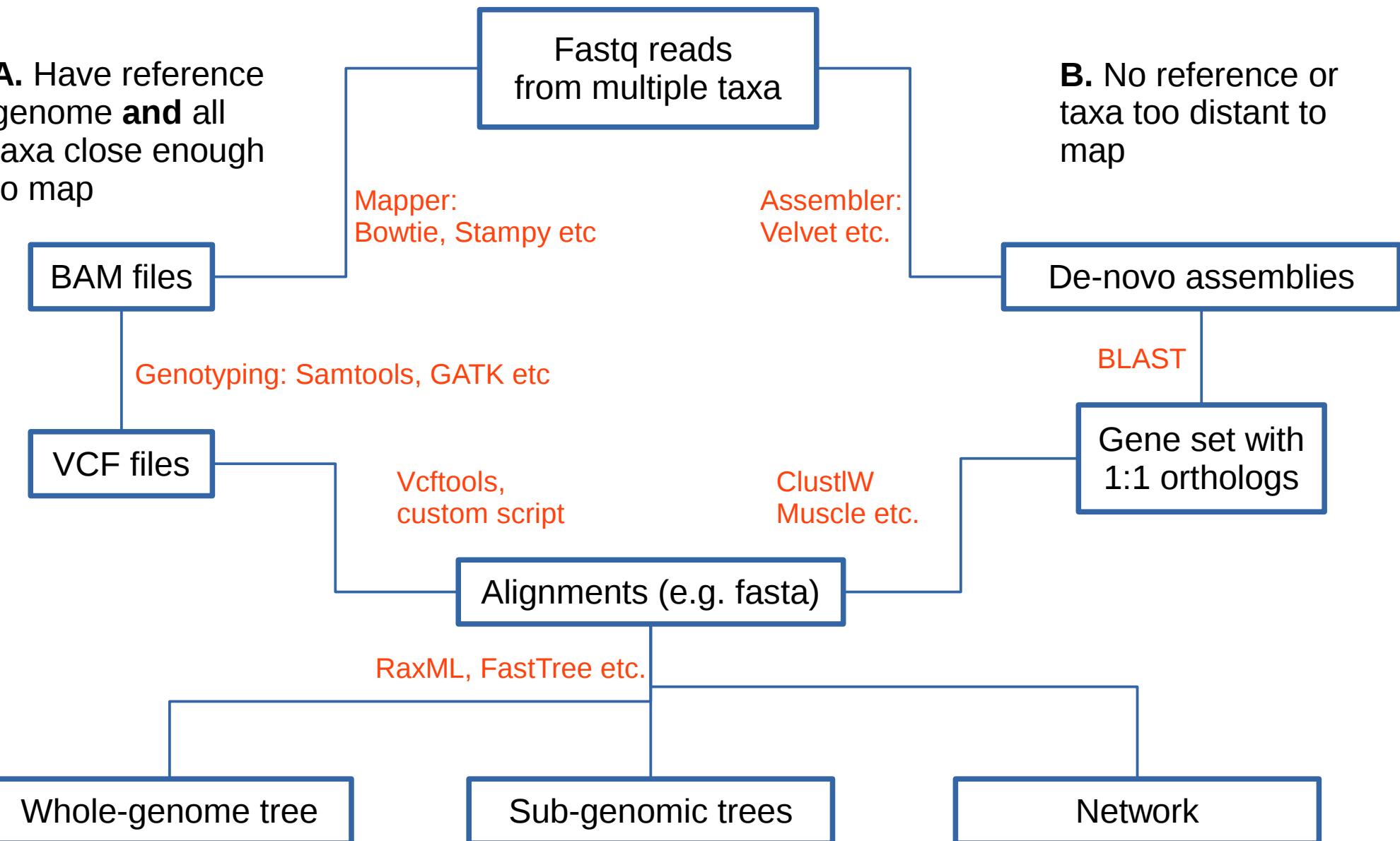


Phylogenomics

- Phylogenetics is a HUGE field – will not be covered in detail
- Instead we will focus on one misconception:
 - Genome-scale data *should* reveal the “true” relationships between species right?
 - Maybe not
- Relationships can differ in different parts of the genome
 - Lateral gene transfer and gene flow
 - Incomplete lineage sorting

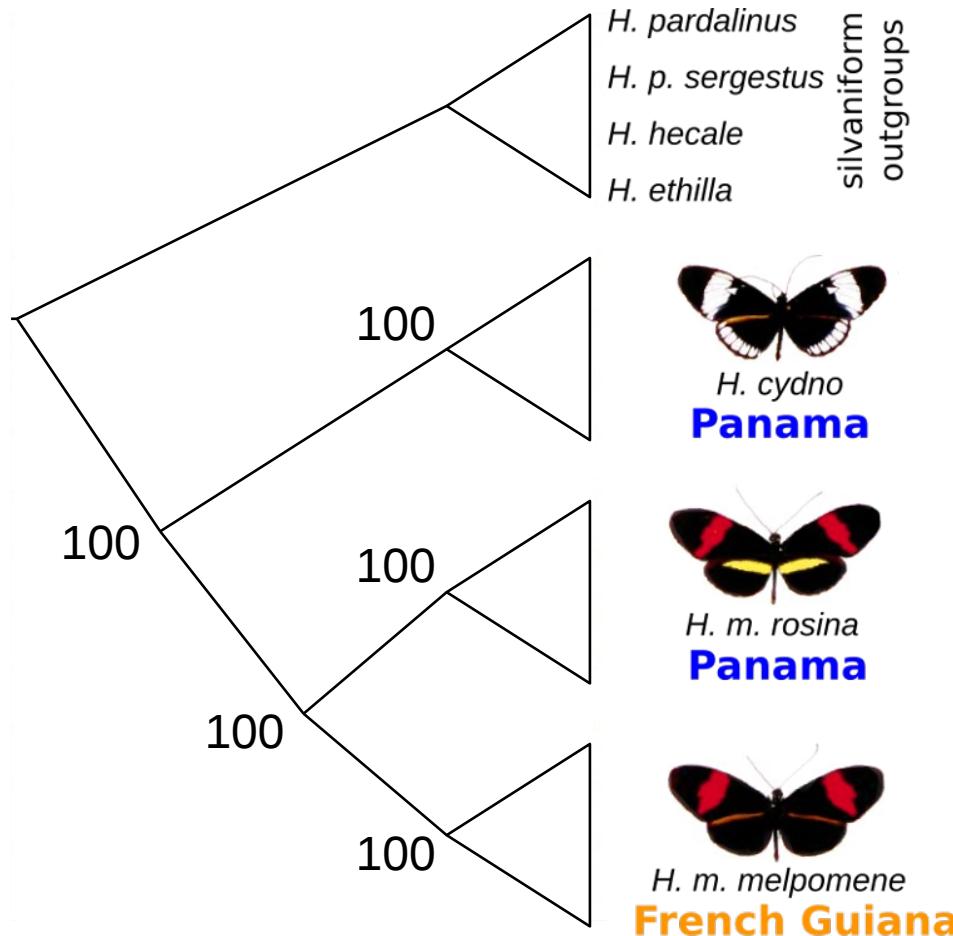
Phylogenomics with next-gen data - pipeline

A. Have reference genome **and** all taxa close enough to map



B. No reference or taxa too distant to map

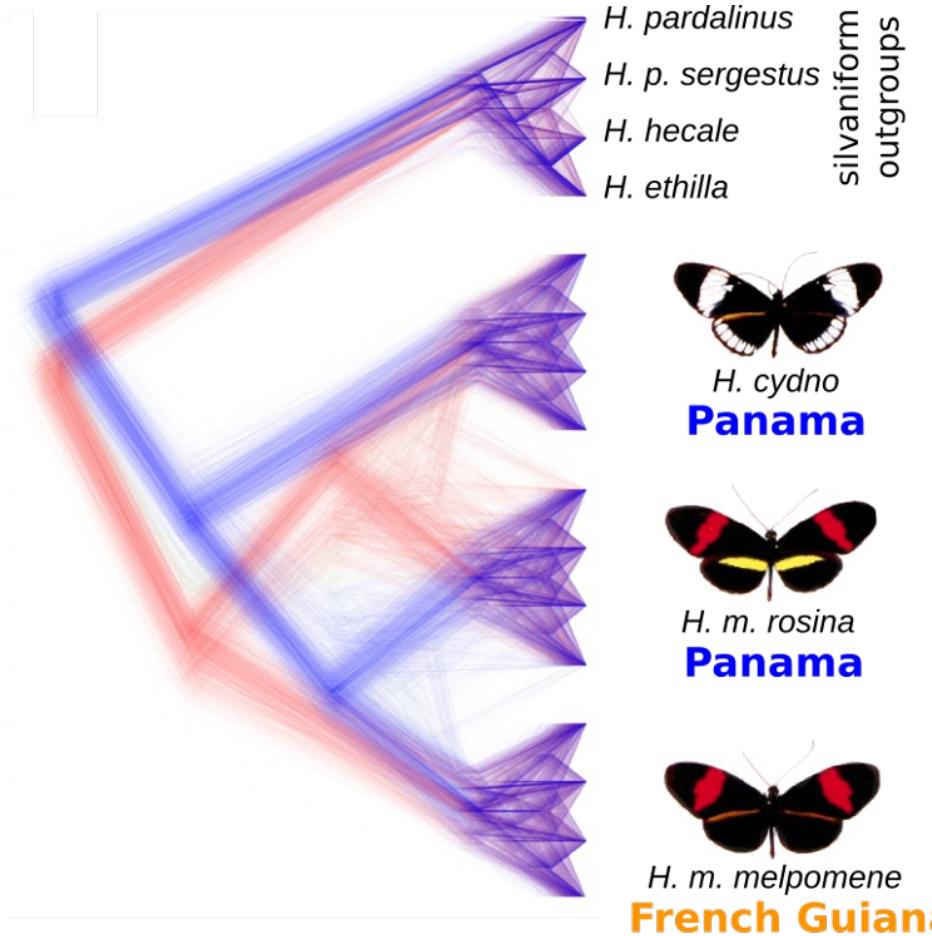
Phylogenomics: beware bootstrapping



Whole Genome Tree

100% bootstrap support for grouping by species.

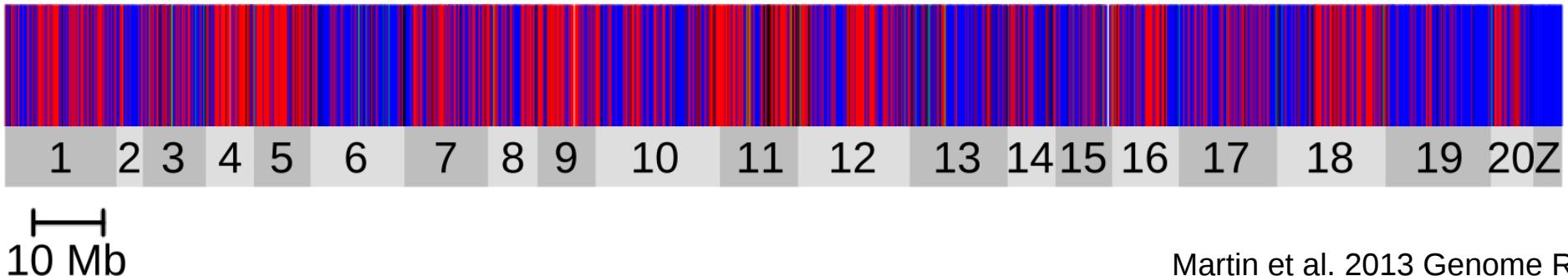
Phylogenomics: beware bootstrapping



Separate trees for 100 kb blocks

Only **53%** group populations by species. **42%** group populations by geography!

Bootstrapping over millions of sites will only ever support the most common tree.



Phylogenomics: YESTERDAY!

Published Online November 27 2014

Science DOI: 10.1126/science.1258524

< Science Express Index

 Leave a comment (0)

RESEARCH ARTICLE

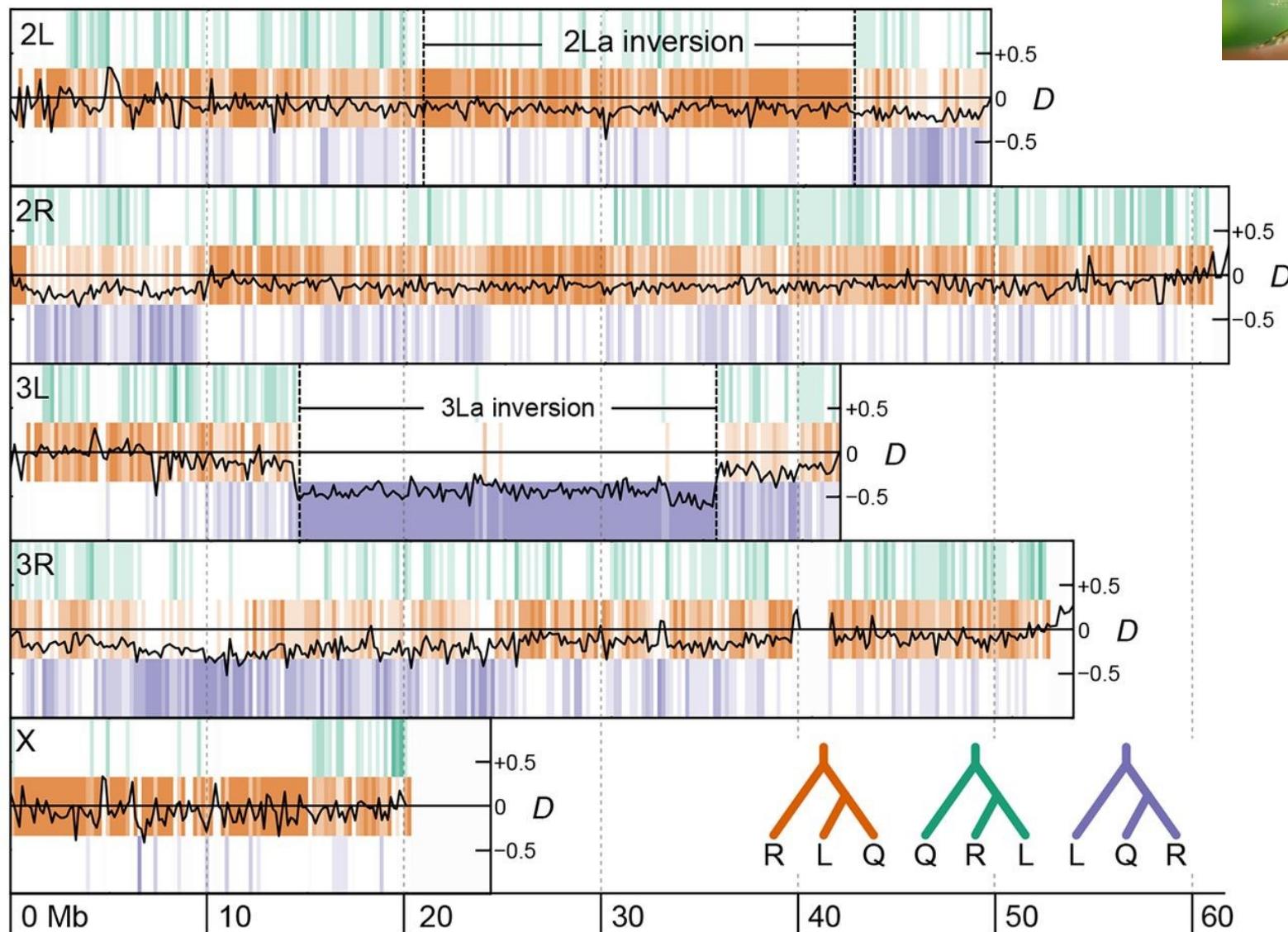
Extensive introgression in a malaria vector species complex revealed by phylogenomics

Michael C. Fontaine^{1,2,*†}, James B. Pease^{3,*}, Aaron Steele⁴, Robert M. Waterhouse^{5,6,7,8}, Daniel E. Neafsey⁶, Igor V. Sharakhov^{9,10}, Xiaofang Jiang¹⁰, Andrew B. Hall¹⁰, Flaminia Catteruccia^{11,12}, Evdoxia Kakanis^{11,12}, Sara N. Mitchell¹¹, Yi-Chieh Wu⁵, Hilary A. Smith^{1,2}, R. Rebecca Love^{1,2}, Mara K. Lawniczak^{13,‡}, Michel A. Slotman¹⁴, Scott J. Emrich^{2,4}, Matthew W. Hahn^{3,15,§}, Nora J. Besansky^{1,2,§}

Phylogenomics: YESTERDAY!



wikipedia.org



Phylogenomics: networks are more honest

