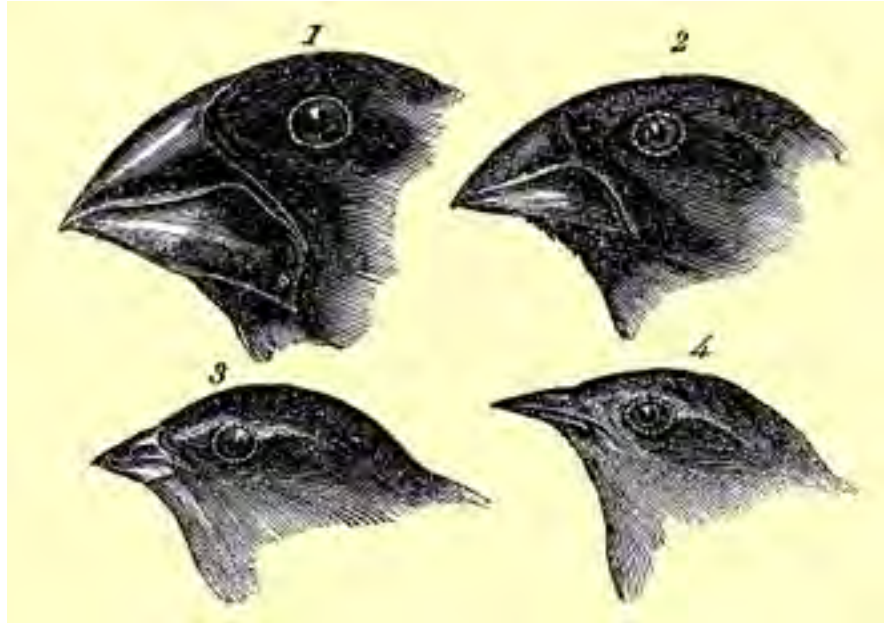


# Genomic studies of speciation and gene flow



# Why study speciation genomics?

# Why study speciation genomics?

Long-standing questions (role of geography/gene flow)

# Why study speciation genomics?

Long-standing questions (role of geography/gene flow)

How do genomes diverge?

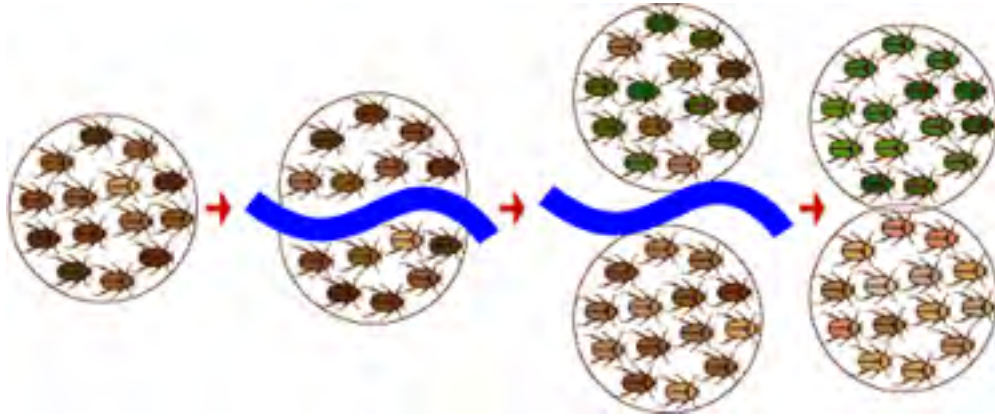
# Why study speciation genomics?

Long-standing questions (role of geography/gene flow)

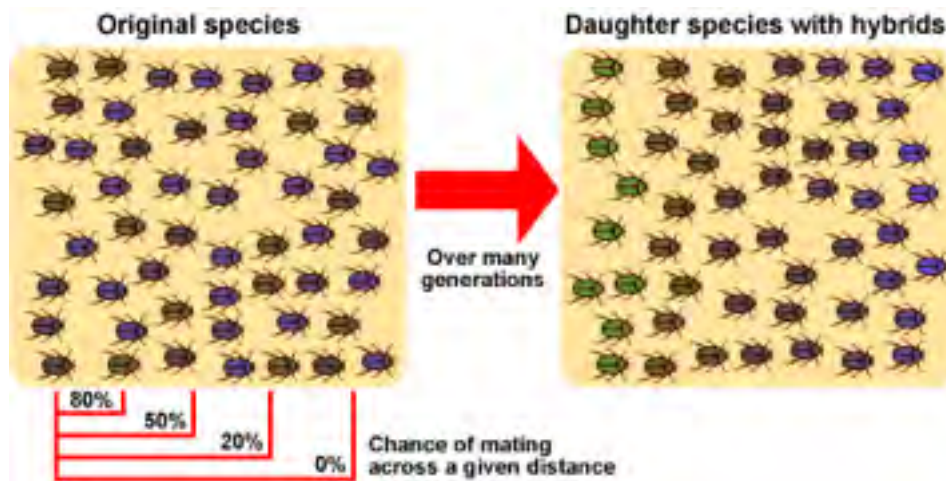
How do genomes diverge?

Find speciation genes

# Genomic divergence during speciation

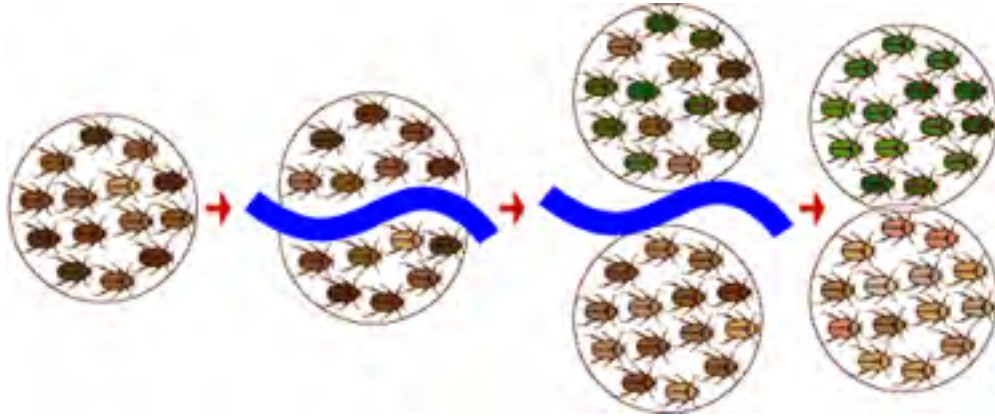


1. Speciation as a bi-product of physical isolation

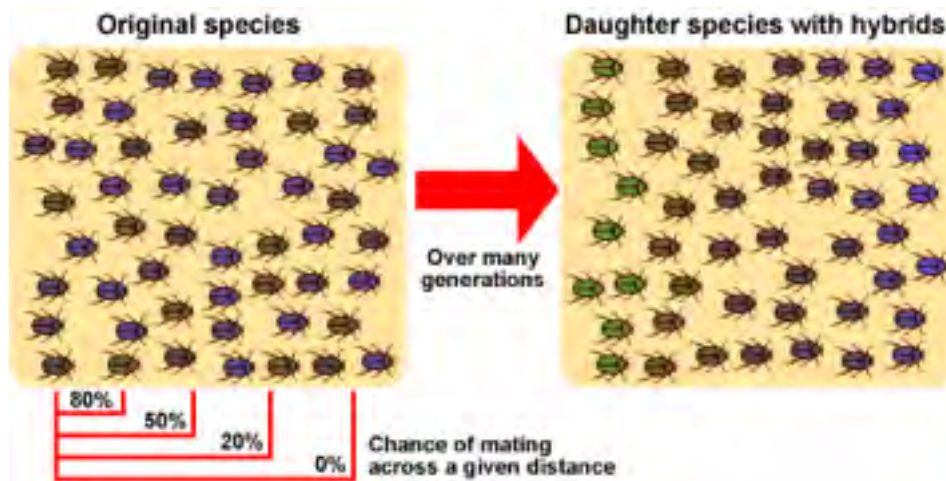


2. Speciation due to selection – without isolation

# Genomic divergence during speciation



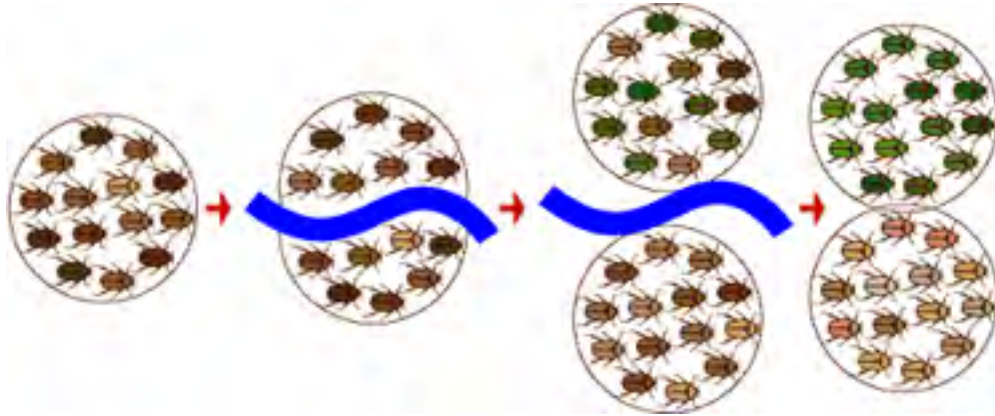
1. Speciation as a bi-product of physical isolation



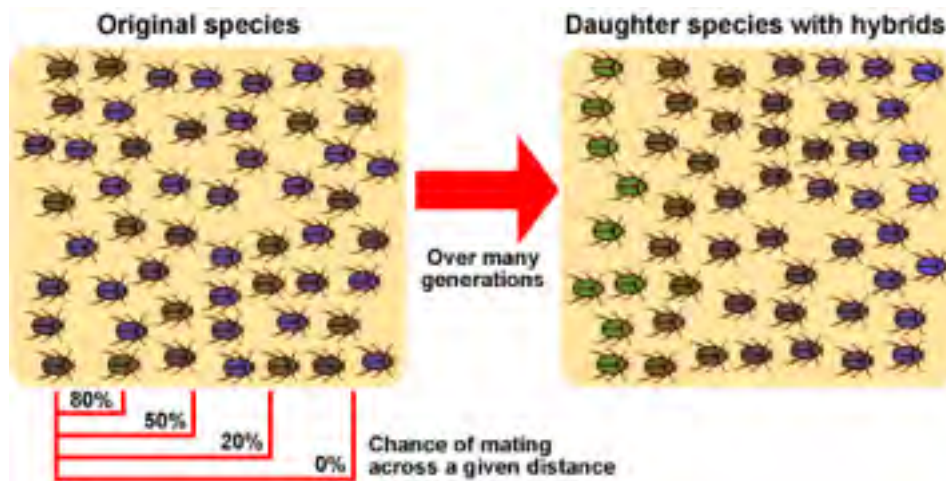
2. Speciation due to selection – without isolation



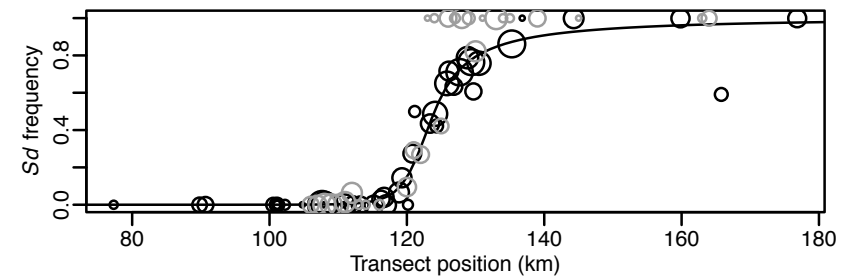
# Genomic divergence during speciation



1. Speciation as a bi-product of physical isolation



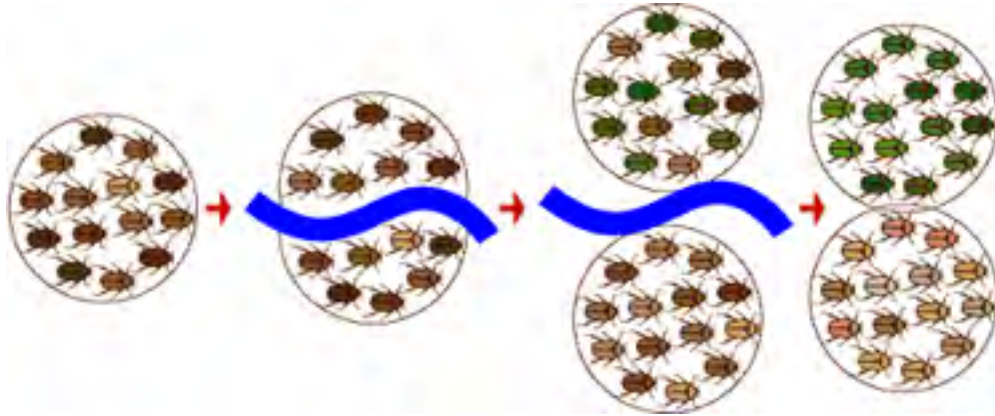
2. Speciation due to selection – without isolation



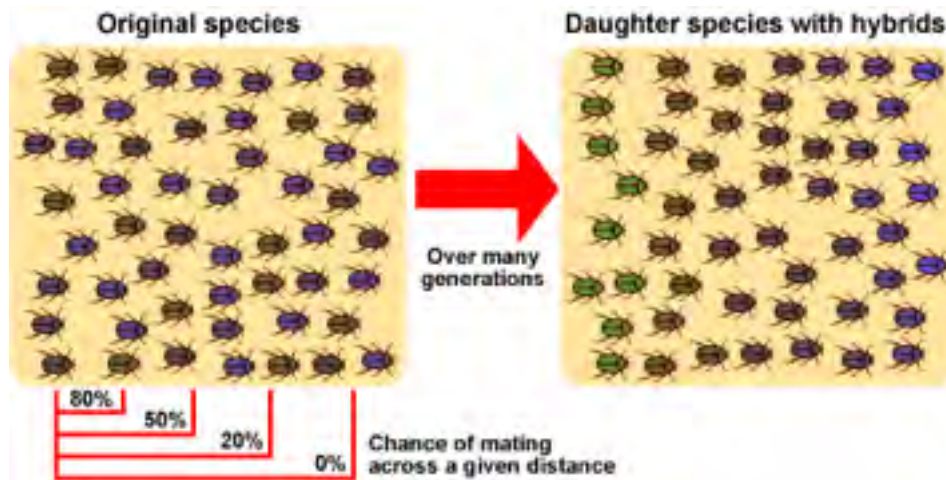
Cline theory - e.g. Barton and Gale 1993



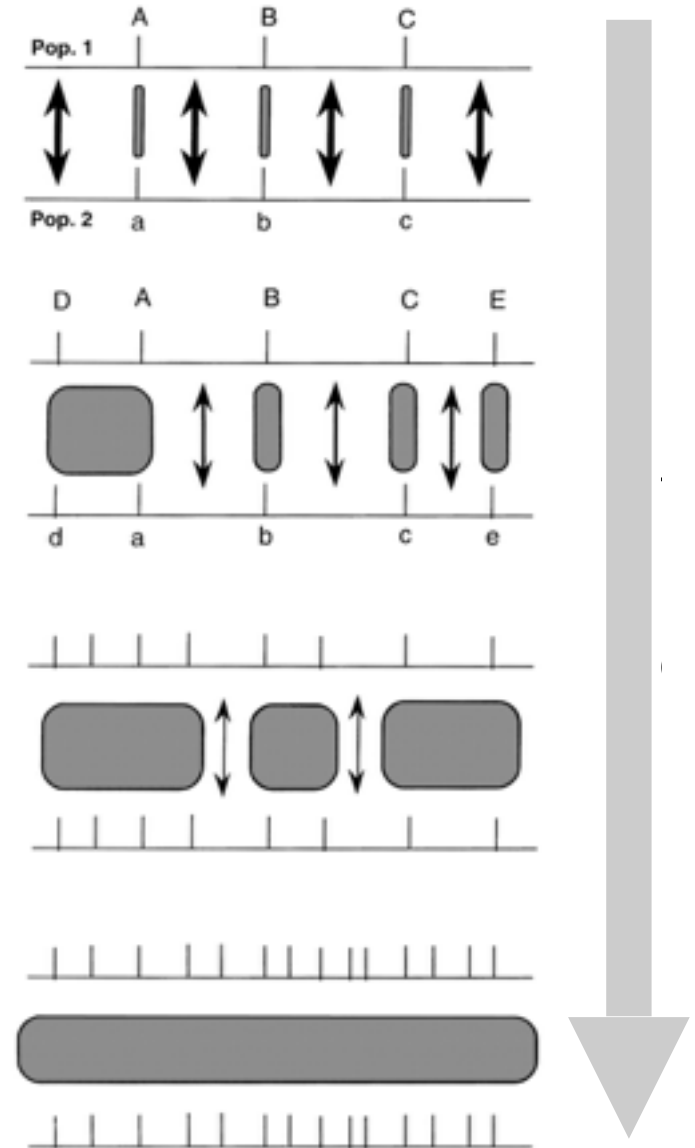
# Genomic divergence during speciation



1. Speciation as a bi-product of physical isolation



2. Speciation due to selection – without isolation



# Stage 1 - one or few loci under disruptive selection

Gene  
under  
selection



Genome

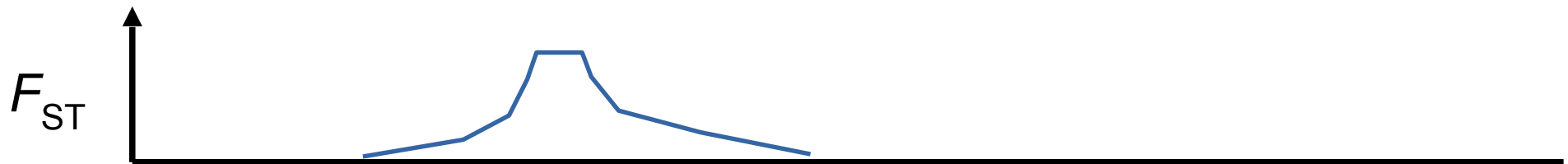


Feder, Egan and Nosil TiG

# Stage 2 - Divergence hitchhiking



Genome



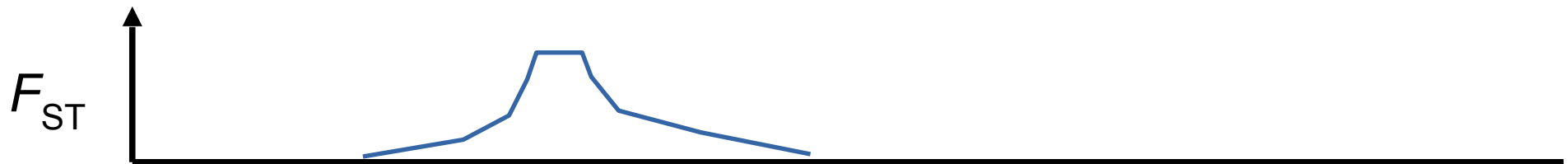
Feder, Egan and Nosil TiG

# Stage 2b - Inversion

Inversion links co-adapted alleles



Genome

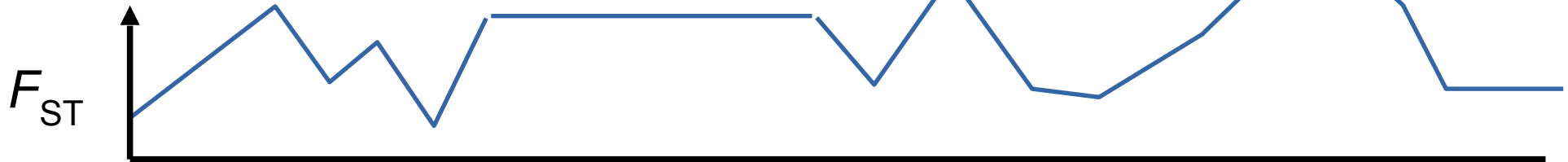


Feder, Egan and Nosil TiG

# Stage 3 - Genome hitchhiking



Genome

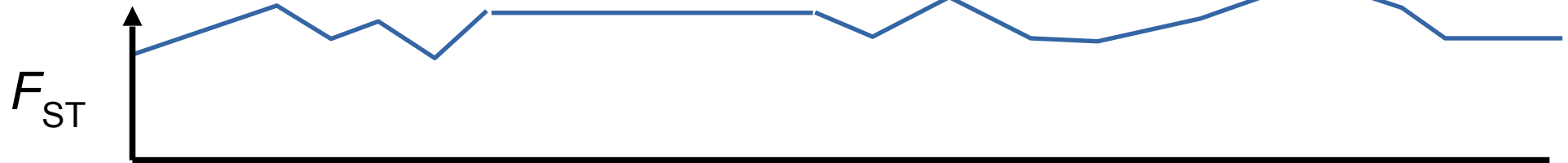


Feder, Egan and Nosil TiG

# Stage 4 - Genome wide isolation



Genome

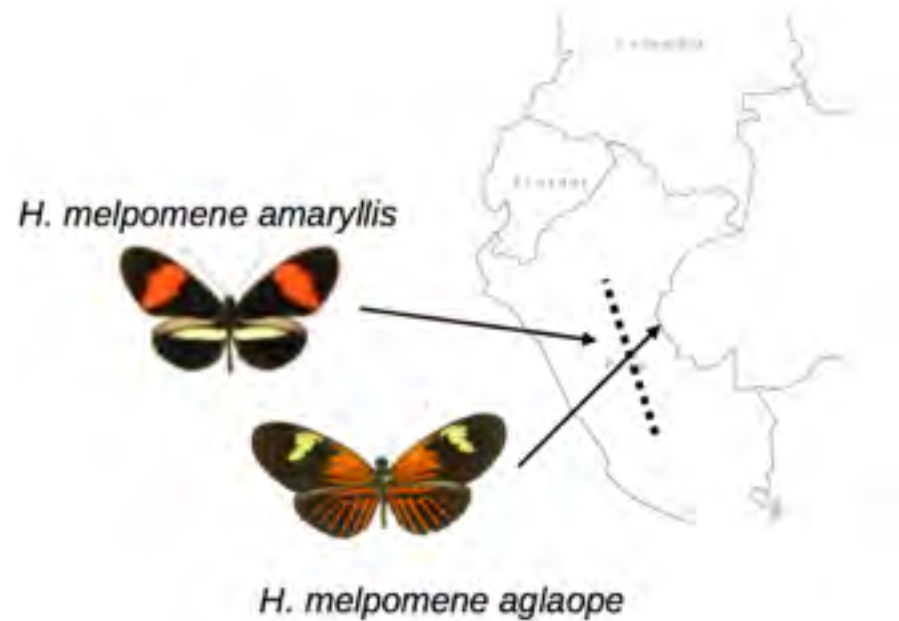
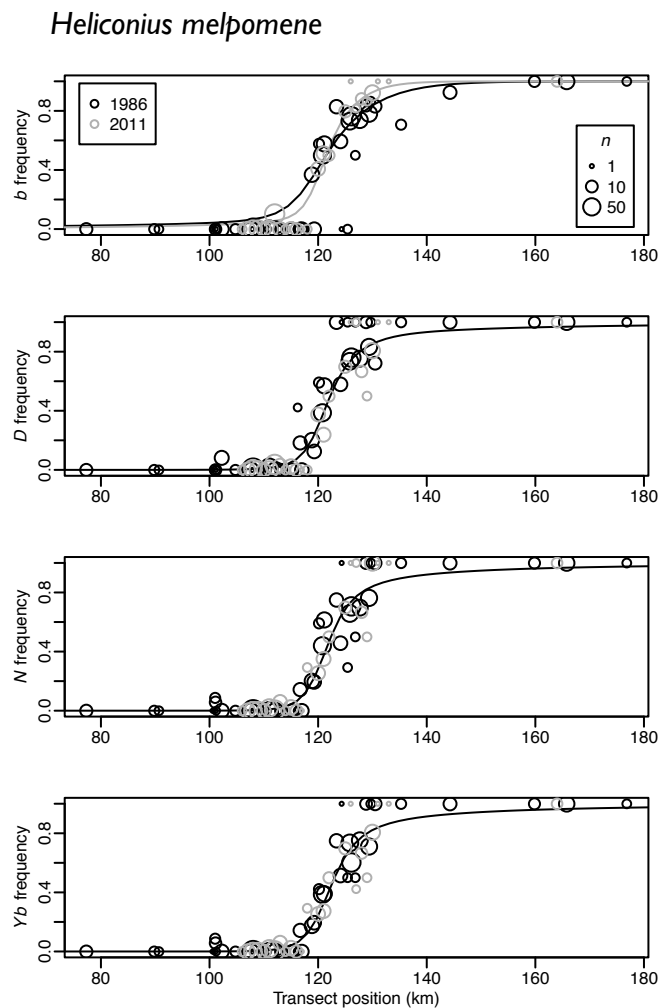


Feder, Egan and Nosil TiG



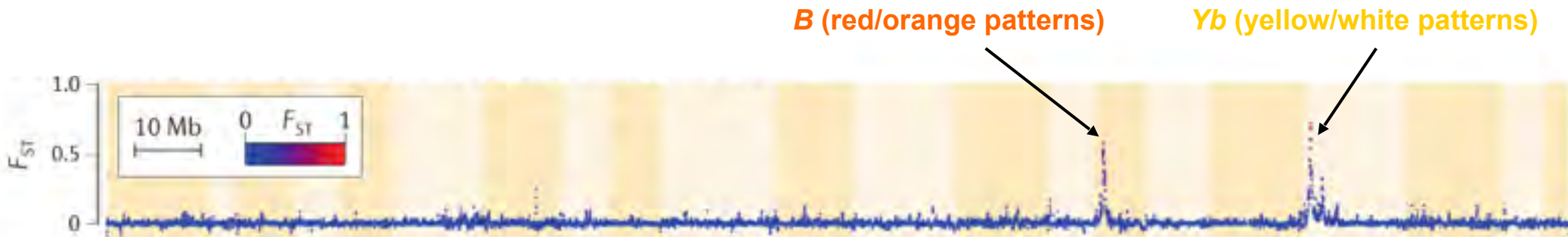
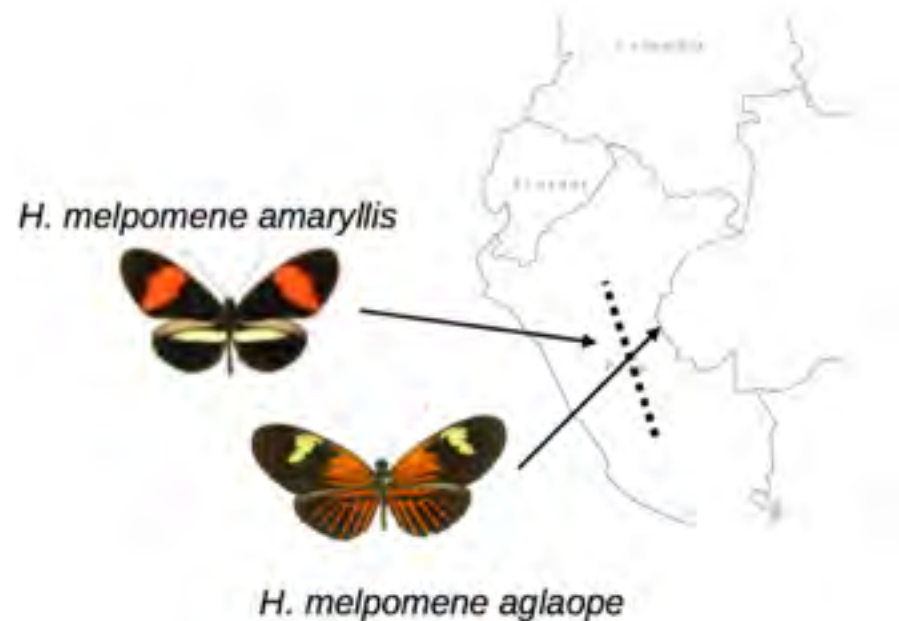
# Some sub-species clearly in stage 1

## Wing pattern “races” of *Heliconius melpomene*



# Some sub-species clearly in stage 1

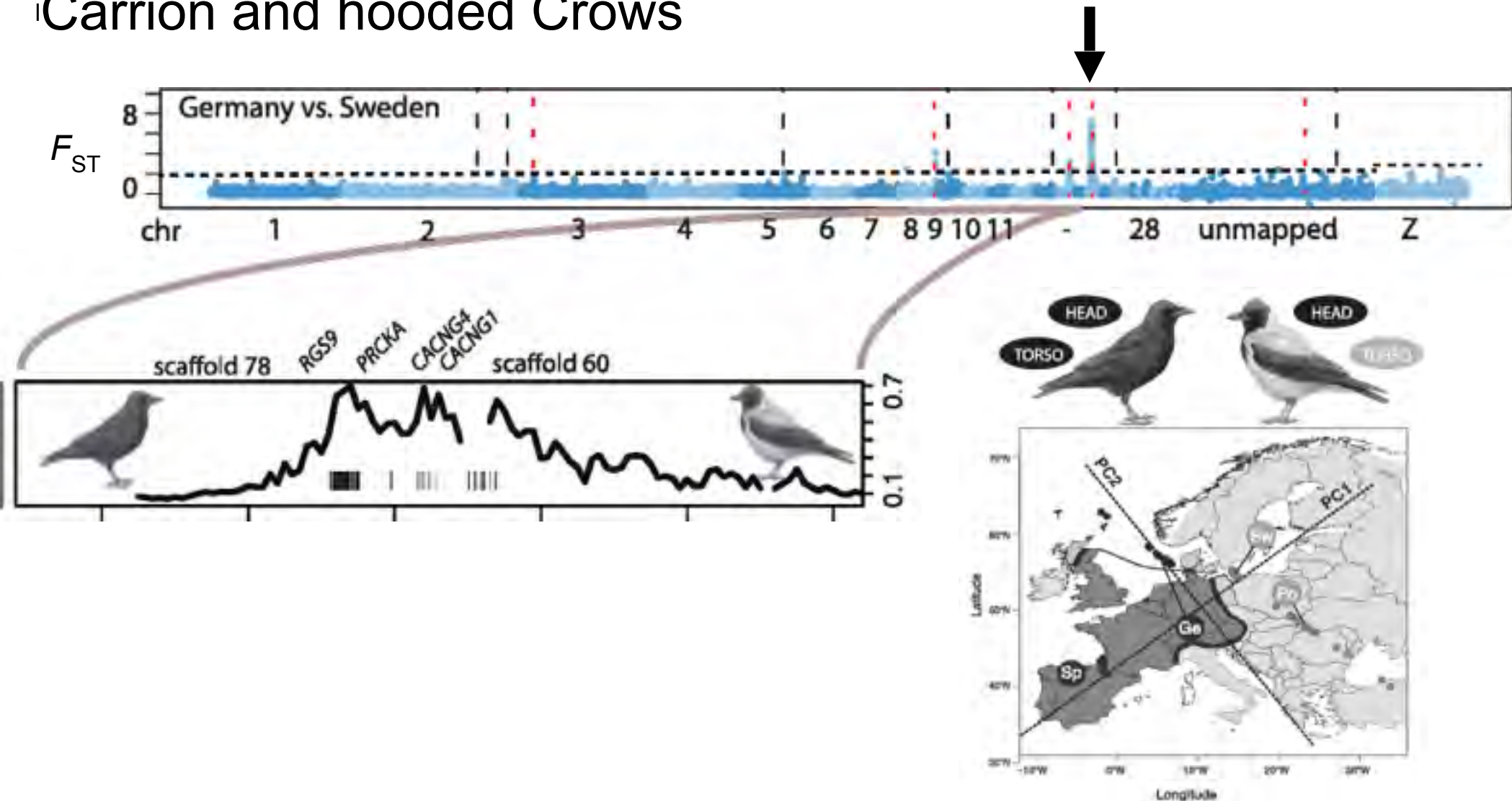
Wing pattern “races” of  
*Heliconius melpomene*



S. H. Martin et al. Genome Res. 23, 1817–1828 (2013).  
O. Seehausen et al. Nat. Rev. Genet. 15, 176–92 (2014).

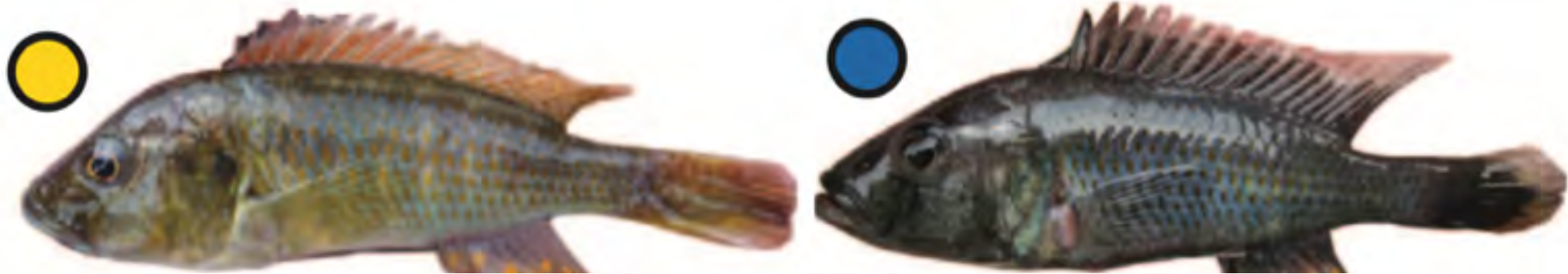
# Some sub-species clearly in stage 1

Carrion and hooded Crows

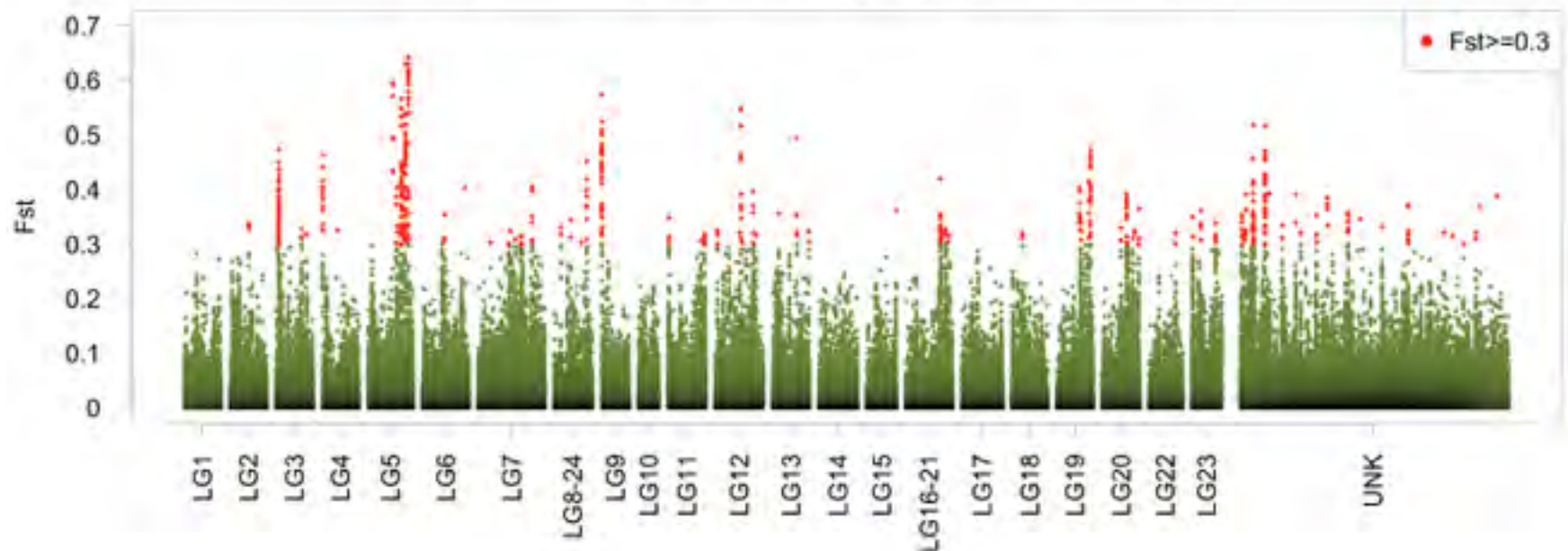


Poelstra, J. W. et al. Science 344, 1410–4 (2014).

# And here is a recent example with multiple islands

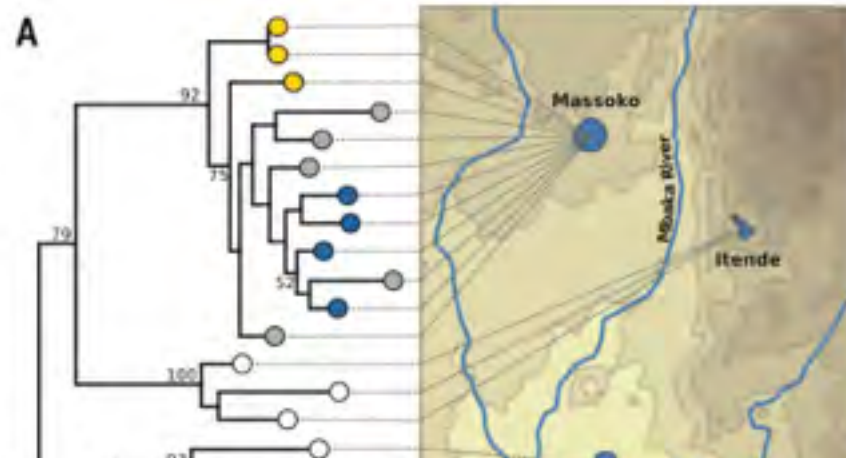


Massoko Benthic/Littoral Fst in 15 variant windows

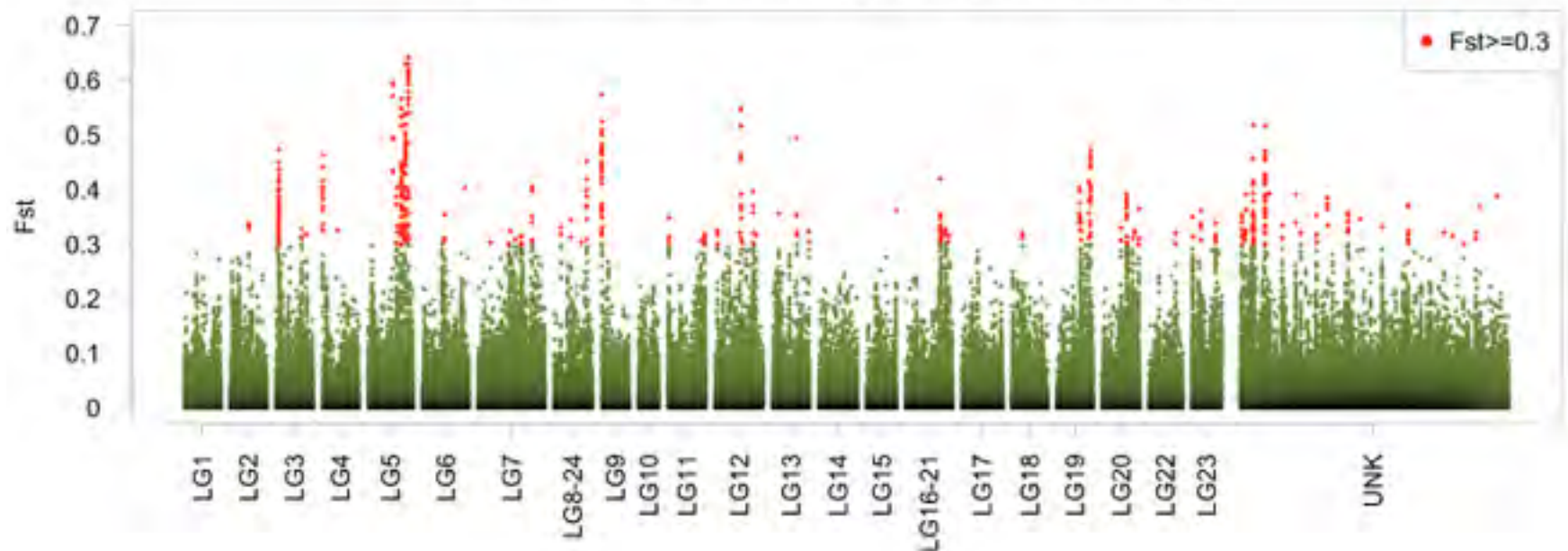




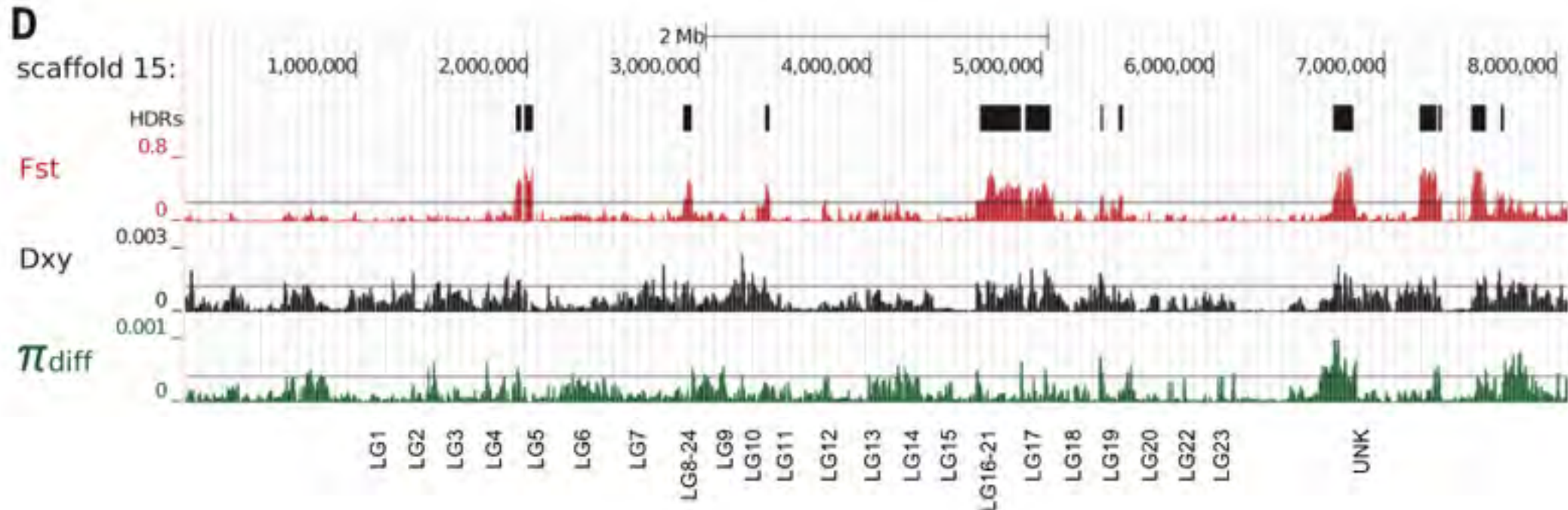
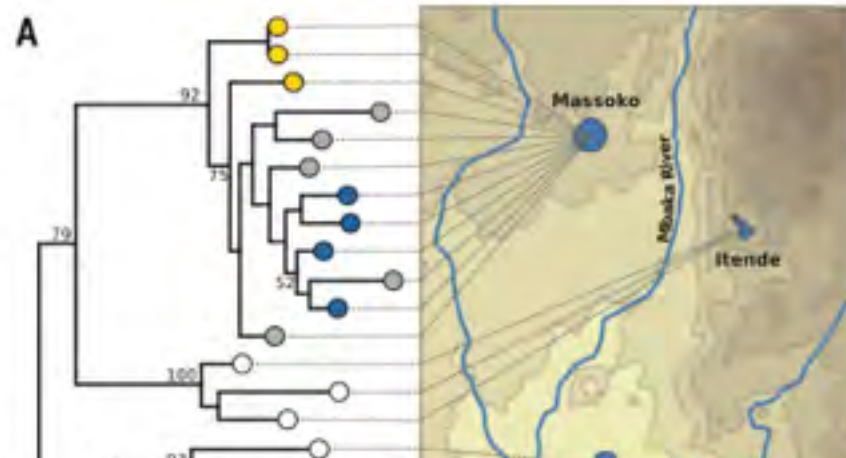
# And here is a recent example with multiple islands



Massoko Benthic/Littoral Fst in 15 variant windows



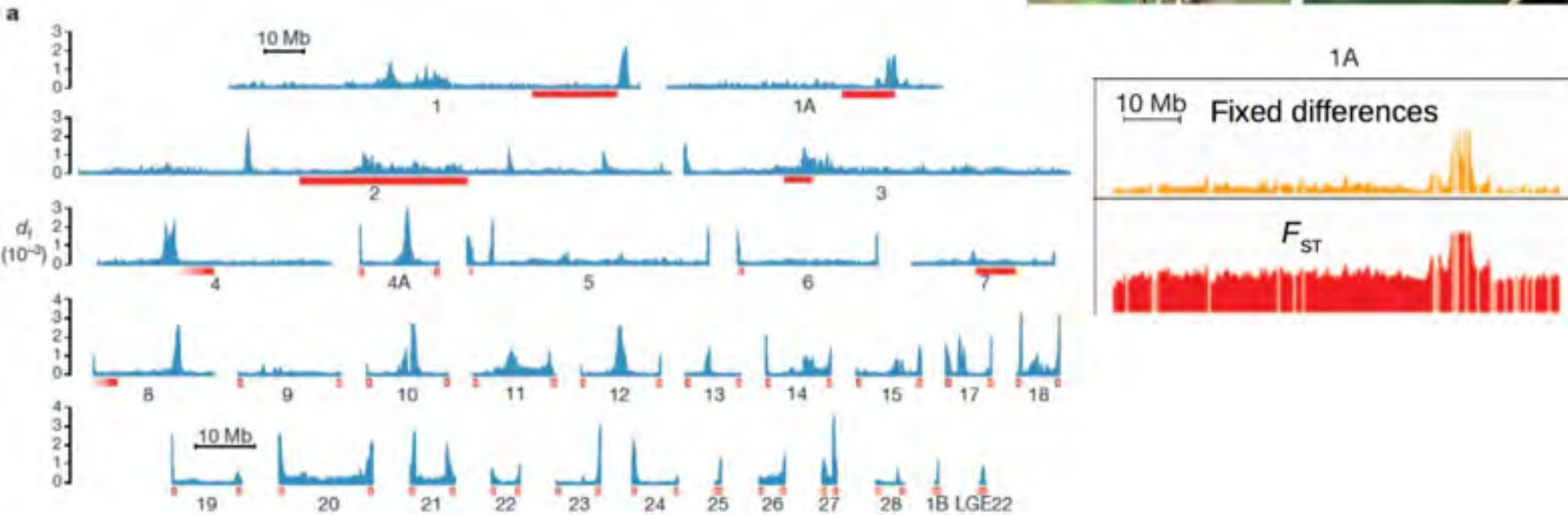
# And here is a recent example with multiple islands





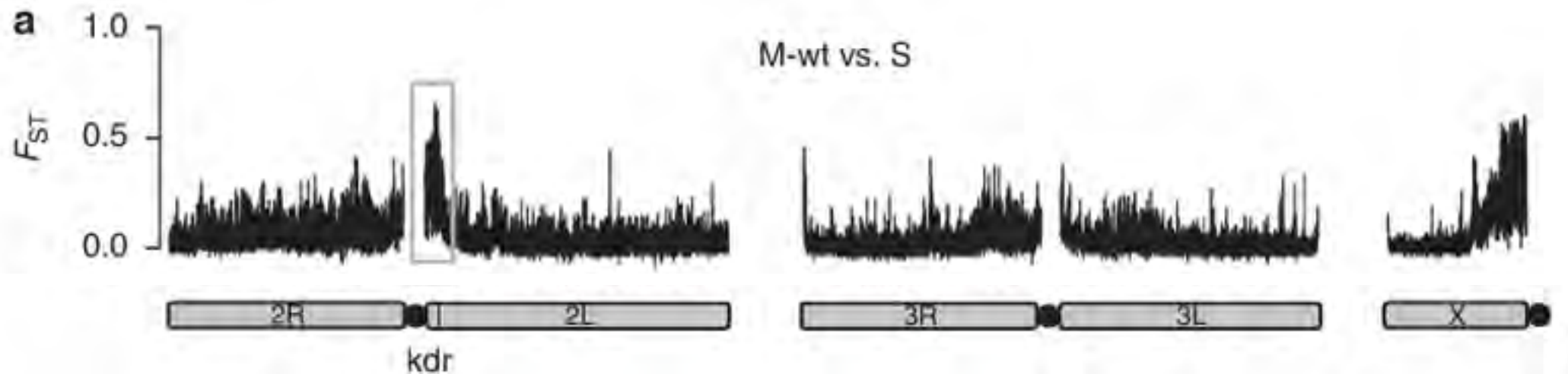
# Other species have islands...but are they real?

## Collared and Pied Flycatchers

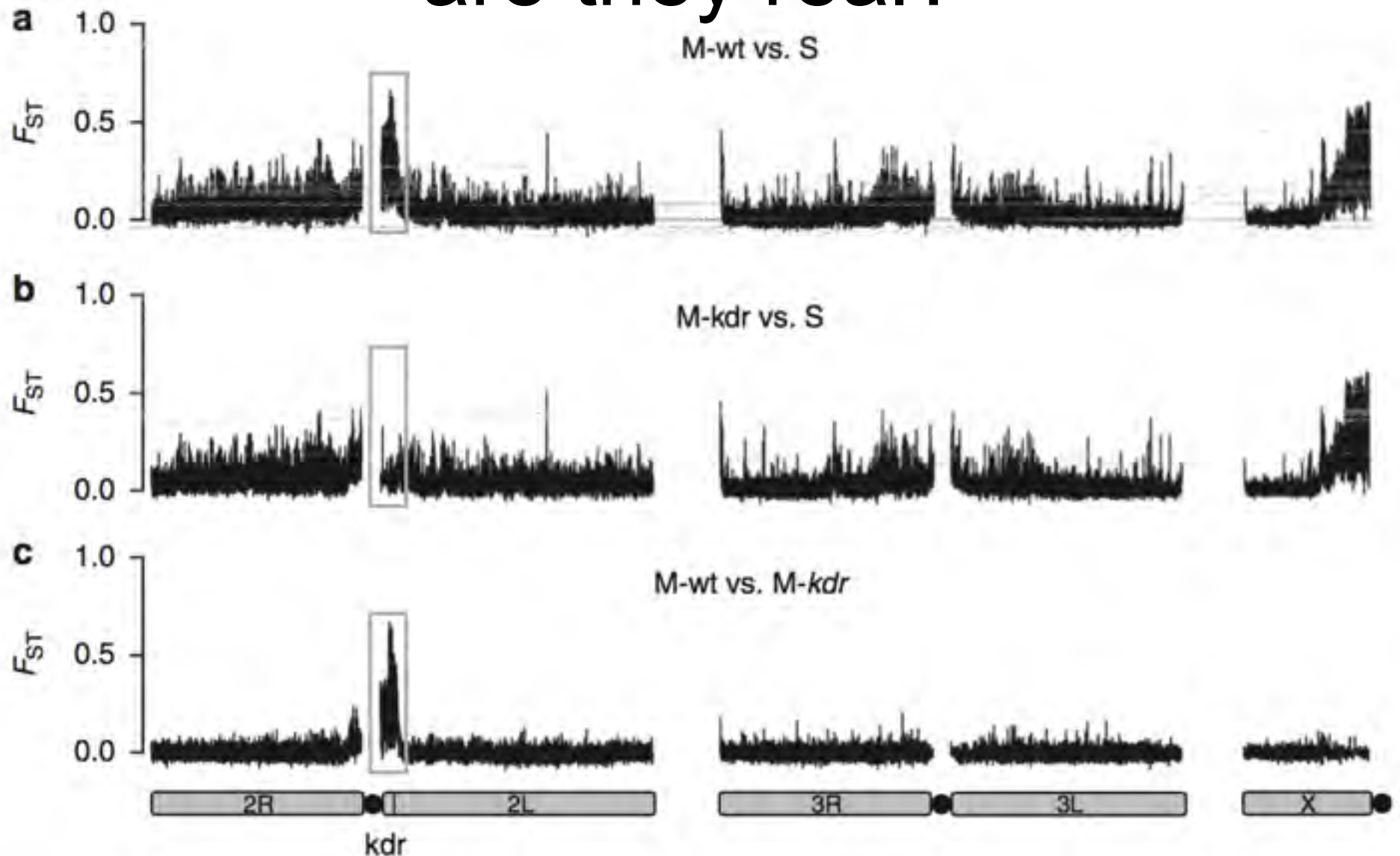


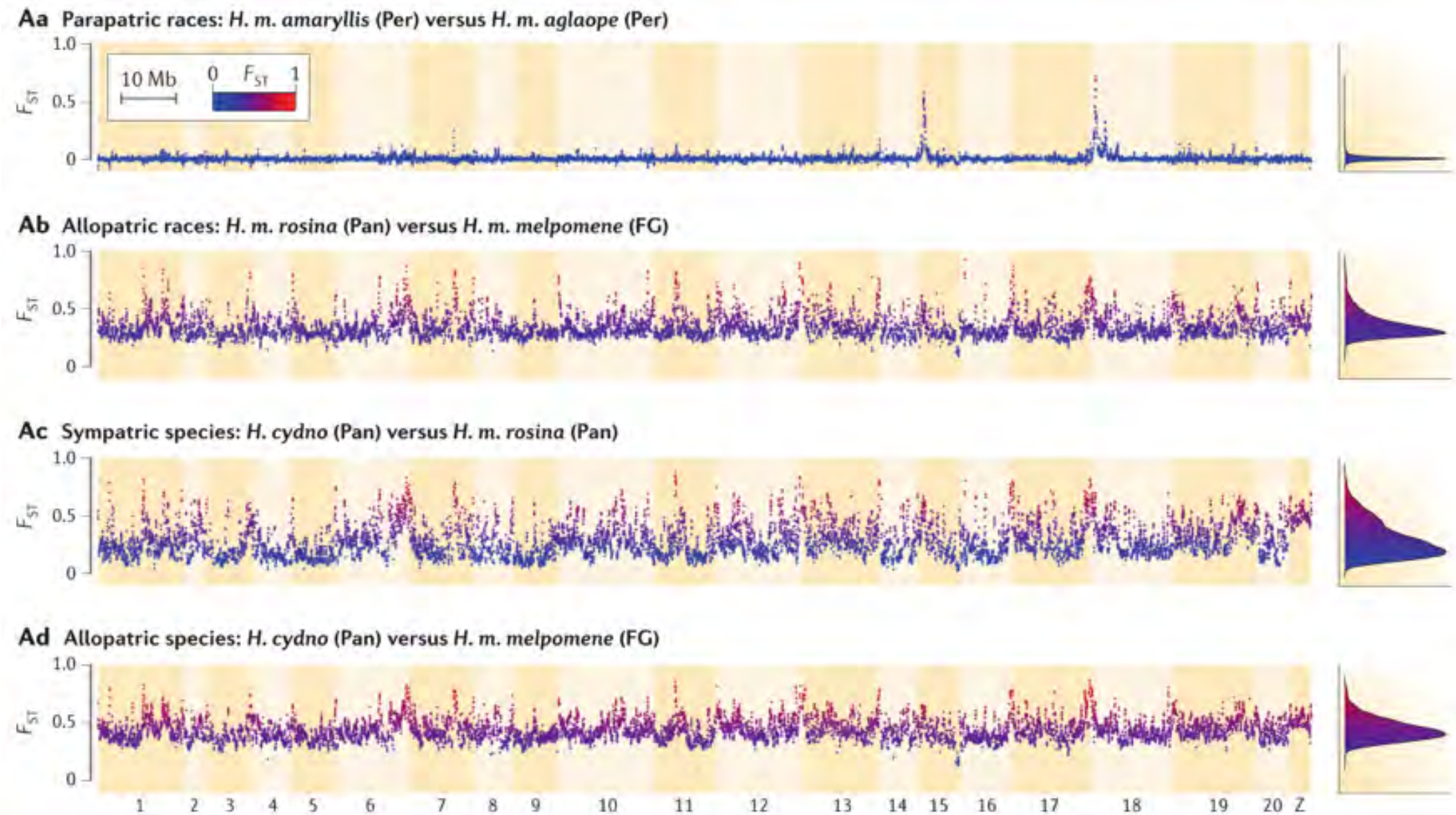
# Other species have islands...but are they real?

*Anopheles gambiae* and *A. coluzzi*  
Formerly *M* and *S* forms of *A. gambiae*



# Other species have islands...but are they real?





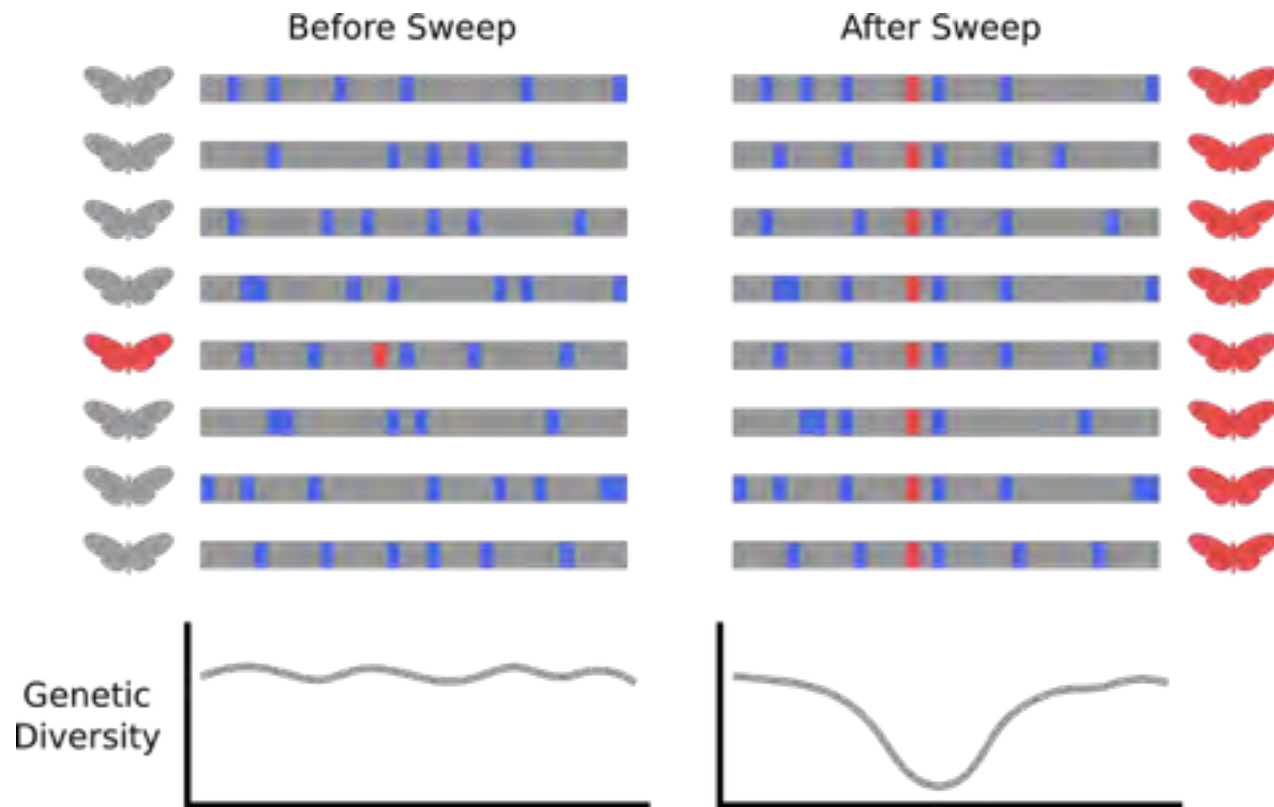


# What do patterns of $F_{st}$ really mean?

- $F_{st}$  measures relative divergence

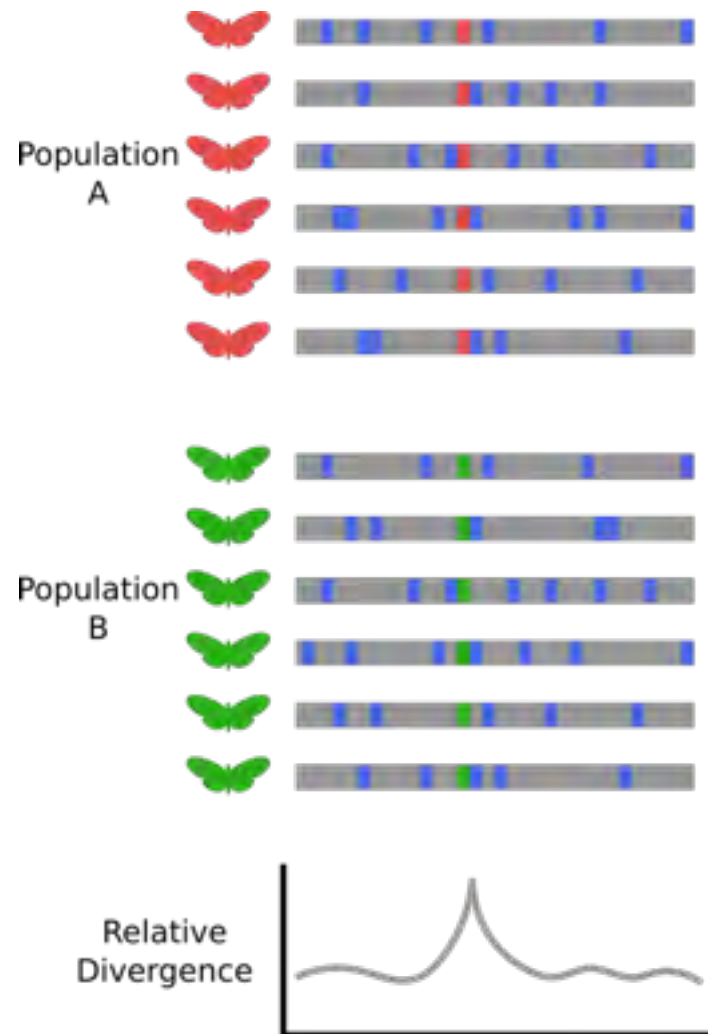
$$F_{ST} = \frac{H_T - H_S}{H_T},$$

- Peaks indicate regions of higher than expected between population divergence, given the within population divergence
- Peaks can therefore result from reduced diversity within species
- This could be due to lower  $N_e$  within species (selective sweeps, background selection)
- So peaks NOT NECESSARILY due to reduced gene flow



Note that sometimes sweeps within species = speciation genes





Sweeps across the species barrier can also lead to  $F_{st}$  peaks

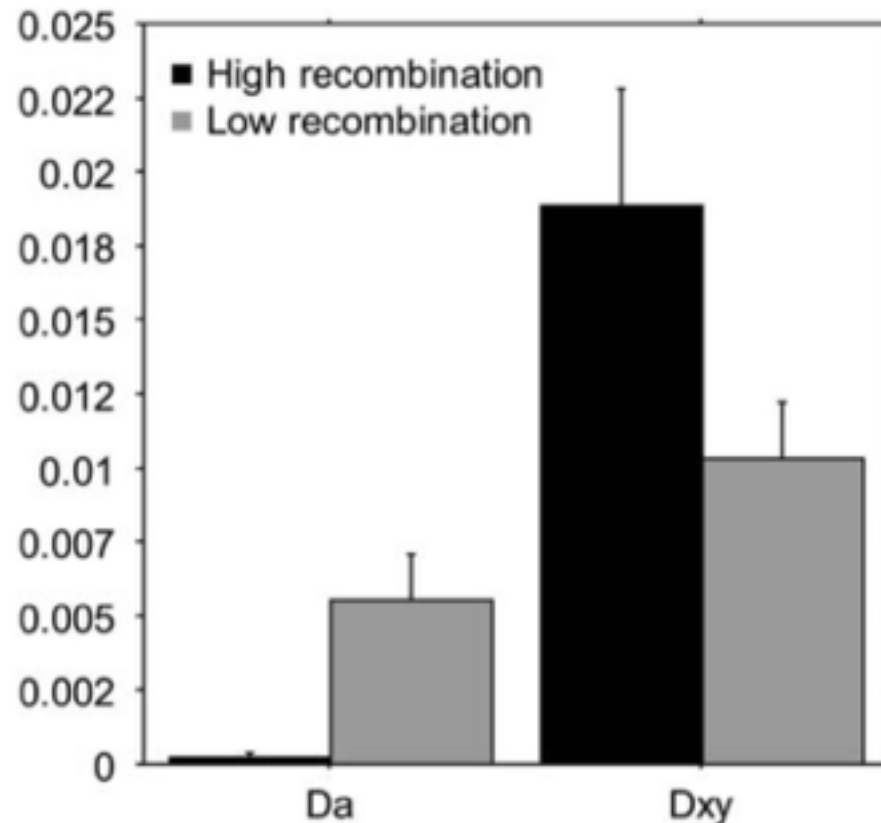
Double peaks??

## REVIEW

# Islands of speciation or mirages in the desert? Examining the role of restricted recombination in maintaining species

MAF Noor and SM Bennett

*Biology Department, Duke University, Durham, NC, USA*



Anopheles M-S  
divergence

Relative divergence  
higher in low  
recombination regions -  
not significant for absolute  
divergence

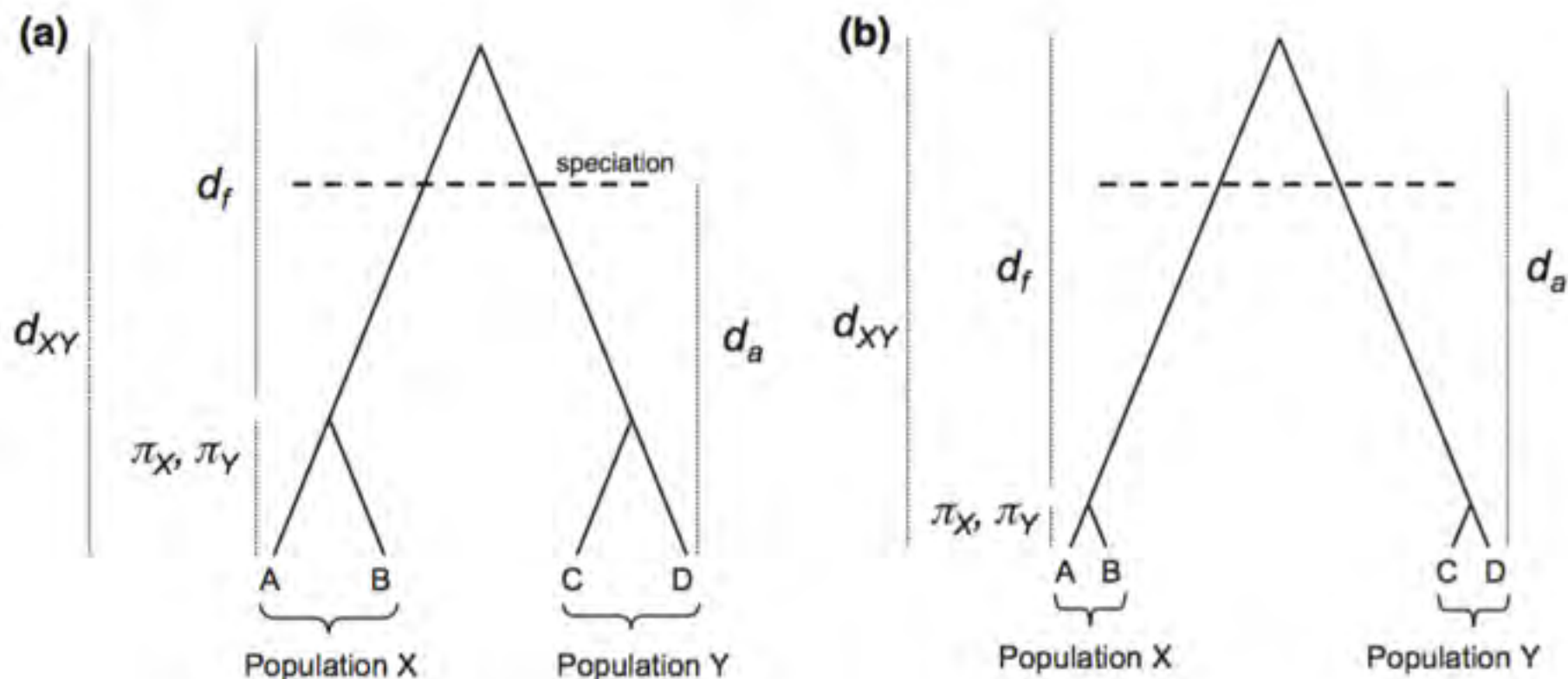
see also: Charlesworth 1998 MBE Measures of  
divergence...

## INVITED REVIEWS AND SYNTHESSES

# Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow

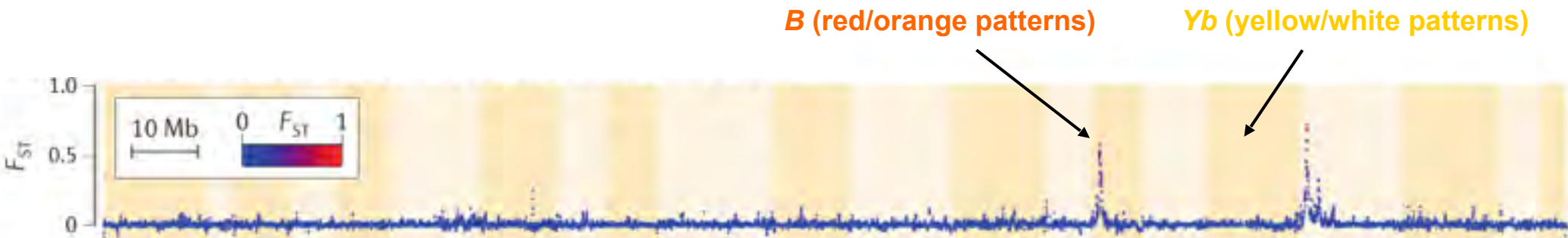
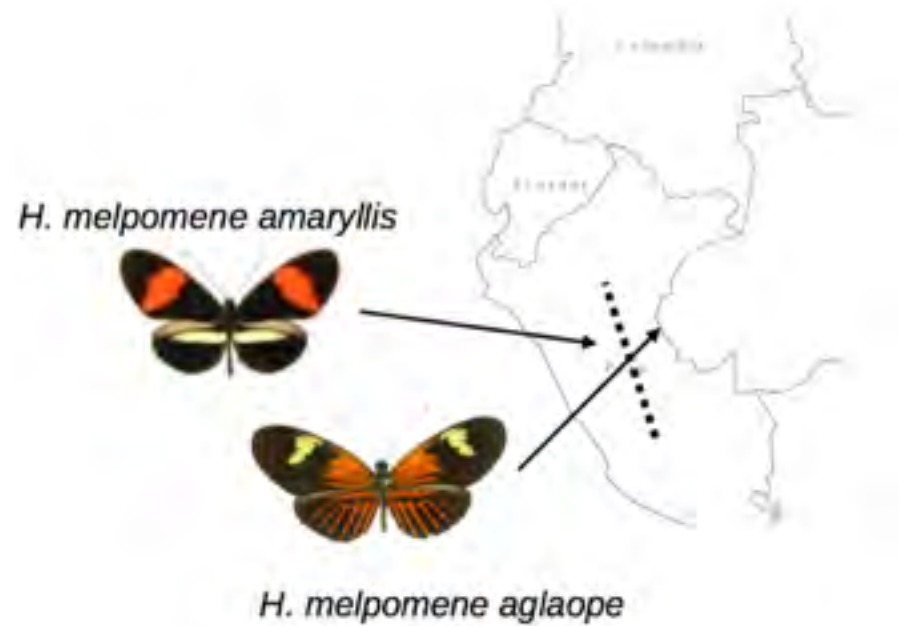
TAMI E. CRUICKSHANK\* and MATTHEW W. HAHN\*†

\*Department of Biology, Indiana University, Bloomington, IN 47405, USA, †School of Informatics and Computing, Indiana University, Bloomington, IN 47405, USA



# No evidence for higher Dxy in wing pattern loci

Wing pattern “races” of *Heliconius melpomene*



S. H. Martin et al. Genome Res. 23, 1817–1828 (2013).  
O. Seehausen et al. Nat. Rev. Genet. 15, 176–92 (2014).



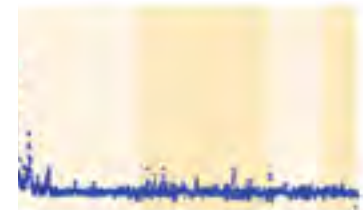
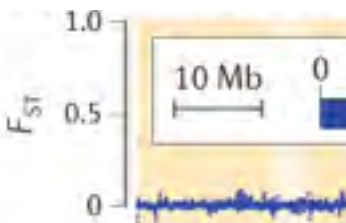
# No evidence for higher Dxy in wing pattern loci

Wing pat  
*Heliconiu*

One further issue with interpreting the data from these two races is whether this comparison relates to speciation at all. There is strong geographic structure involving the wing colour patterns that define these morphs as races, largely due to selection determined by colour morphs in the Müllerian mimic, *H. erato* (Mallet *et al.* 1990). But the races are not separate species: they do not show evidence of hybrid sterility or inviability and appear to be randomly mating in the narrow zone where the colour morphs overlap (Mallet *et al.* 1990). This raises the possibility that the colour-patterning loci contain locally adapted alleles within a largely panmictic (or at least continuously distributed) population and that gene flow outside of these regions represents nothing more than the normal movement of alleles within a species. In this case, there should be

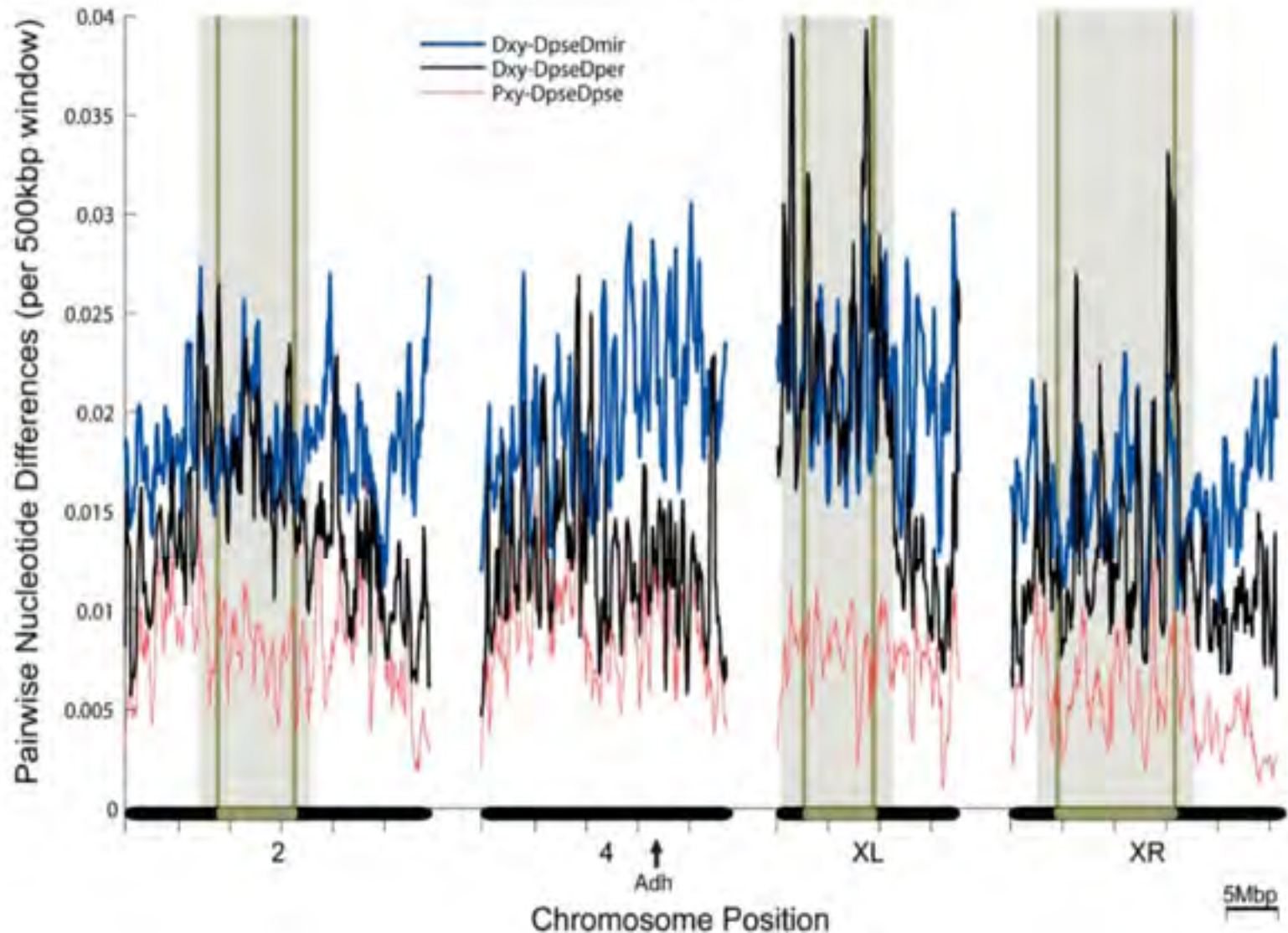


v/white patterns)



8 (2013).  
-92 (2014).

# Suggestion that we use absolute measures of divergence?





# Understanding genomic divergence

No single statistic will capture the complex history of mutation, migration and selection

Patterns need to be interpreted in the specific context of the study species

# Much better to use explicit tests for gene flow

Need to design sampling so the expectations in the absence of gene flow are clear and testable

# Much better to use explicit tests for gene flow

Need to design sampling so the expectations in the absence of gene flow are clear and testable

The key is to identify 'control' populations that are not influenced by admixture

# Explicit tests for gene flow: Neanderthal genome

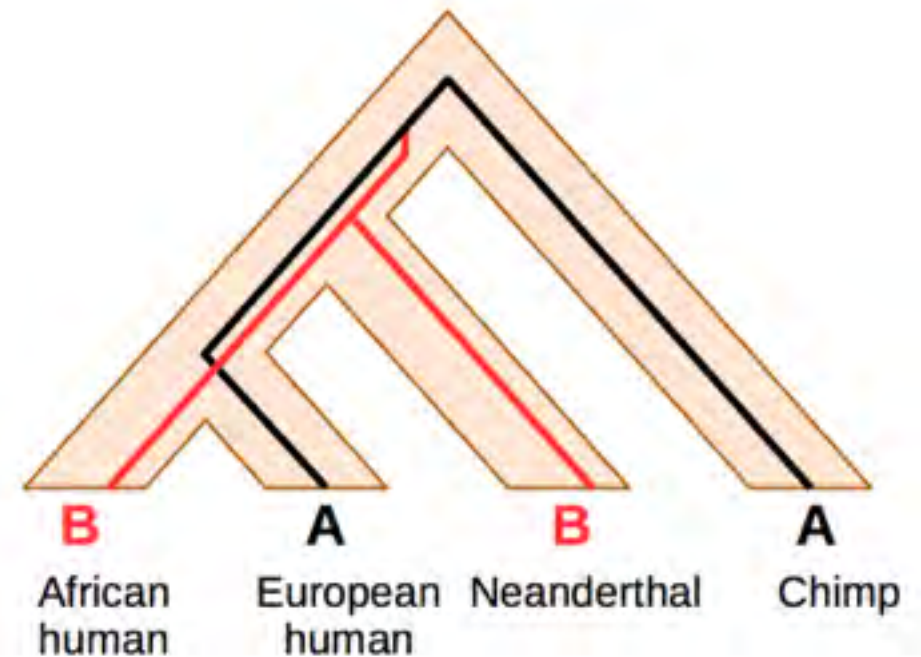
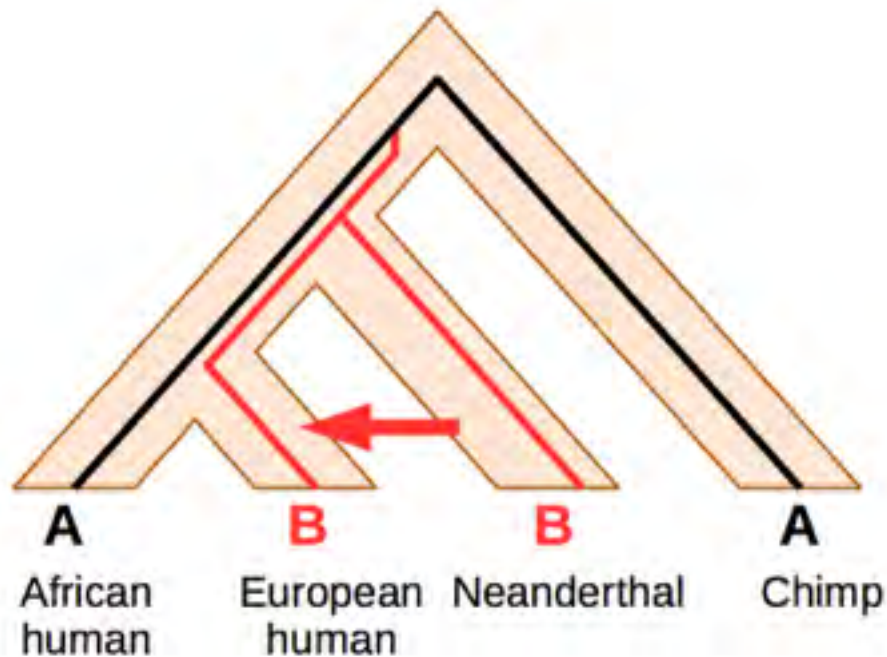


- Isolated DNA from bones 38,000 yrs old in Croatia
- We diverged from Neanderthals around 270-440,000 yrs ago
- Evidence for gene exchange with humans (1-4% of genome?)



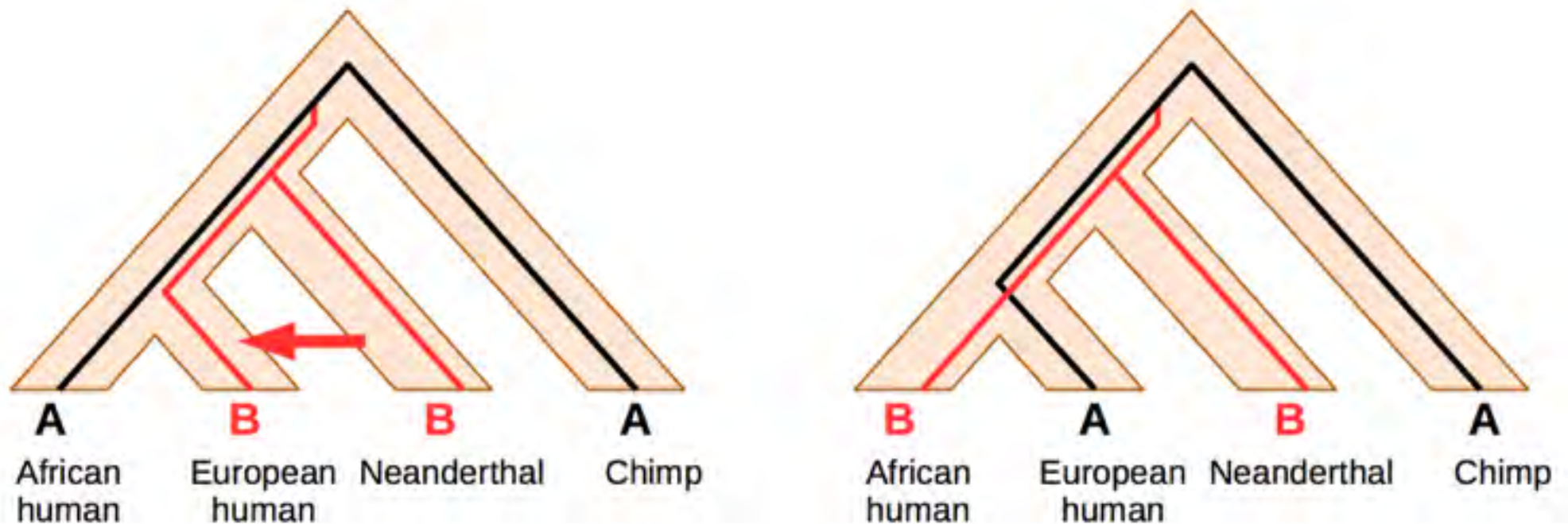
Green et al., 328:710 Science 2010

# Explicit tests for gene flow: ABBA-BABA test



$$D(P_1, P_2, P_3, O) = \frac{\sum C_{ABBA}(i) - C_{BABA}(i)}{\sum C_{ABBA}(i) + C_{BABA}(i)} \quad (1)$$

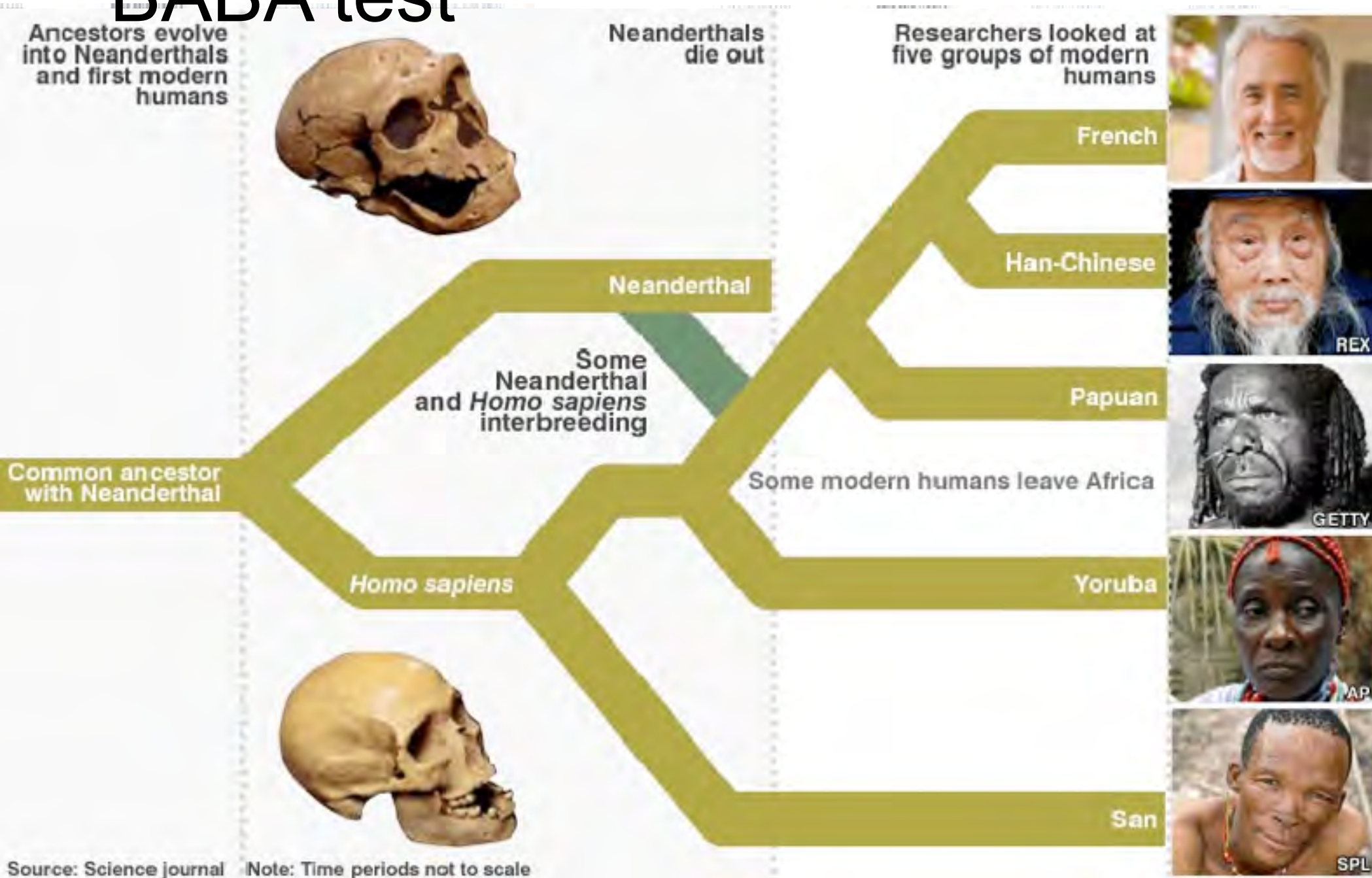
# Explicit tests for gene flow: ABBA-BABA test



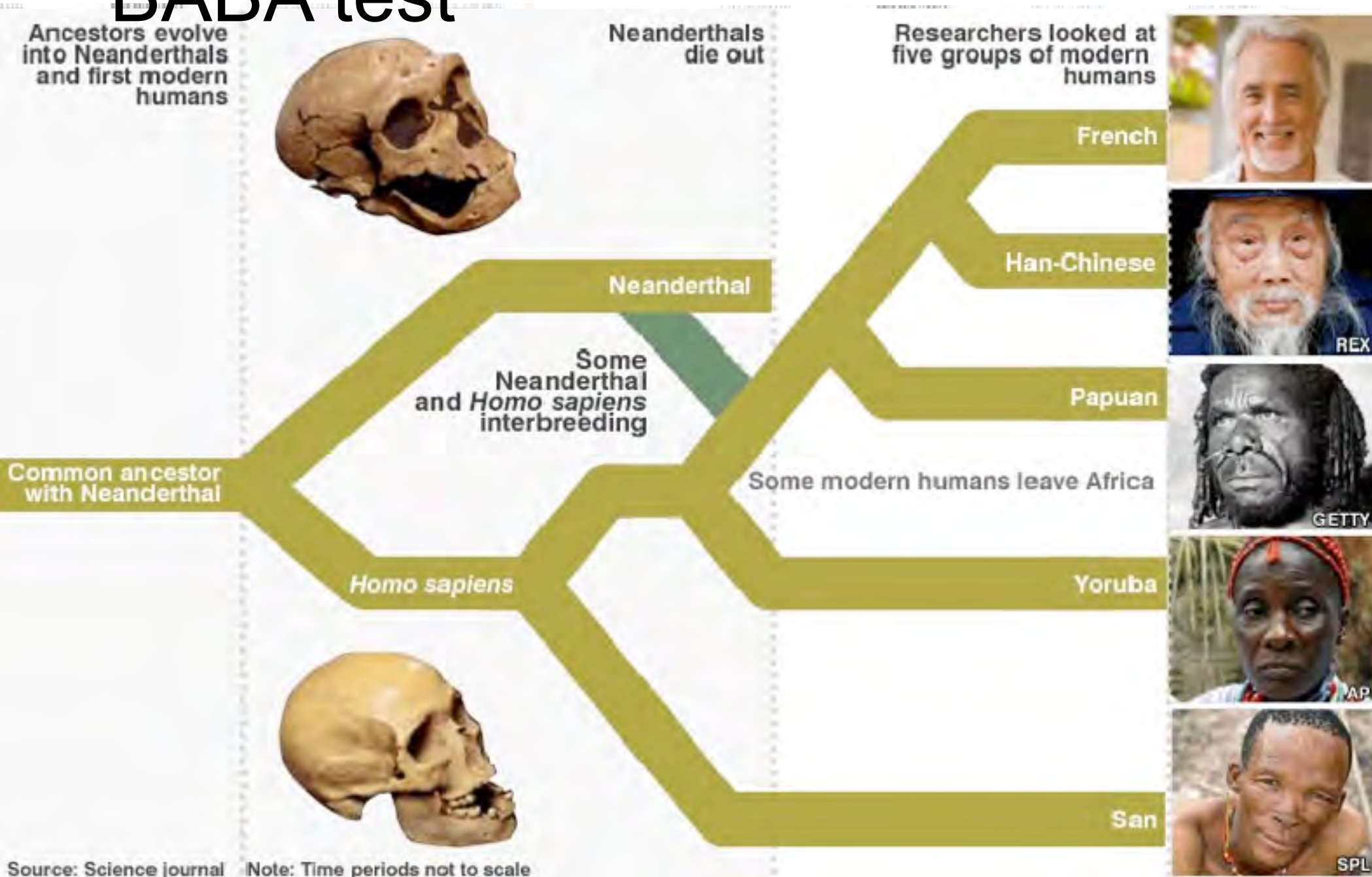
$$D(P_1, P_2, P_3, O) = \frac{\sum C_{ABBA}(i) - C_{BABA}(i)}{\sum C_{ABBA}(i) + C_{BABA}(i)} \quad (1)$$



# Explicit tests for gene flow: ABBA-BABA test

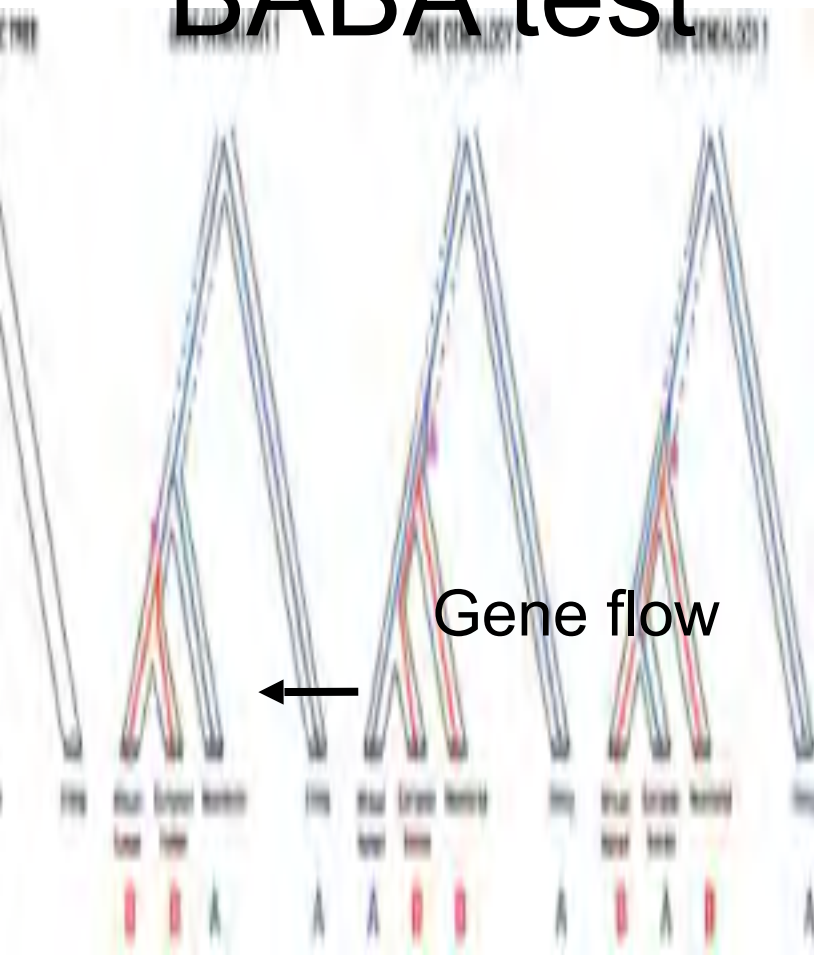


# Explicit tests for gene flow: ABBA-BABA test

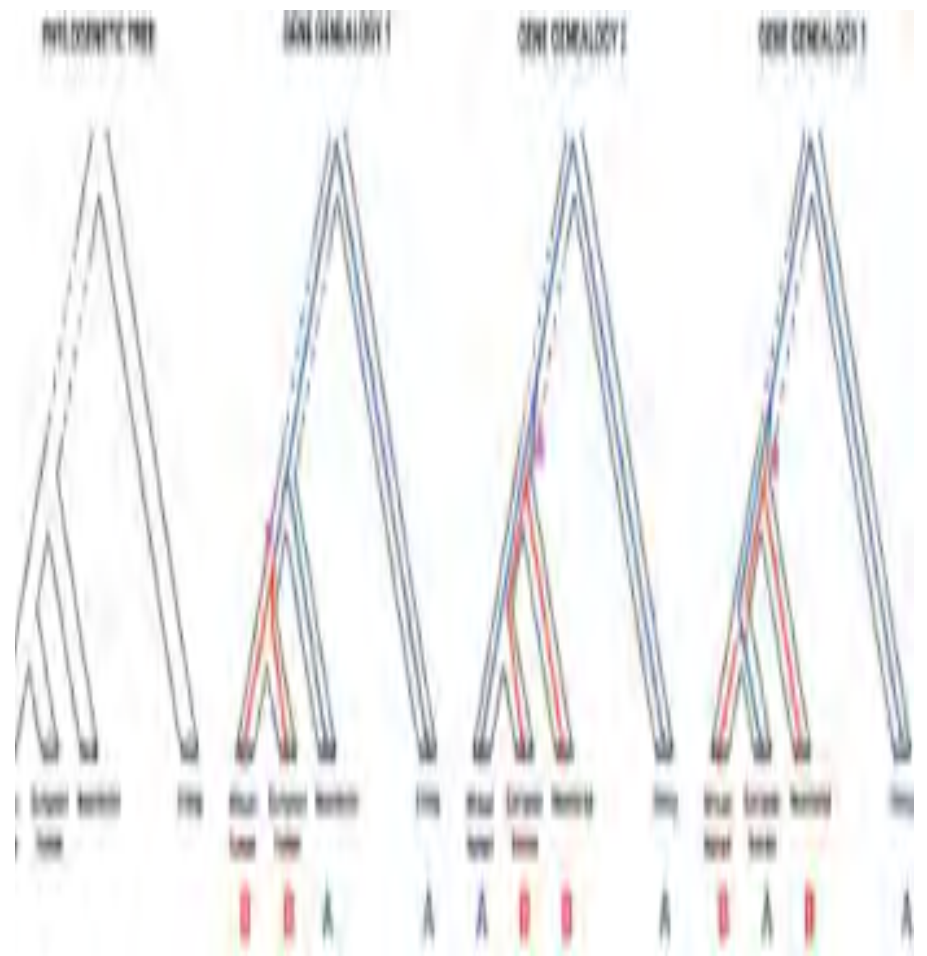




# Explicit tests for gene flow: ABBA-BABA test



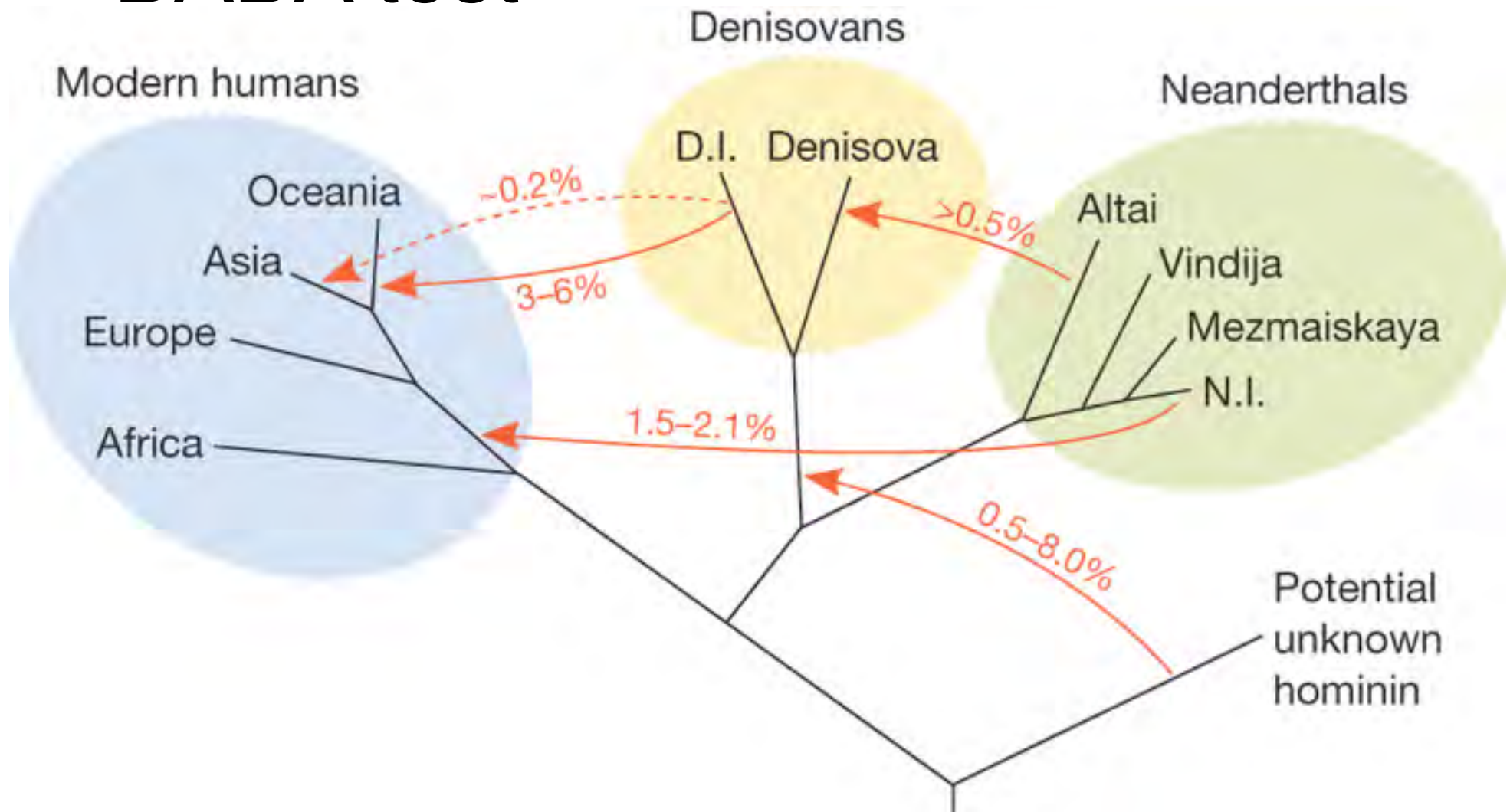
EXPECT:  
50% ABBA  
50% BABA



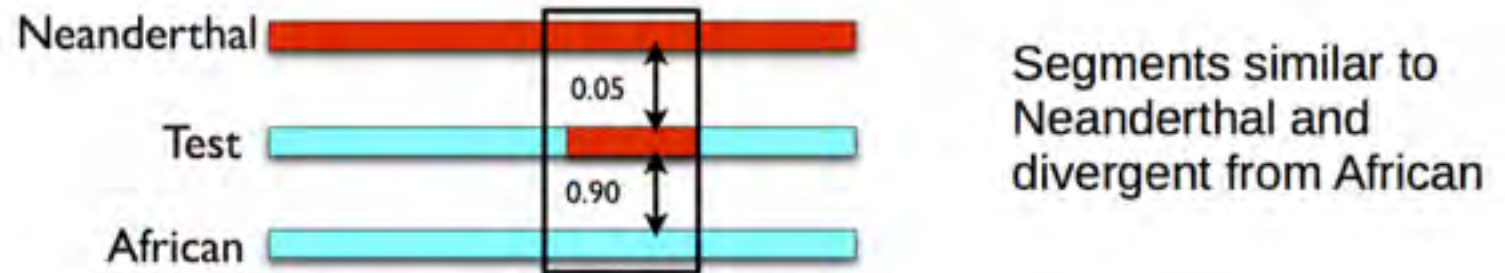
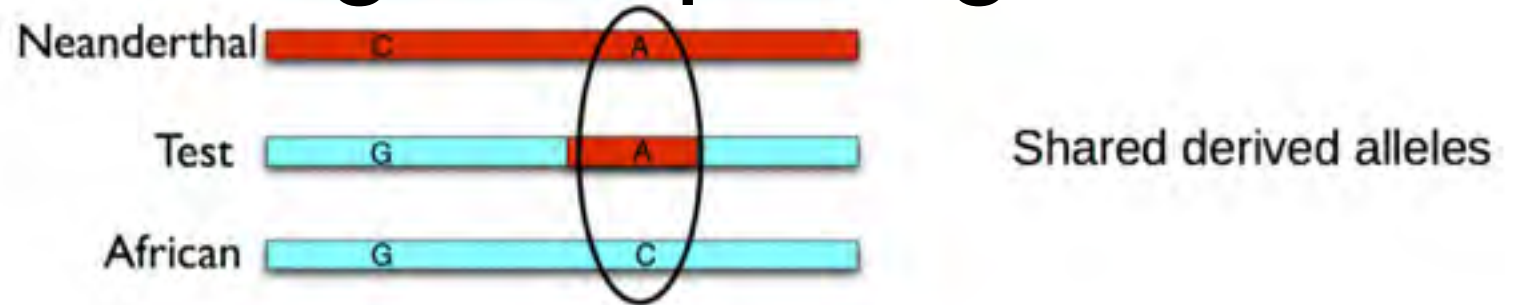
OBSERVE:  
103612 ABBA  
94029 BABA

Green *et al.* 2010  
*Science* 328:710-722

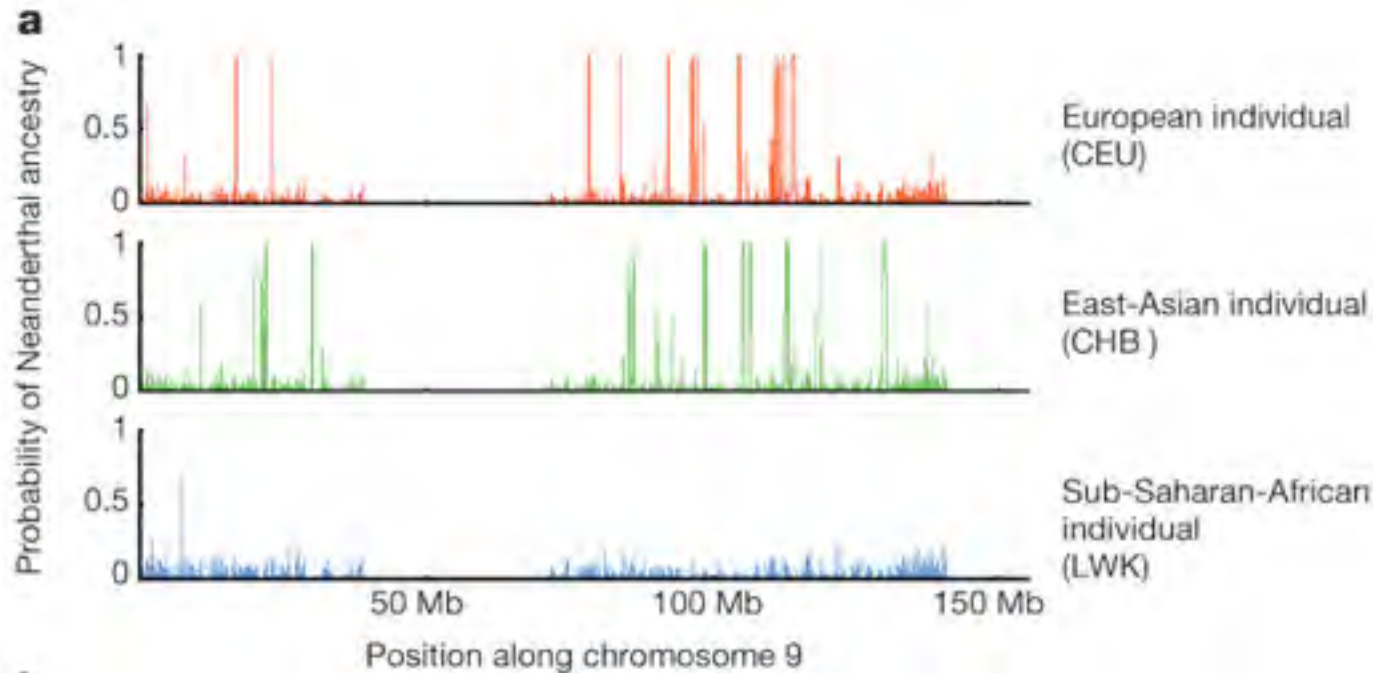
# Explicit tests for gene flow: ABBA-BABA test



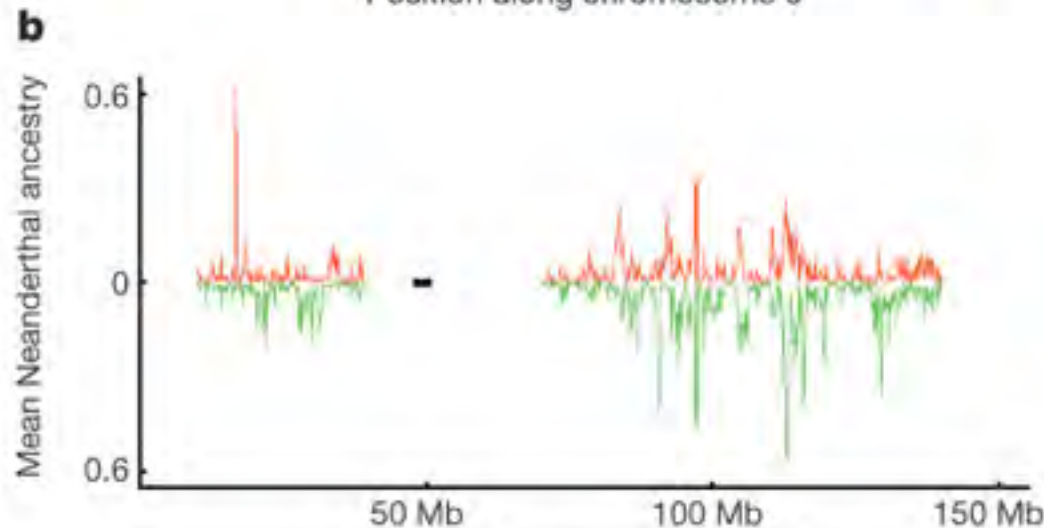
# Explicit tests for gene flow: Combining multiple signals



# Explicit tests for gene flow: Combining multiple signals



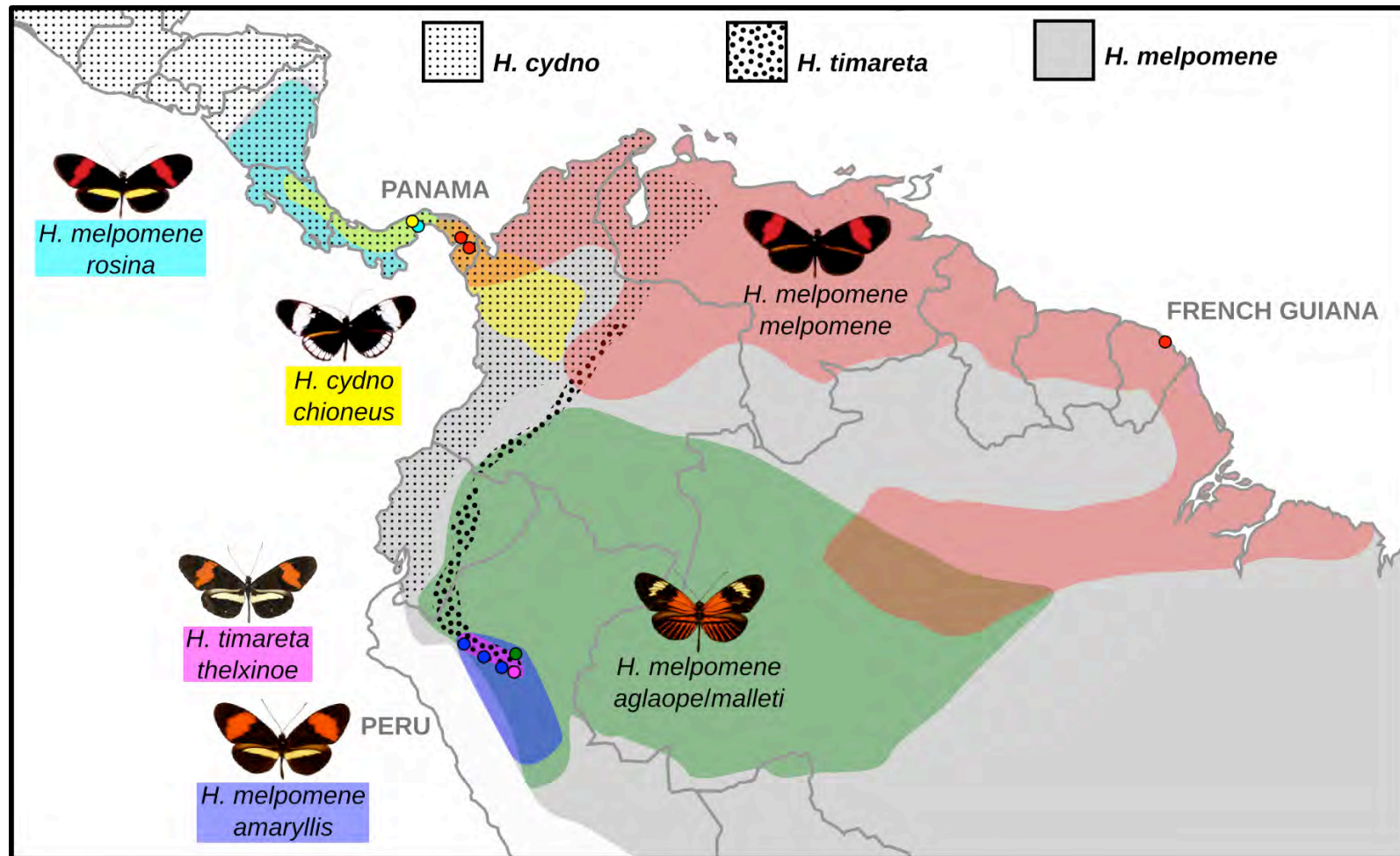
- 1) Derived alleles at high frequency shared with Neanderthal
- 2) High divergence to Africa but low to Neanderthal
- 3) Long haplotype blocks



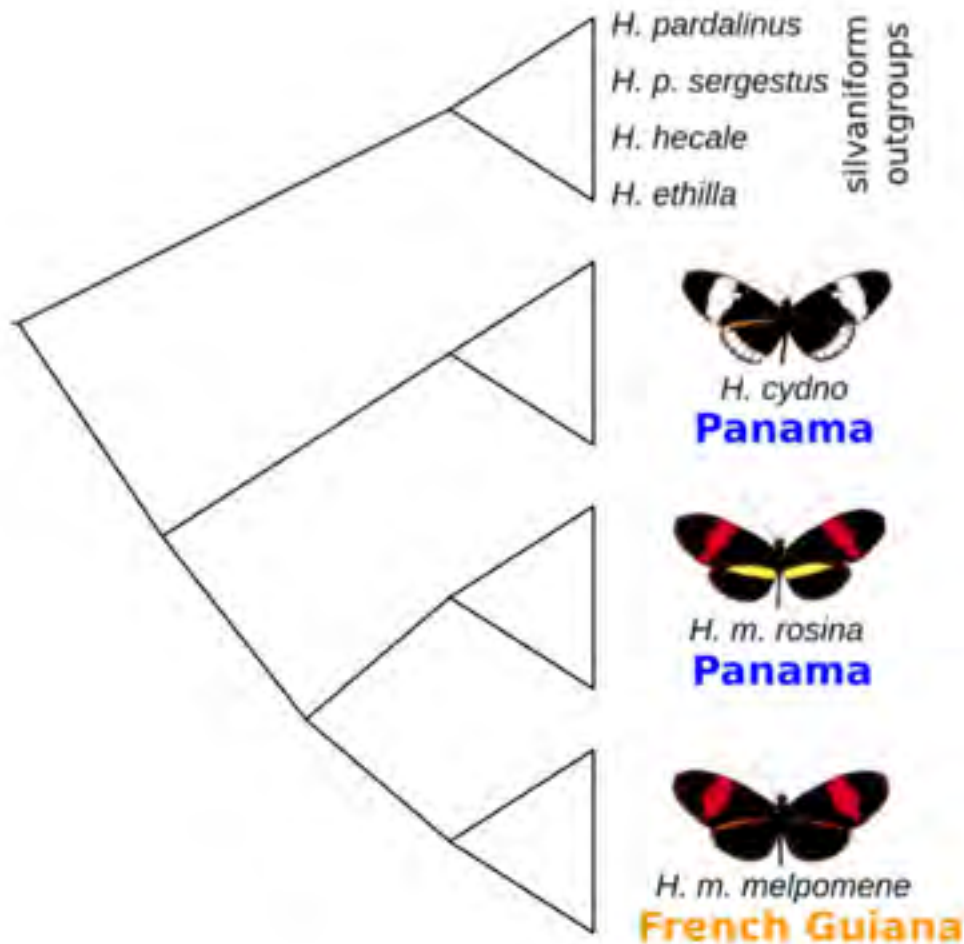
**The genomic landscape of  
Neanderthal ancestry in  
present-day humans -  
Sankararaman et al. Nature  
2014**



# Explicit tests for gene flow: *Heliconius* butterflies



# Explicit tests for gene flow: *Heliconius* butterflies



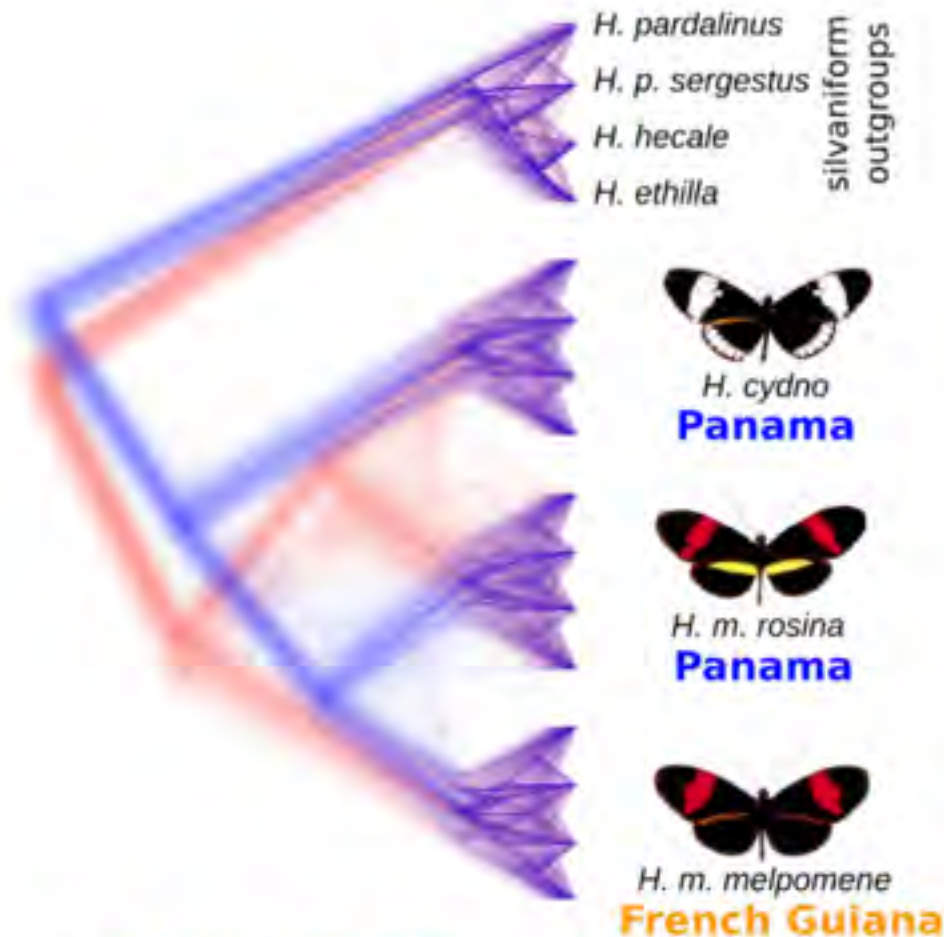
Whole-genome  
phylogeny supports  
grouping by species

Many sources of reproductive  
isolation:

Female hybrids are sterile  
Different host plant use  
Different habitat preference  
Strong assortative mating



# Explicit tests for gene flow: *Heliconius* butterflies

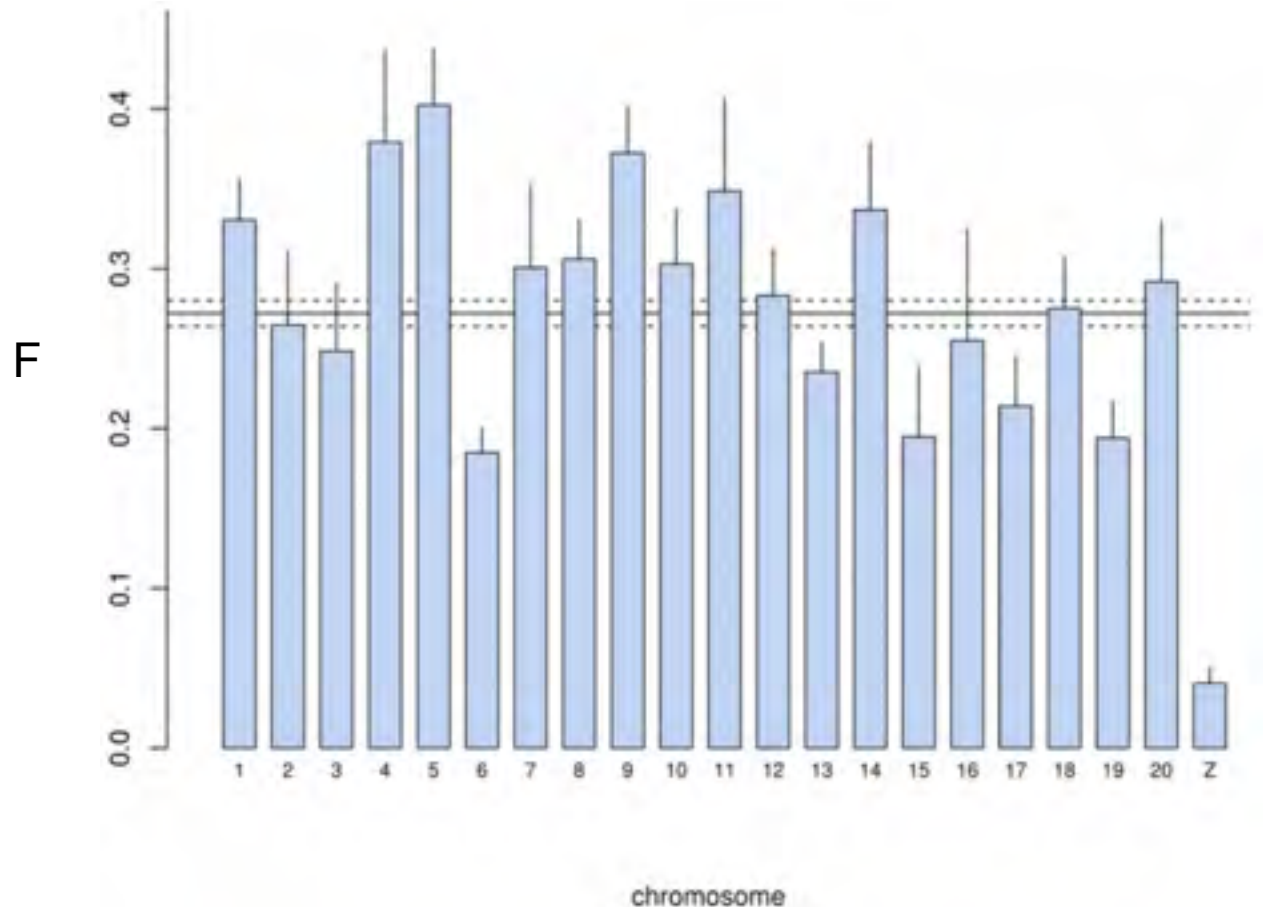


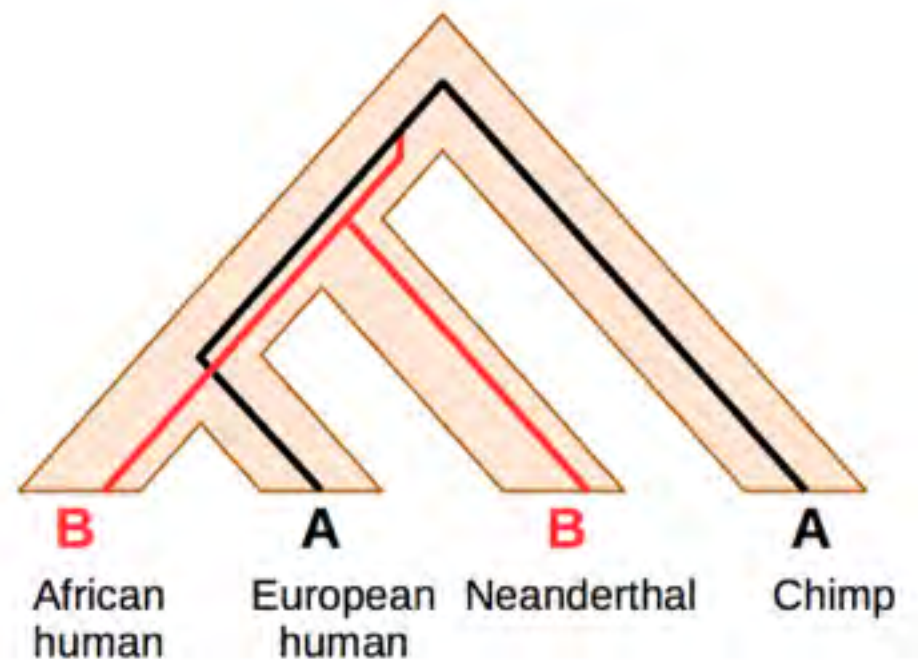
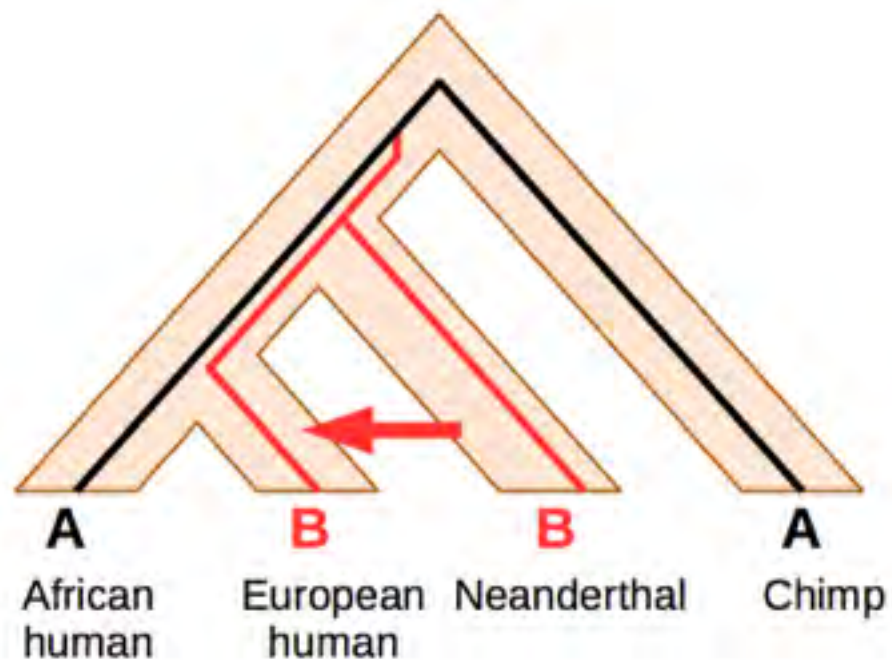
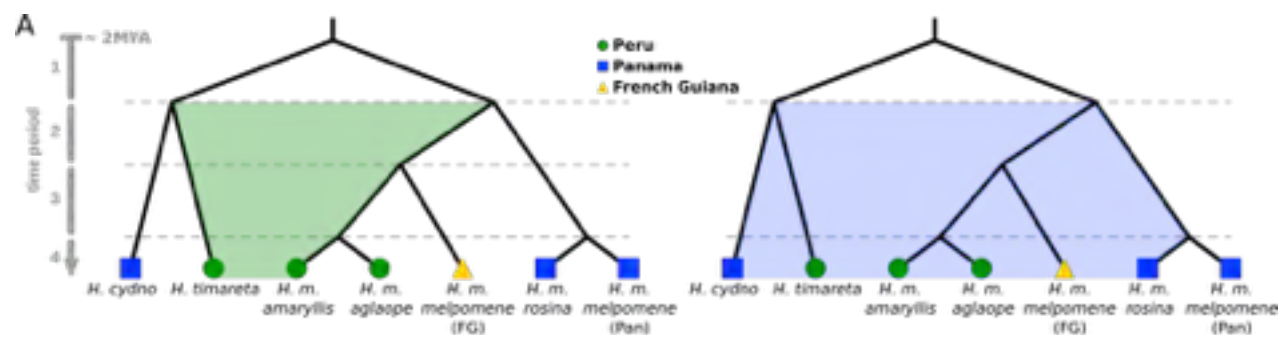
- 100 kb trees
- Only 53% group by species
- 42% group by geography!

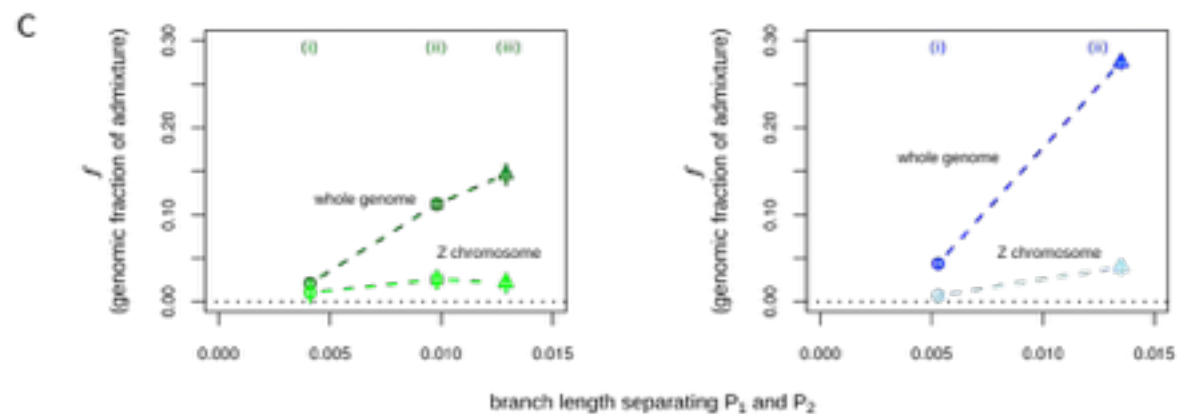
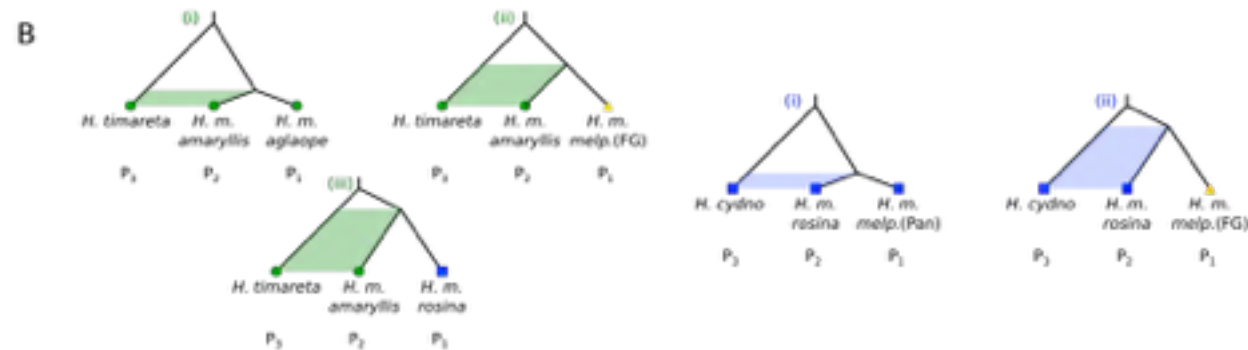
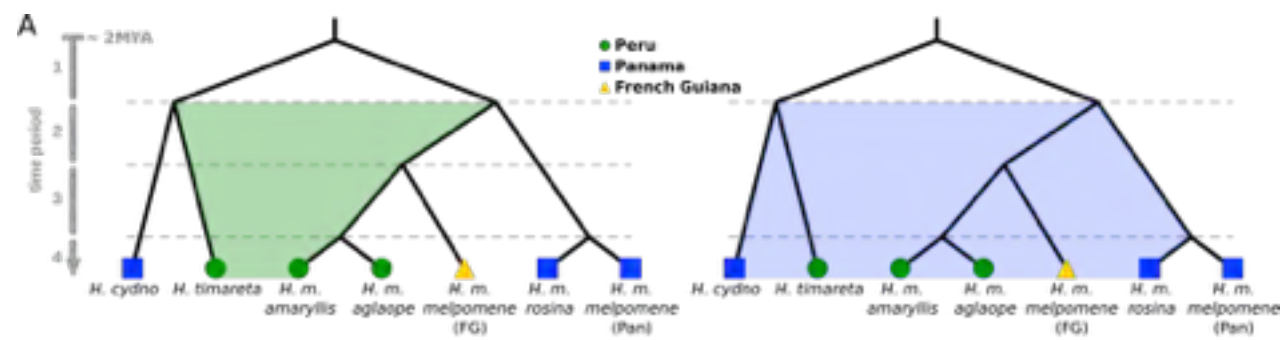


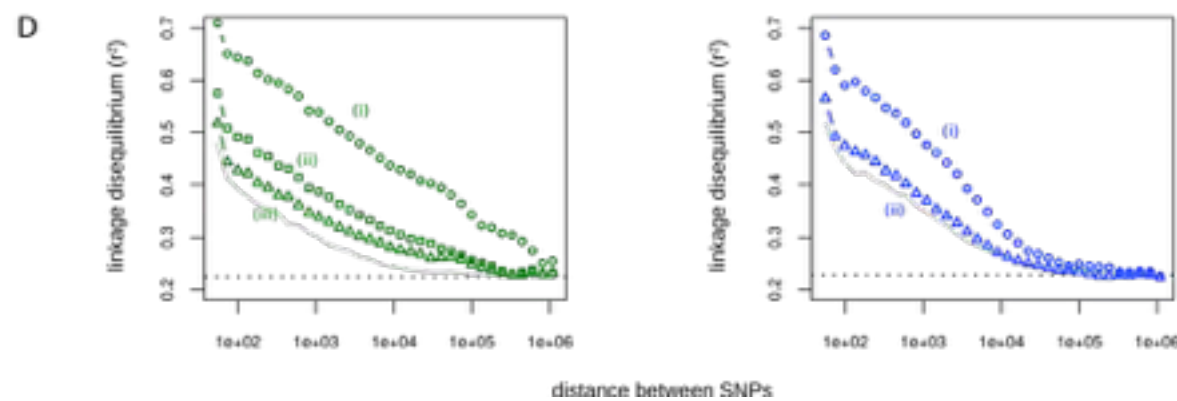
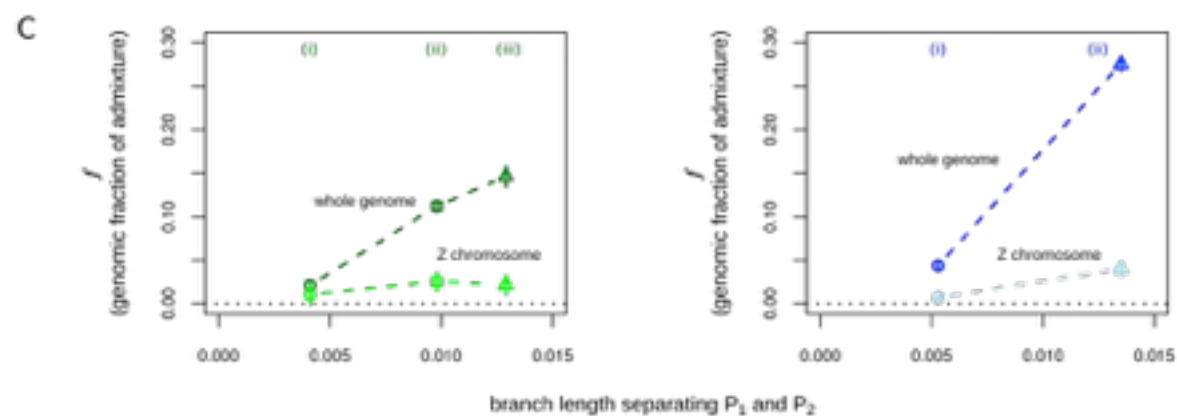
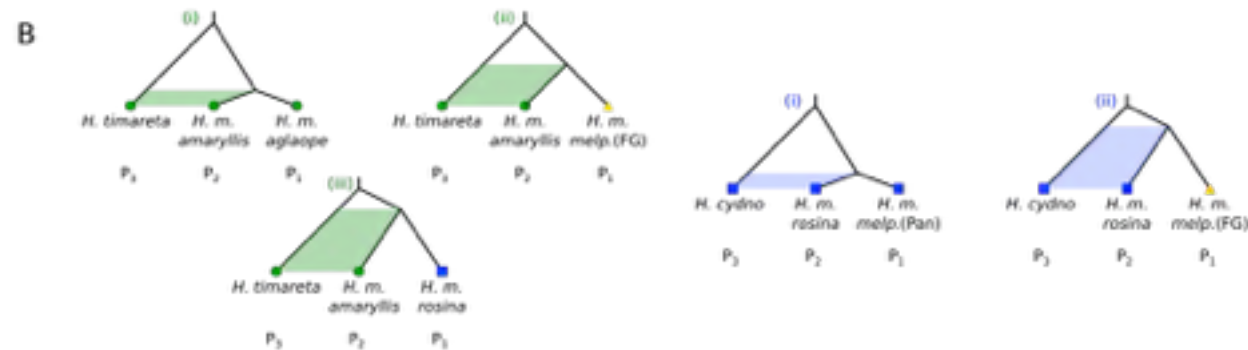
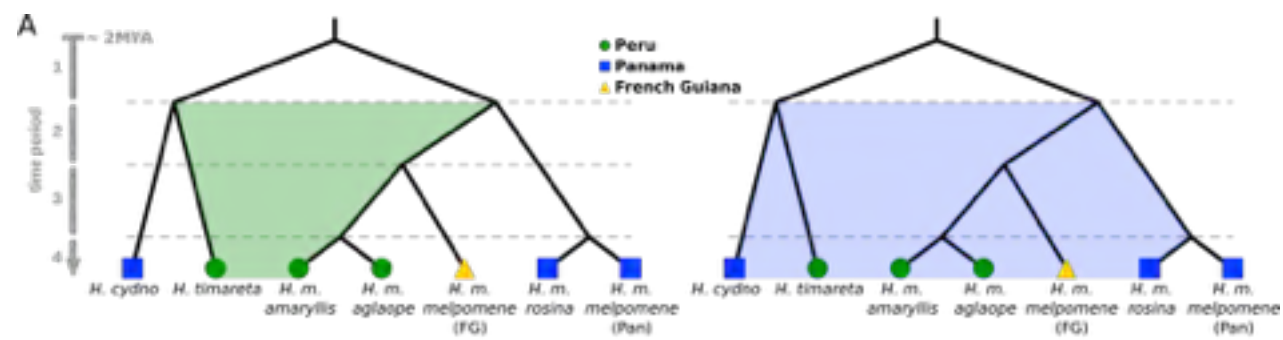
# Explicit tests for gene flow: *Heliconius* butterflies

- Much larger proportion of genome is flowing as compared to Neanderthals
- Similarly strong effect on sex chromosome

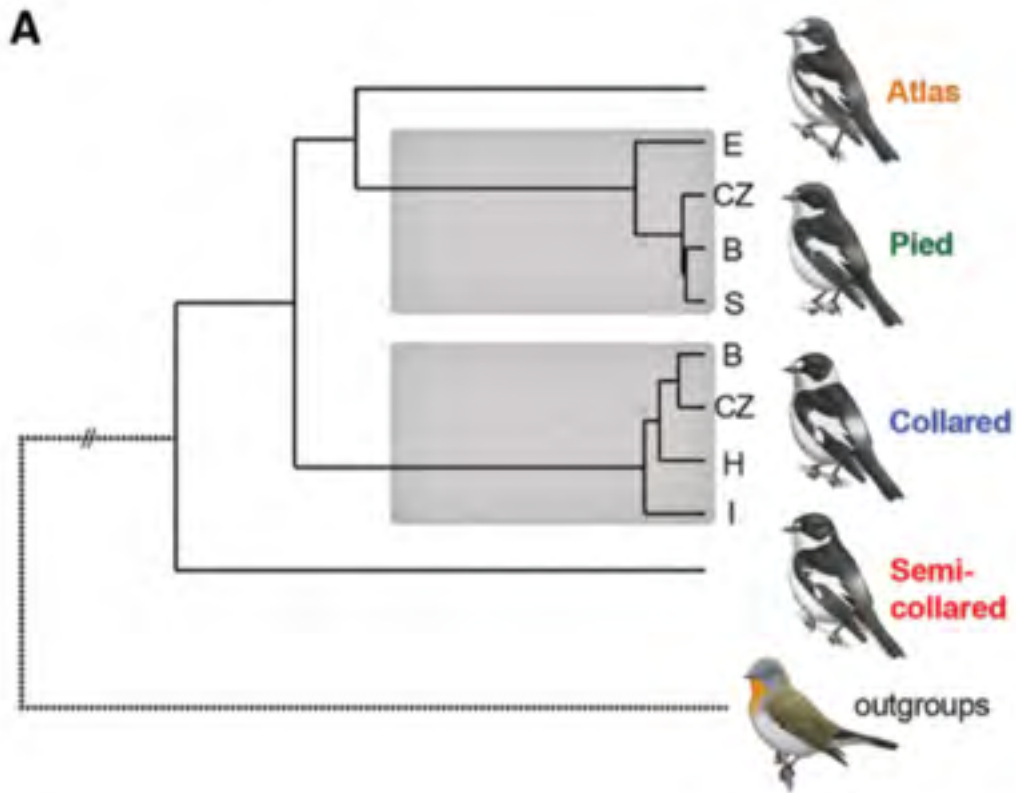






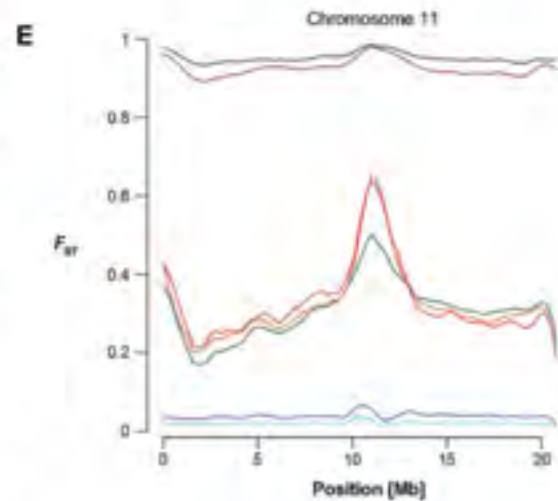
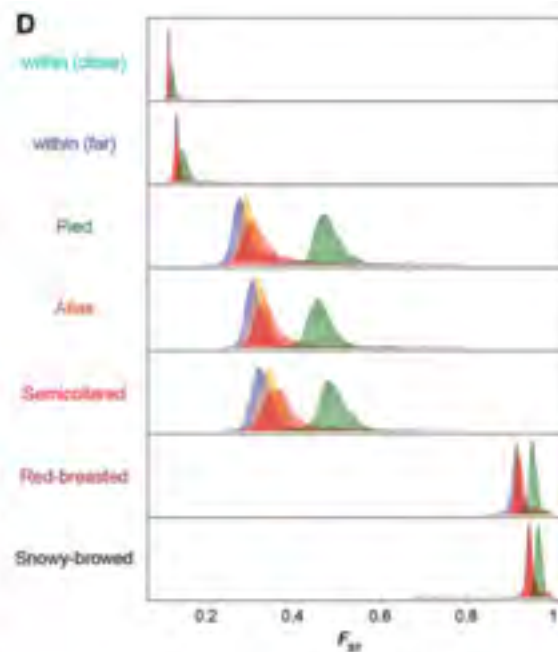
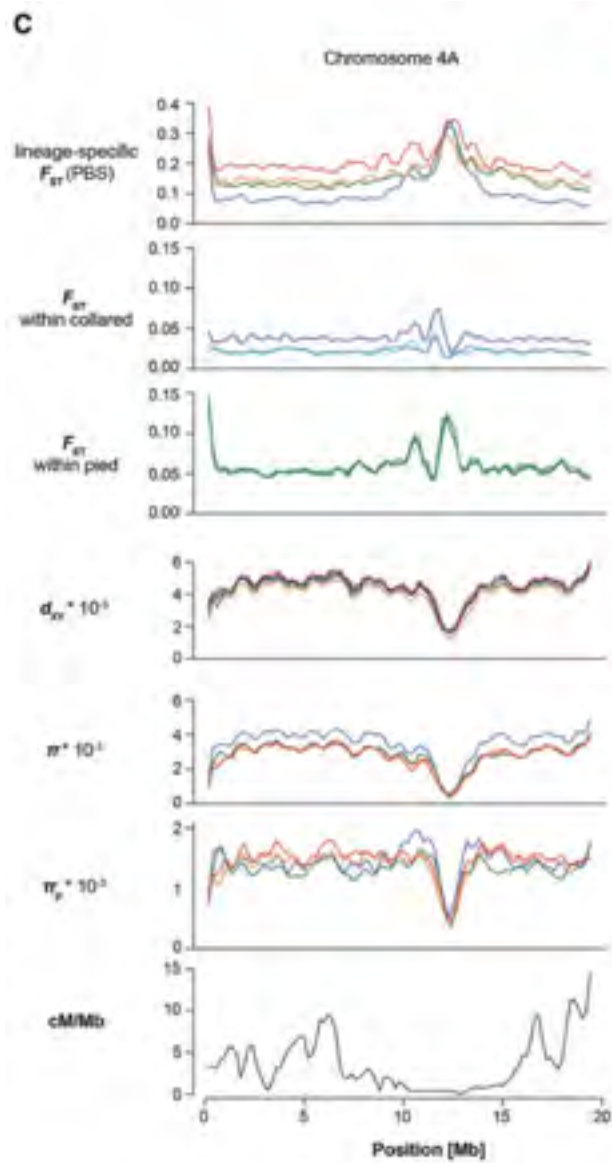
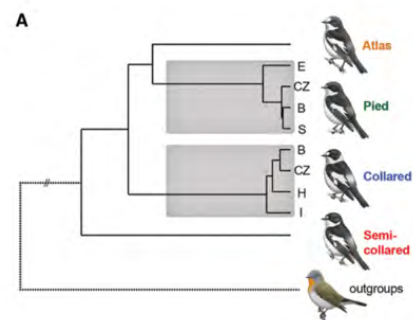






Sequenced 20 individuals per population at 20x coverage

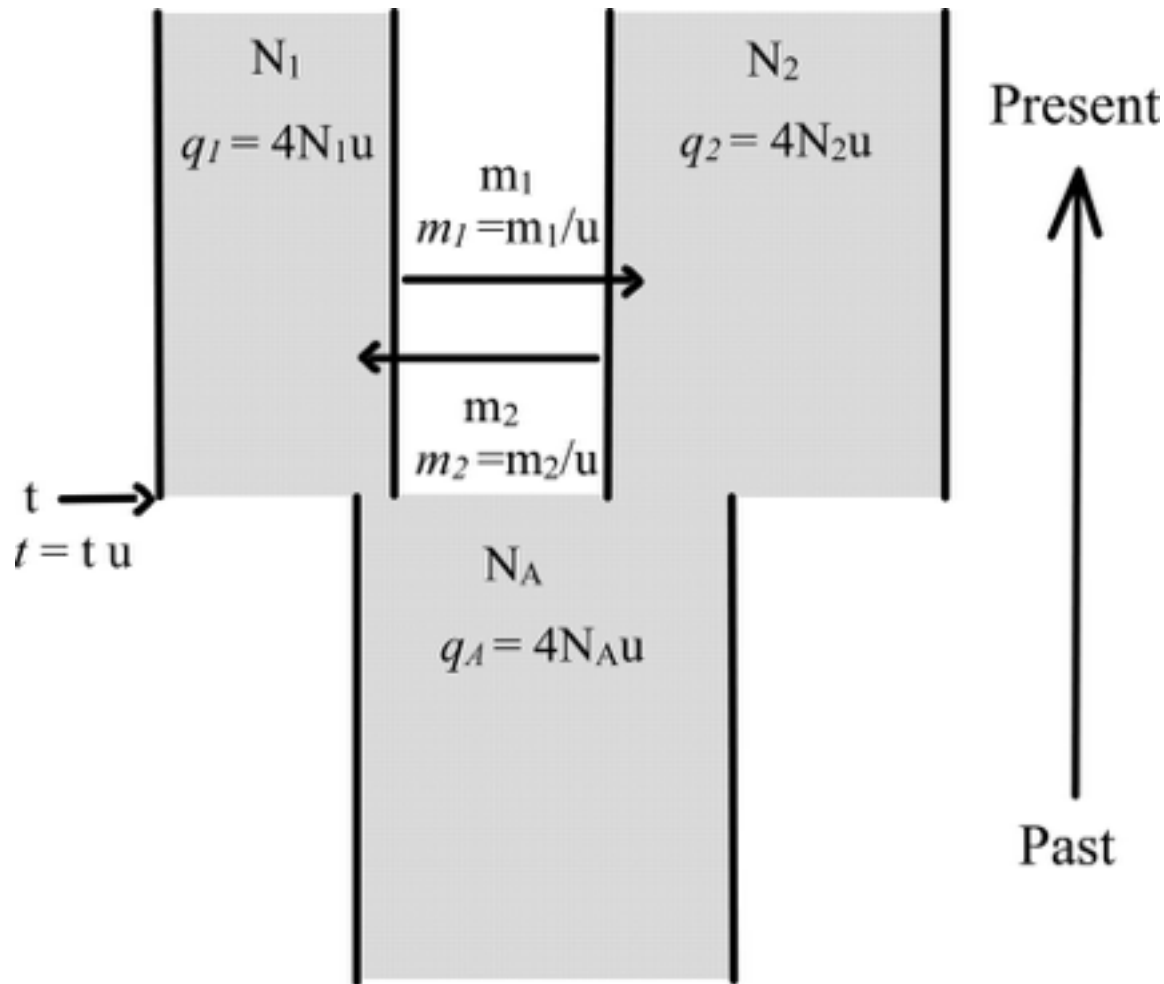




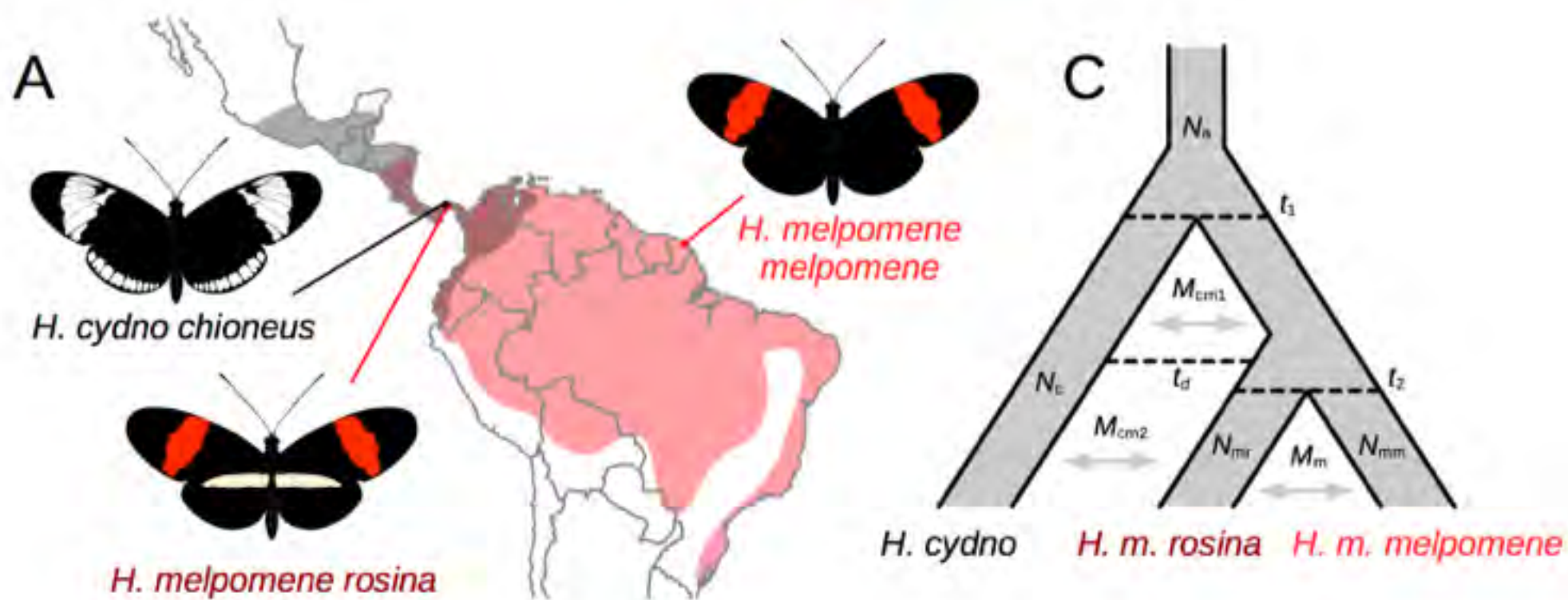
**Supplemental Table S4.** ABBA-BABA tests for gene flow. Populations/species among which the test indicates gene flow are highlighted in bold.

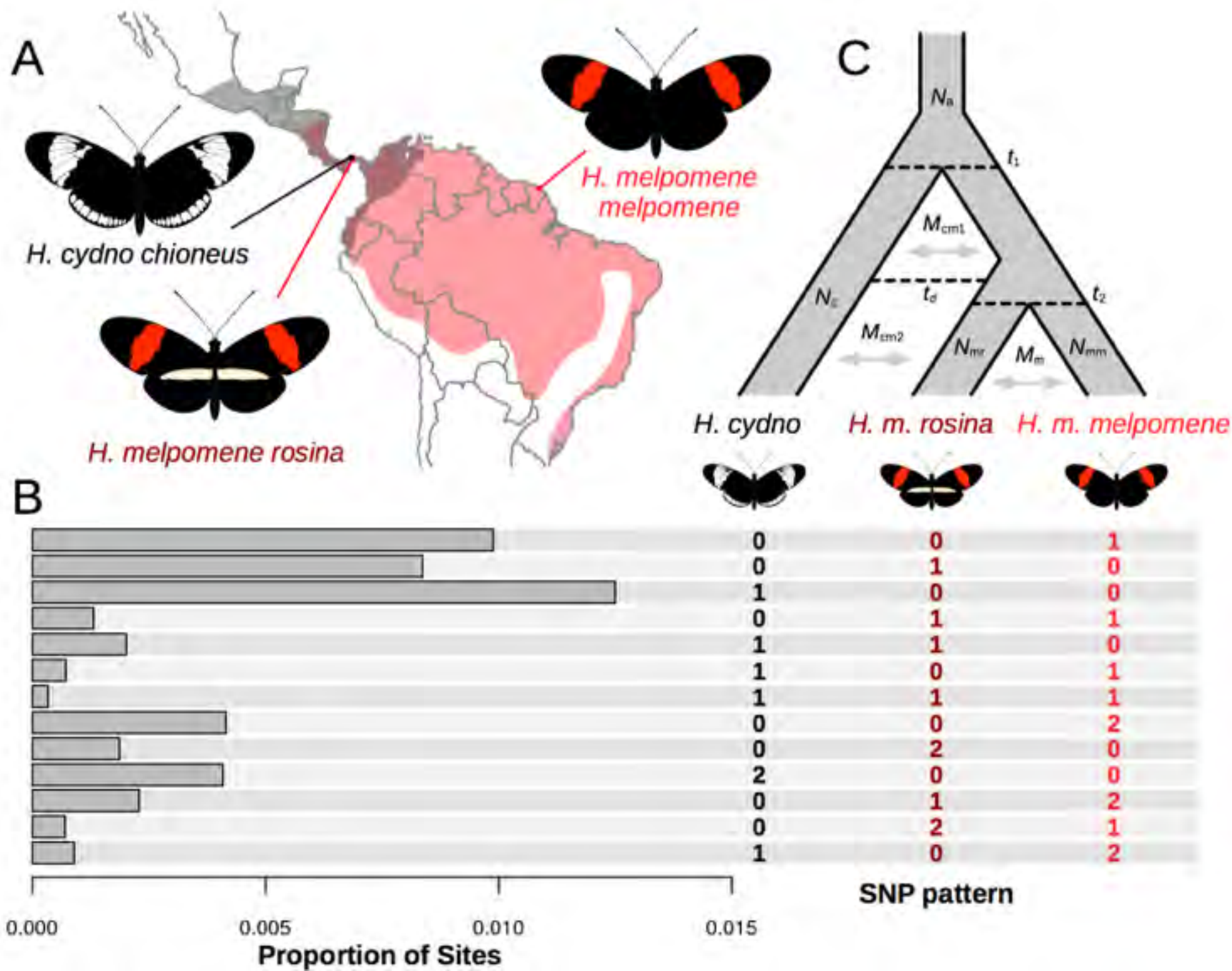
1. Inner	2. Inner	1. Outgroup	Mean(D)	SE(D)	p-value
collared Italy	collared CZ	pied CZ	0.0010	0.0010	0.3344
pied Spain	pied CZ	collared CZ	0.0004	0.0005	0.4186
<b>pied Spain</b>	Atlas	<b>collared Italy</b>	-0.1648	0.0027	$<10^{-4}$
<b>pied Spain</b>	Atlas	<b>semicollared</b>	-0.0108	0.0016	$<10^{-4}$
pied Spain	<b>collared Italy</b>	<b>semicollared</b>	0.1162	0.0018	$<10^{-4}$
Atlas	<b>collared Italy</b>	<b>semicollared</b>	0.1242	0.0016	$<10^{-4}$

# An alternative is to take an explicit modelling approach



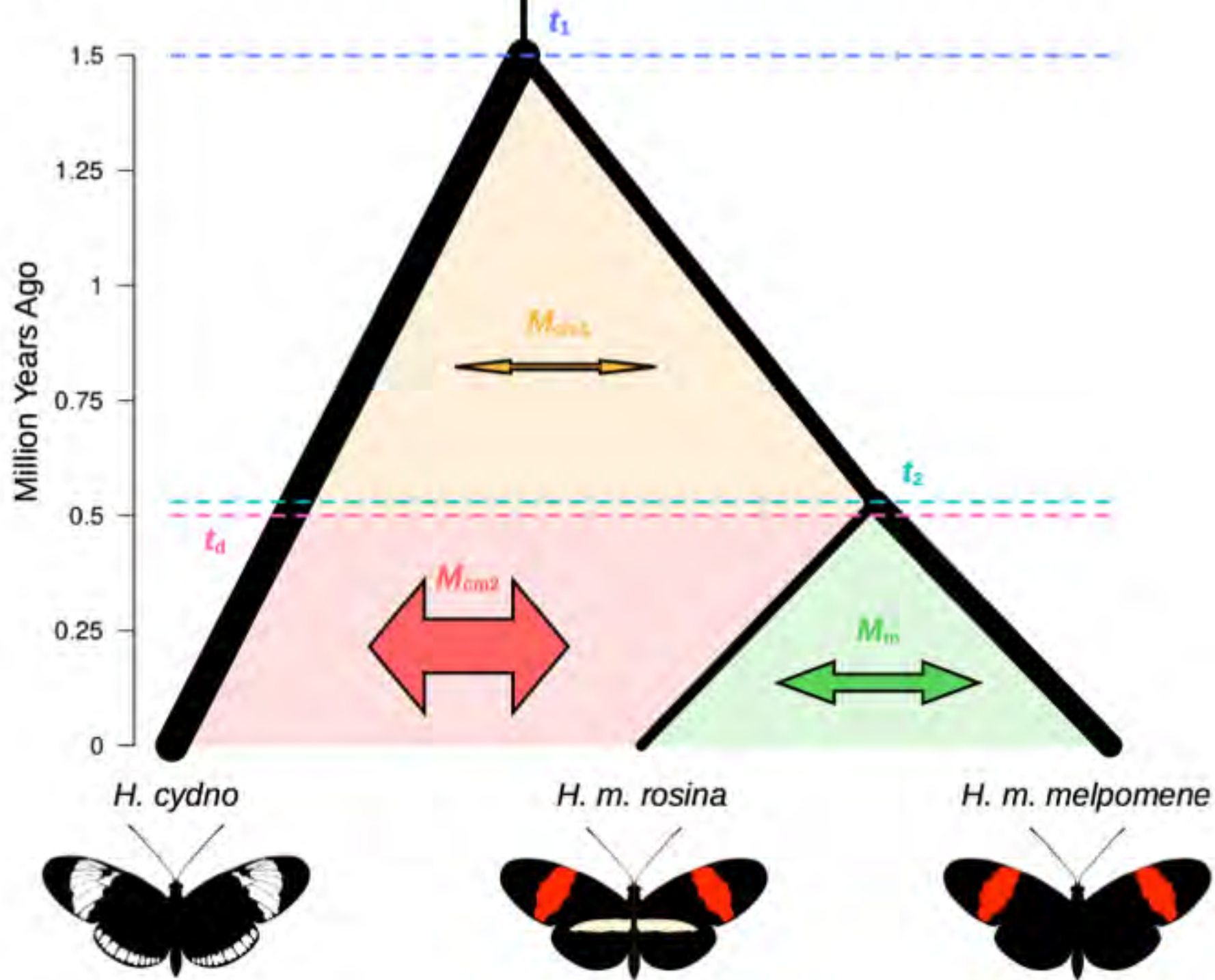
IM and IMa  
Jody Hey



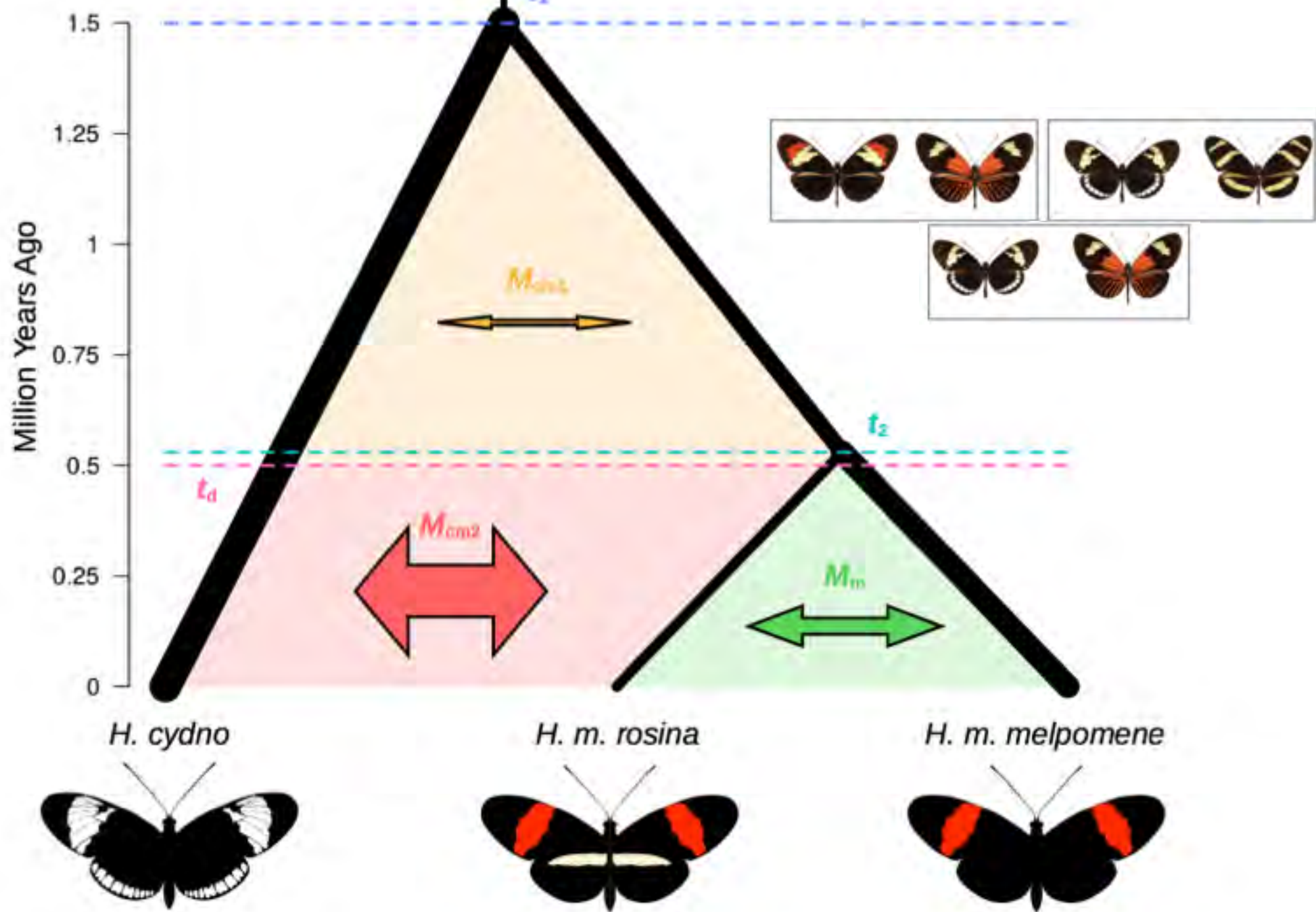




B



B



So far models have mostly just estimated genome-wide parameters...assuming the genome is homogenous

Where we need to go next is to incorporate genome heterogeneity in selection and recombination

So far models have mostly just estimated genome-wide parameters...assuming the genome is homogenous

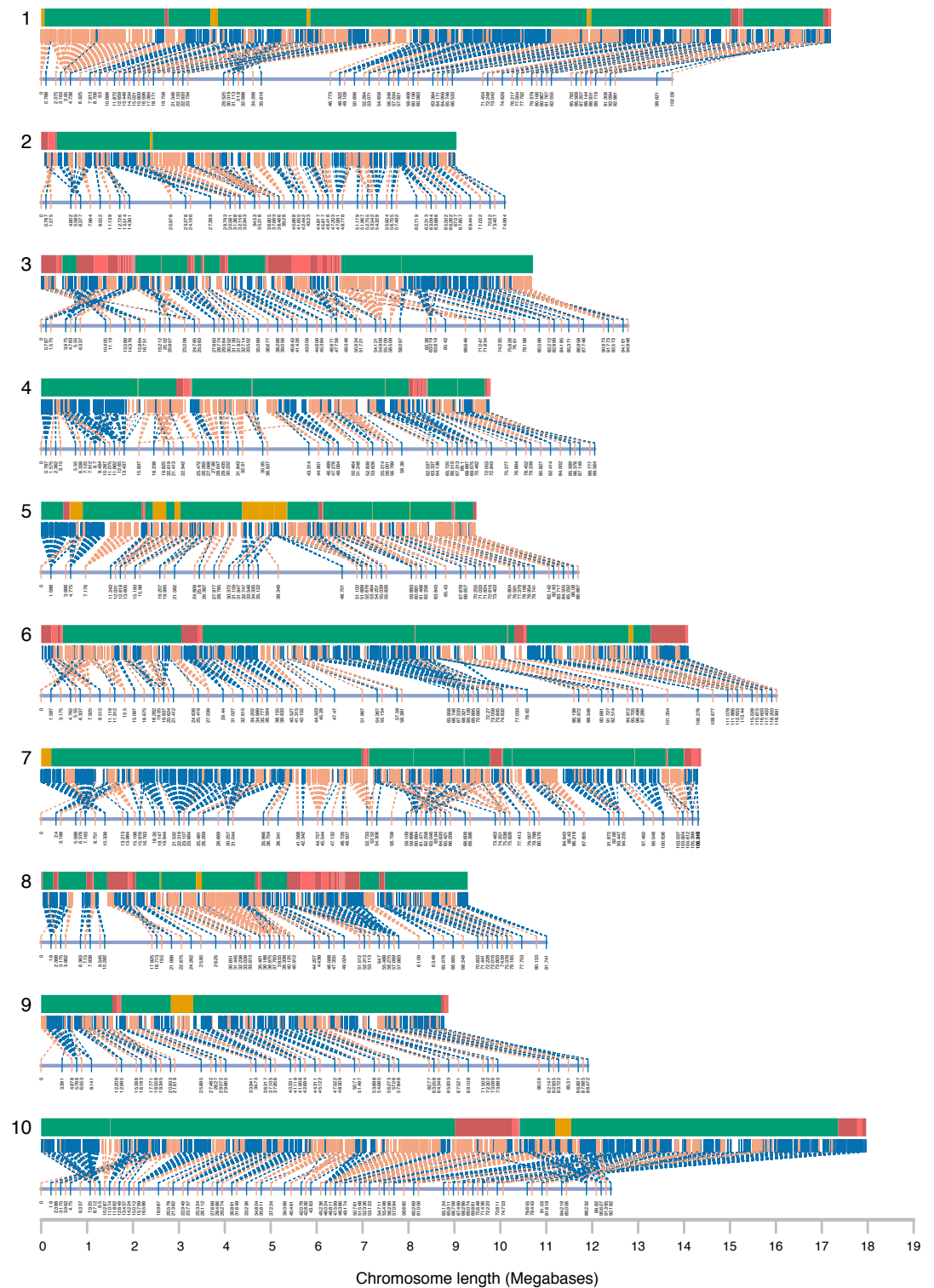
Where we need to go next is to incorporate genome heterogeneity in selection and recombination

## **Identifying Loci Under Selection Against Gene Flow in Isolation-with-Migration Models**

**Vitor C. Sousa,<sup>\*,1,2</sup> Miguel Carneiro,<sup>†</sup> Nuno Ferrand,<sup>†</sup> and Jody Hey<sup>\*,1</sup>**

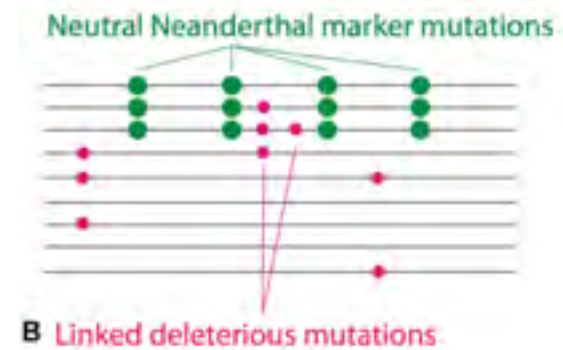
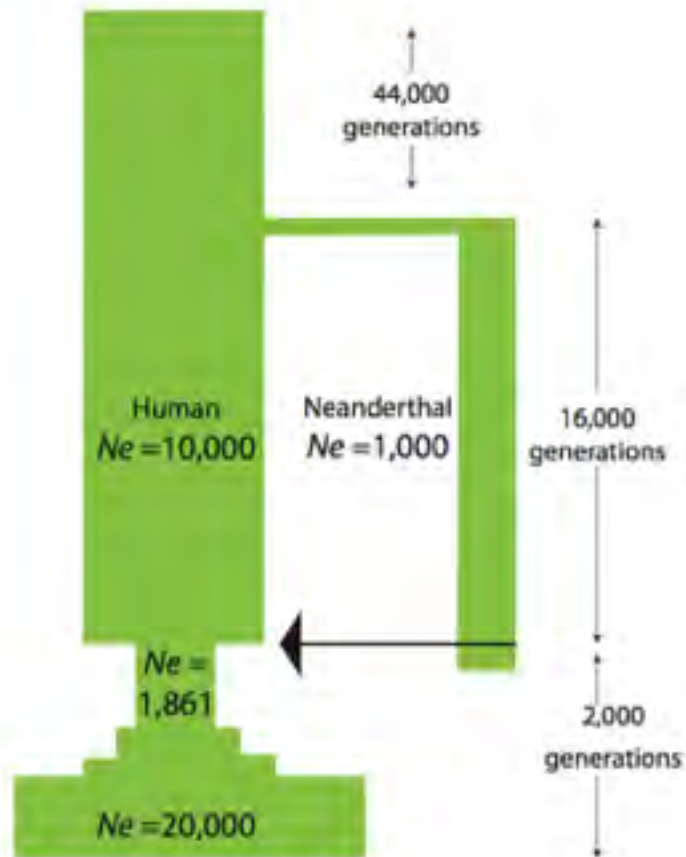
<sup>\*</sup>Department of Genetics, Rutgers, The State University of New Jersey, Piscataway, New Jersey 08854, and <sup>†</sup>CIBIO, Centro de Investigação em Biodiversidade e Recursos Genéticos, Universidade do Porto, 4099-002 Porto, Portugal

# High density linkage maps to map the recombination landscape

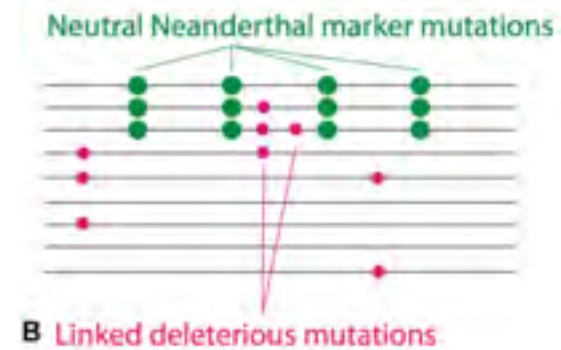
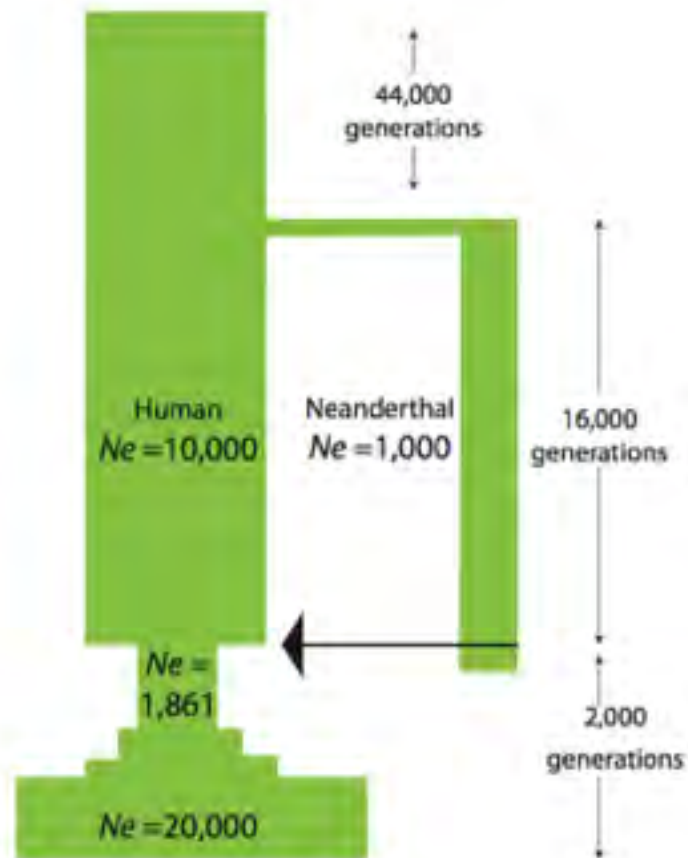




# The effect of background selection on introgression in humans



# The effect of background selection on introgression in humans



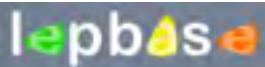
Admixture is less in gene rich regions supporting this model.....

# Population and speciation genomics: Conclusions

- Great power to detect subtle signals of selection and gene flow
- Can make more general observations about genes and regions involved in adaptation
- BUT genomic processes complicate the picture
- Best approaches combine multiple signals to infer process
- Eventually we need to combine background selection, recombination, positive selection

And finally a shameless plug....

# And finally a shameless plug....



BLAST | Downloads | WebApollo | Blog | Source code

Search LepBase...

## the Lepidopteran genome database

The Lepidoptera comprises over 170,000 species, including major agricultural pests, important plant pollinators and the first domesticated insect. The Lepidoptera have played a pivotal role in the development of ecological and evolutionary biology and includes 'model' organisms for a variety of disciplines, including conservation biology, theoretical ecology, systematics, developmental biology, genetics and evolutionary theory.

As research questions in the Lepidoptera are increasingly being approached using genomic data, Lepbase offers a platform that integrates these data, focusing on the specific needs of the Lepidopteran research community to open up this diverse clade to comparative analysis.

### Available genomes



### Heliconline DISCOVAR assemblies



## What's new

This is version 1.0 of the Lepbase ensembl genome browser. New features include a dedicated BLAST server, Lepidoptera-specific orthologue predictions & gene trees, and WebApollo for community annotation. If there is something missing that you would like to see then please [contact us](#).

### New species/assemblies in version 1.0:

- *Chilo suppressalis* CsuOGS1.0
- *Heliconius melpomene* Hmel2
- *Lerema accius* v1.1
- *Manduca sexta* Msex\_1.0
- *Papilio glaucus* v1.1
- *Plodia interpunctella* v1
- 18 Heliconine DISCOVAR assemblies

## More from Lepbase...

We aim to provide a comparative genomics resource for the Lepidoptera research community, with BLAST and WebApollo servers in addition to this Ensembl instance, visit [leabase.org](http://leabase.org) or follow @leabase to find out more about the project.

## Coming soon

- BioMart
- Rfam annotations
- Variations
- Whole genome alignments
- *Bicyclus anynana* v1.0

LepBase is funded by a BBSRC Bioinformatics and Biological Resources fund award (BB/K020181/1, BB/K019945/1, BB/K020129/1) to Prof. Mark Blaxter (University of Edinburgh), Prof. Chris Jiggins (University of Cambridge), Dr. Karancho Dasrathapada (University of York) and maintained by two post-doctoral bioinformaticians, Dr. Richard Challin and Dr. Sujal Kumar, based in the Blaxter lab. Reuben Nowell, another member of the Blaxter lab, also contributes to LepBase as part of his involvement in the [Bicyclus anynana](#) genome project.



THE UNIVERSITY  
of EDINBURGH



UNIVERSITY of York



## Contact us

We want to work with the Lepidoptera research community to build Lepbase into a genuinely useful resource. If you have more data that you would like to see included or want advice on how to use Lepbase in your research, please [contact us](#).



# Adaptive introgression













Photo credit  
Andrei Sourakov







Photo credit  
Andrei Sourakov

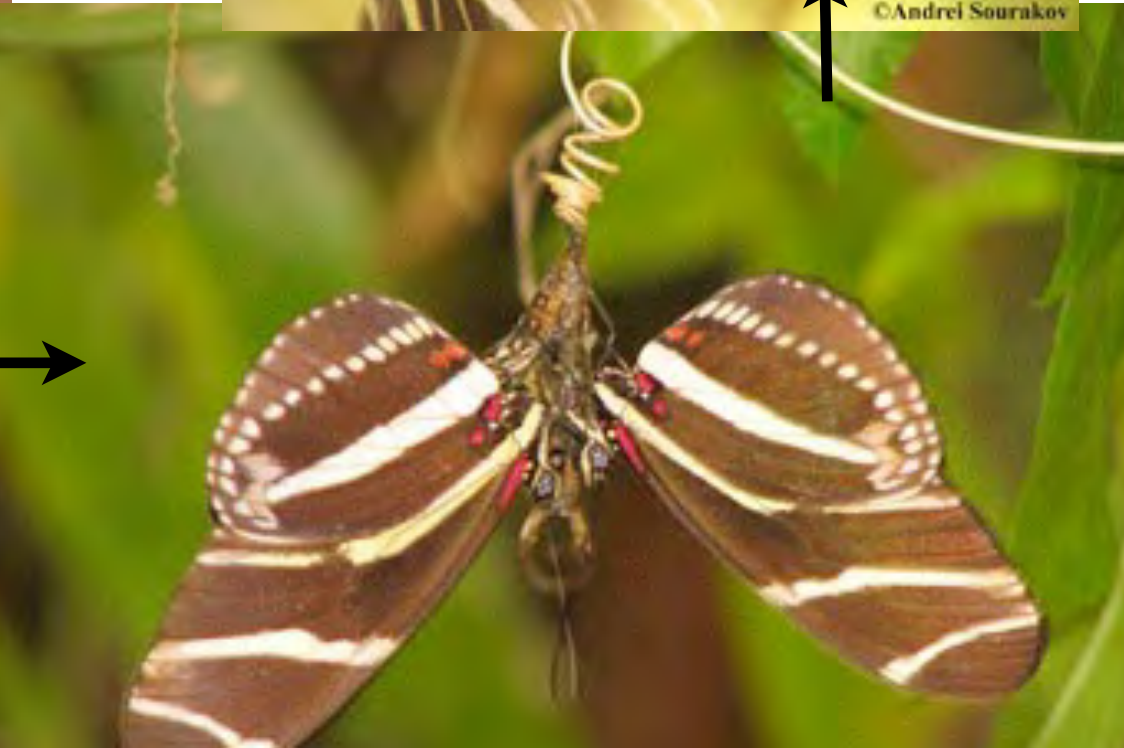




Photo credit  
Andrei Sourakov





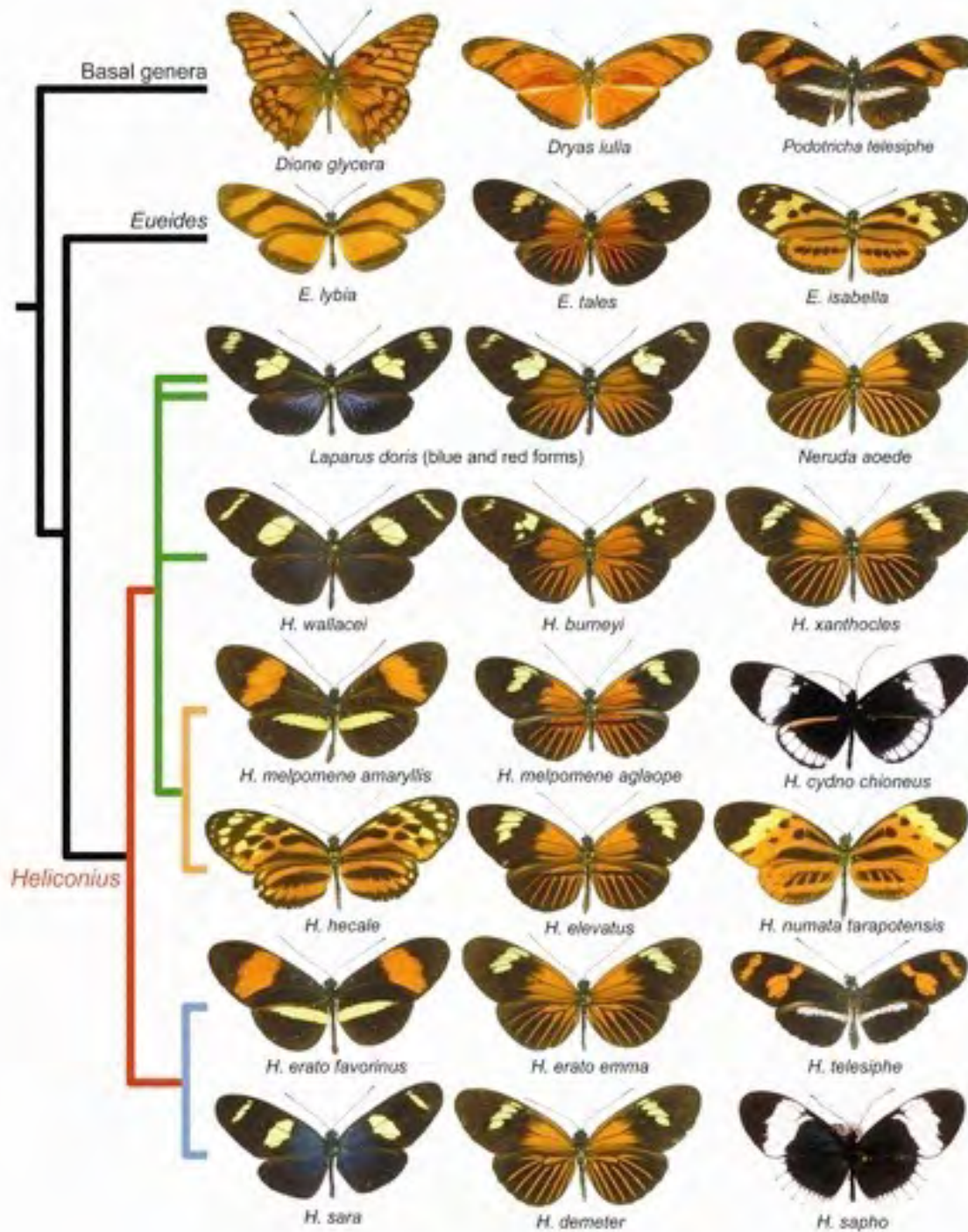












43 species

77 species



Fritz Müller

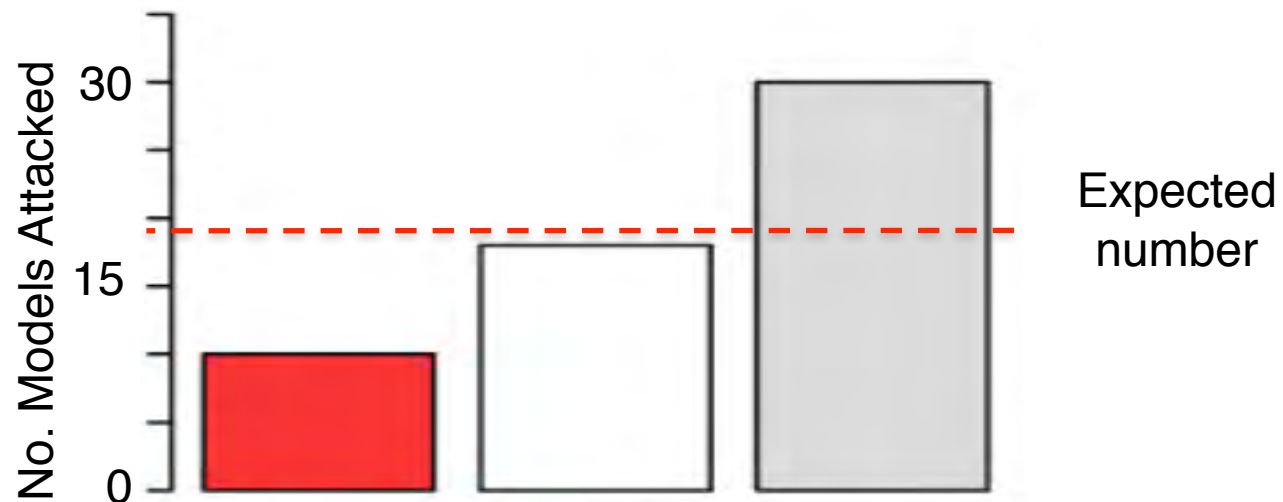


© Richard G.

Rufous-tailed Jacamar



*H. melpomene*   *H. cydno*   F1 hybrid



– G-test:  $G = 7.25$ , d.f. = 1,  $p = 0.007$

Merrill et al., Proc. Roy. Soc 2012





102

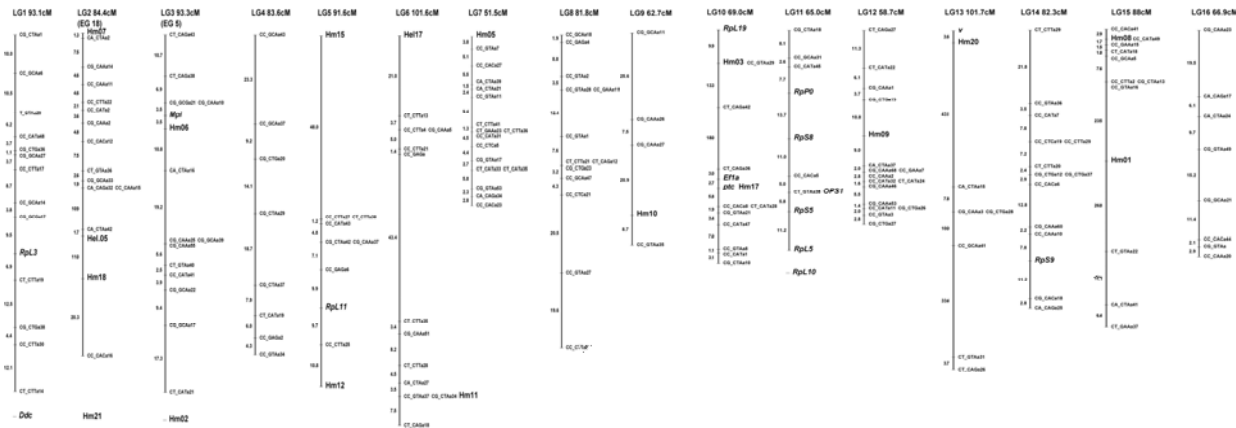


204



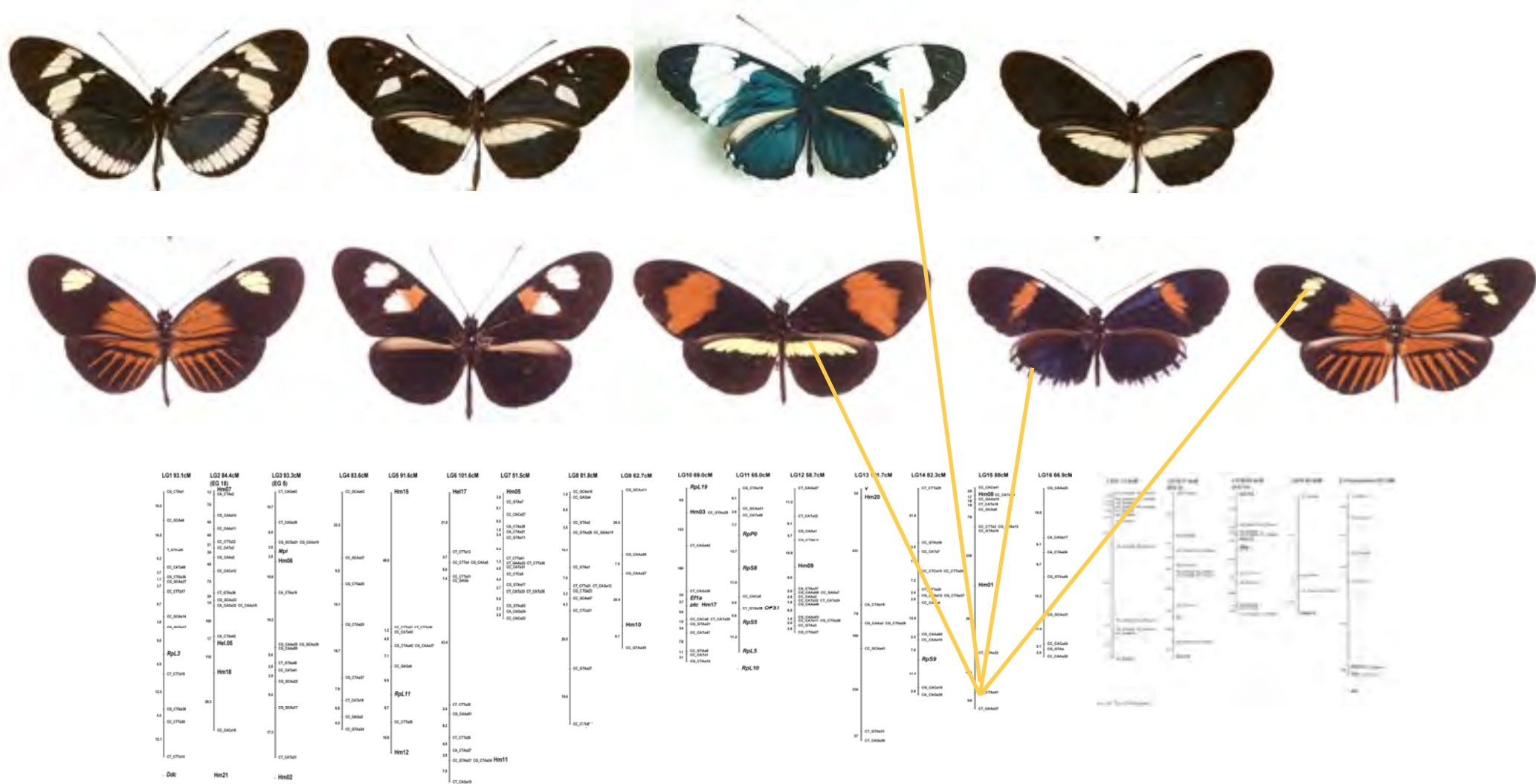
76

# Several major loci control *Heliconius* patterns

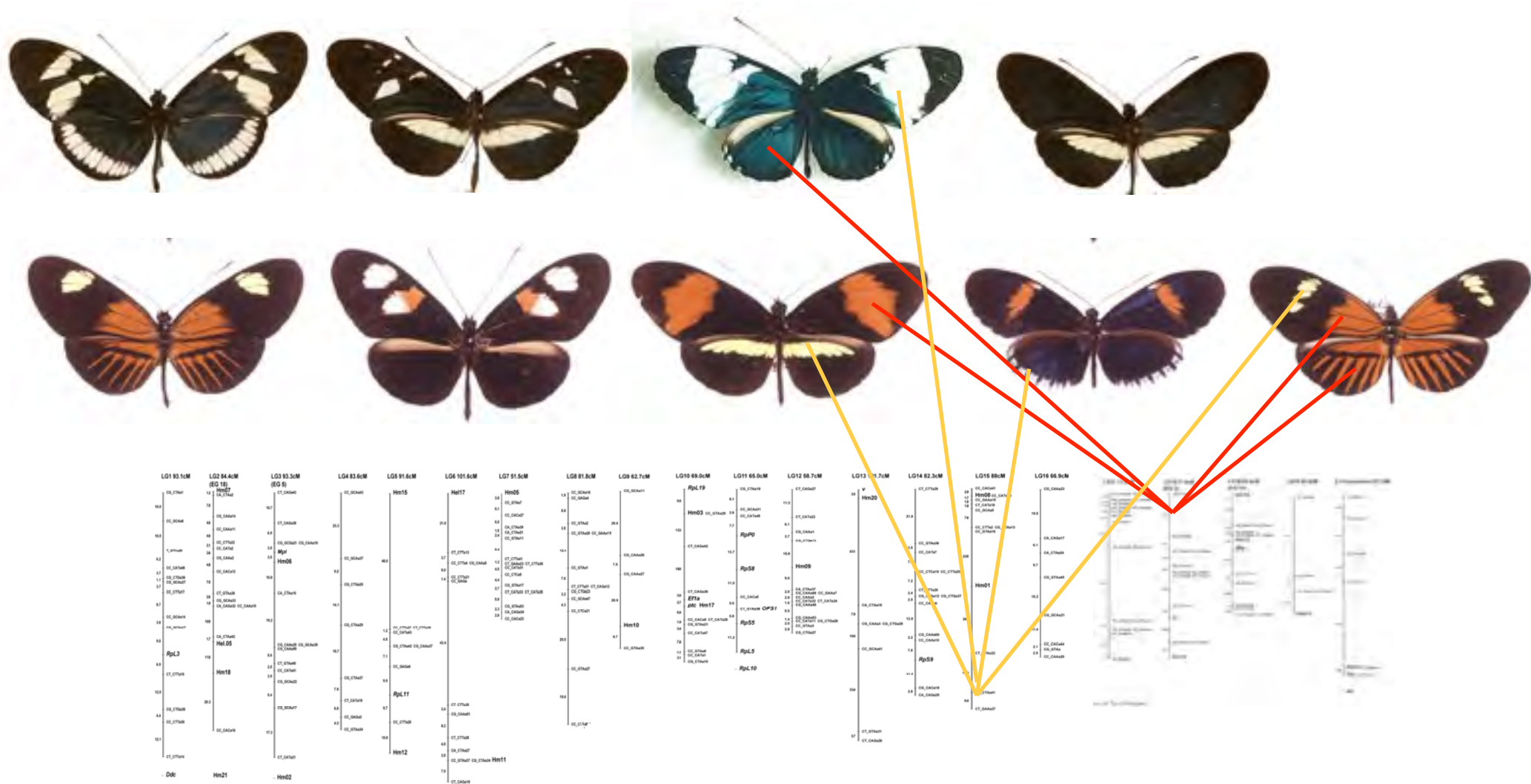




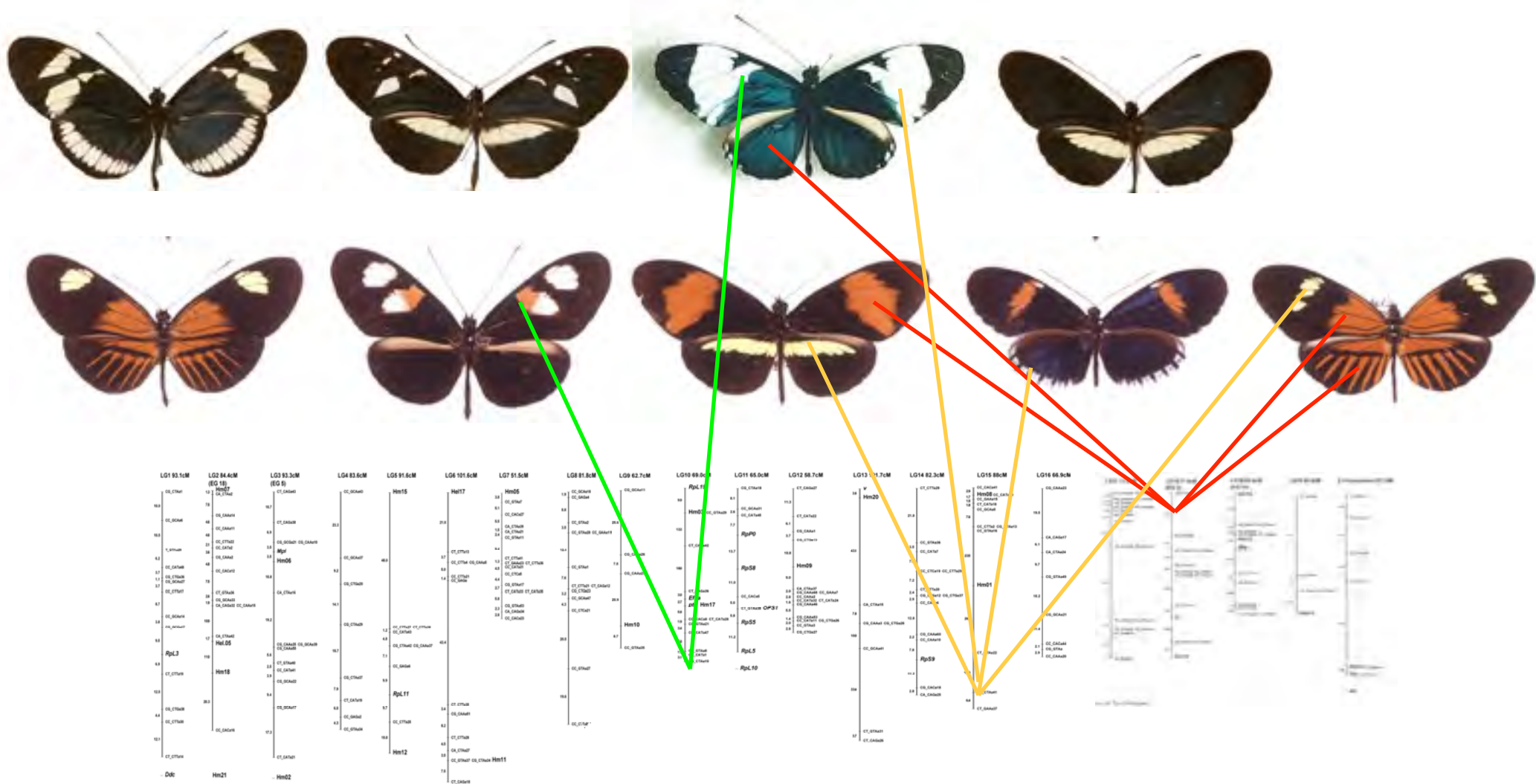
# Several major loci control *Heliconius* patterns



# Several major loci control *Heliconius* patterns

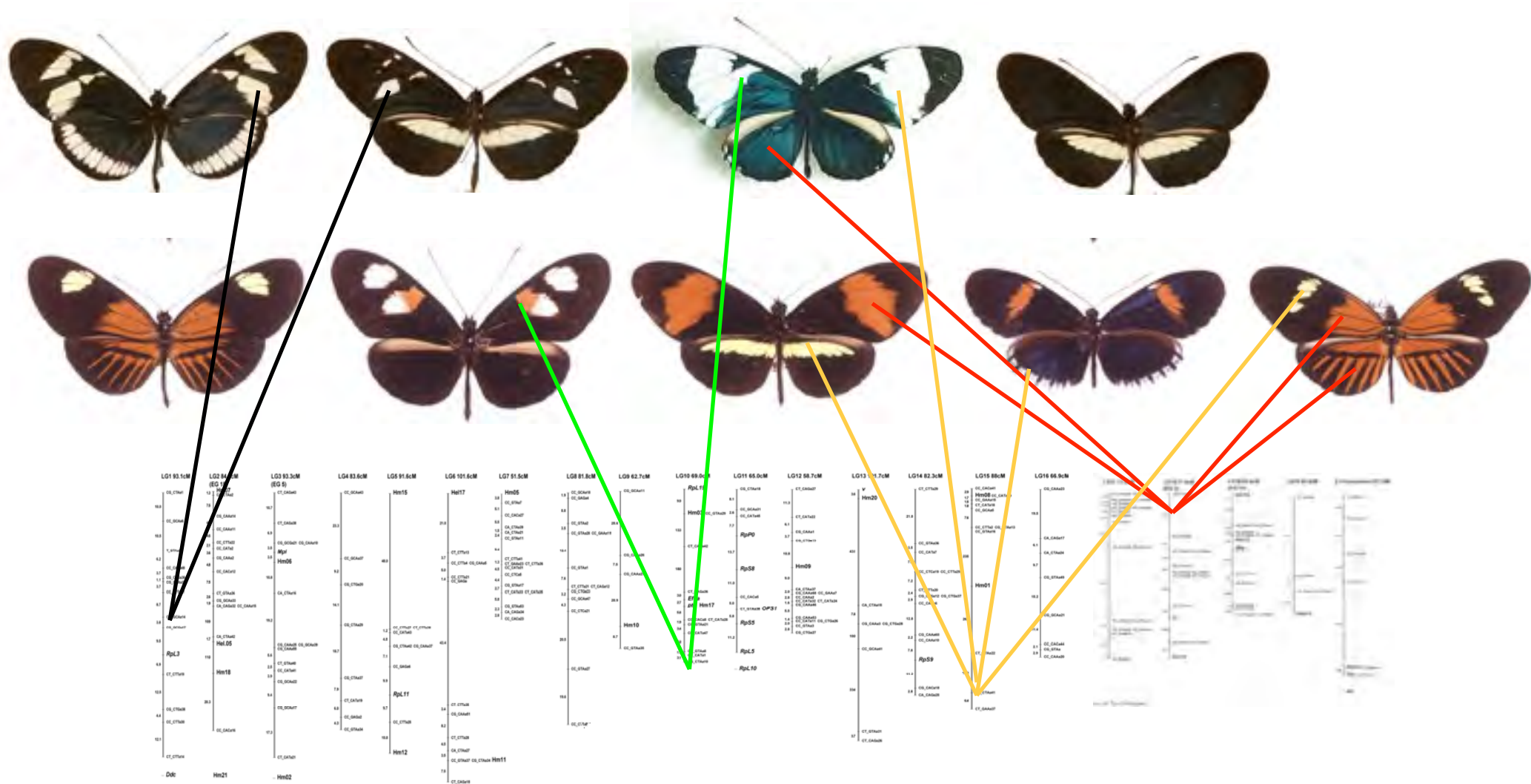


# Several major loci control *Heliconius* patterns



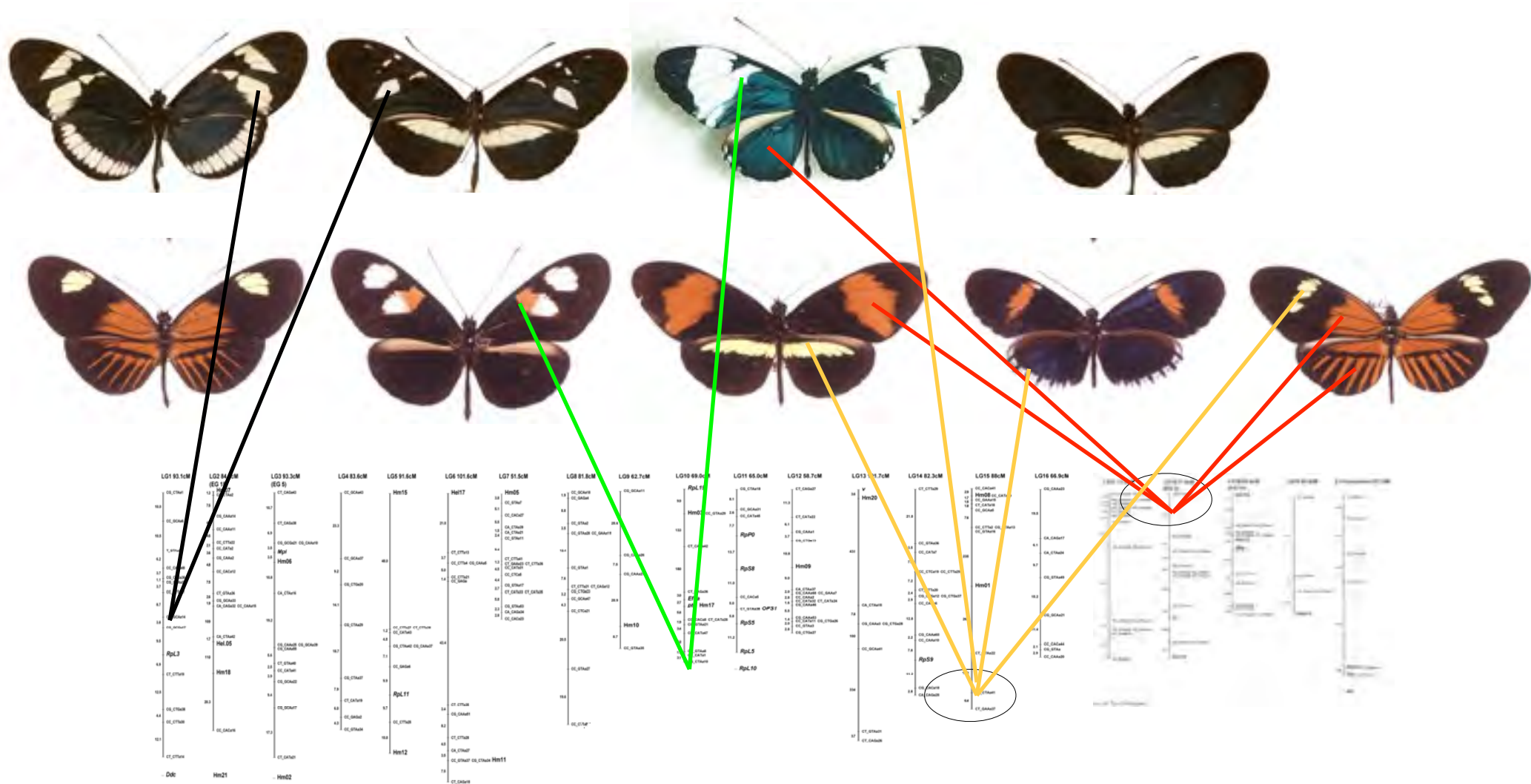


# Several major loci control *Heliconius* patterns



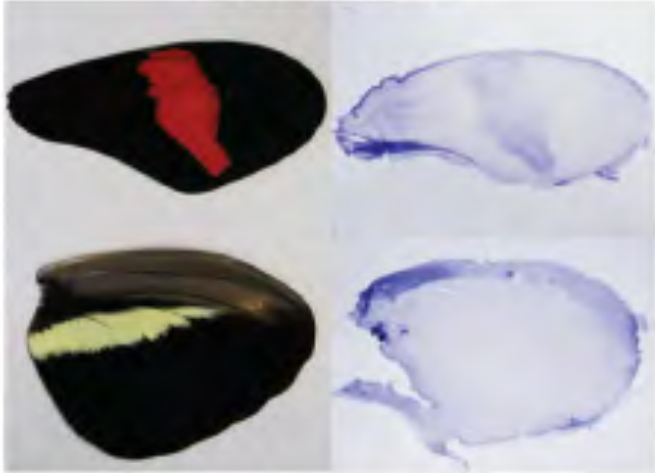


# Several major loci control *Heliconius* patterns



B

*H. erato petiverana*



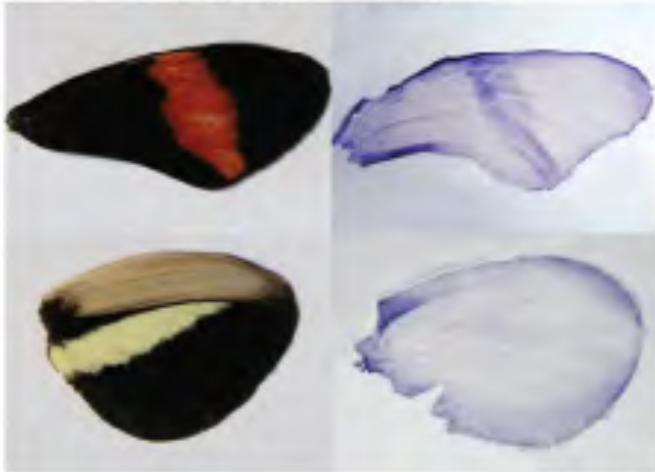
*H. erato erato*



*H. cydno galanthus*



*H. melpomene rosina*



*H. melpomene malleti*



*H. melpomene plesseni*



Reed et al., 2011 Science

E

Homothorax  
DAPI

Optix  
DAPI

*H. elevatus*

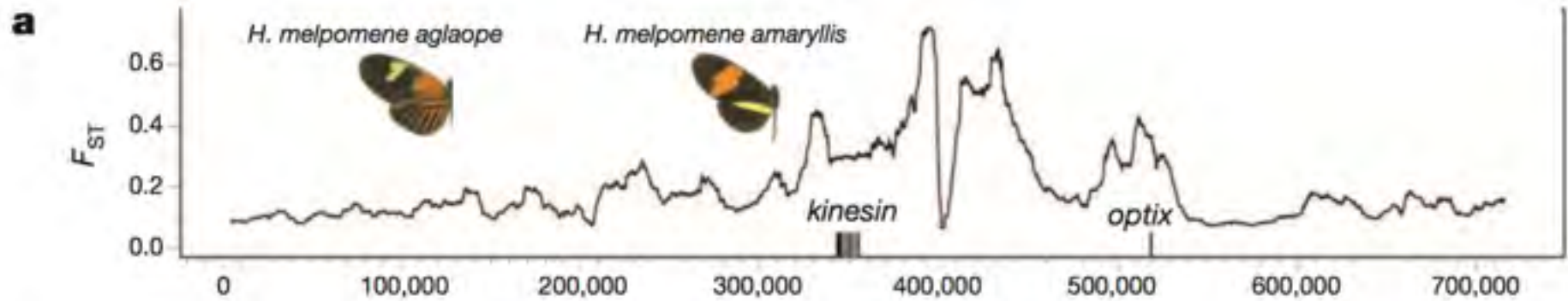
1000µm

*H. m rosina*



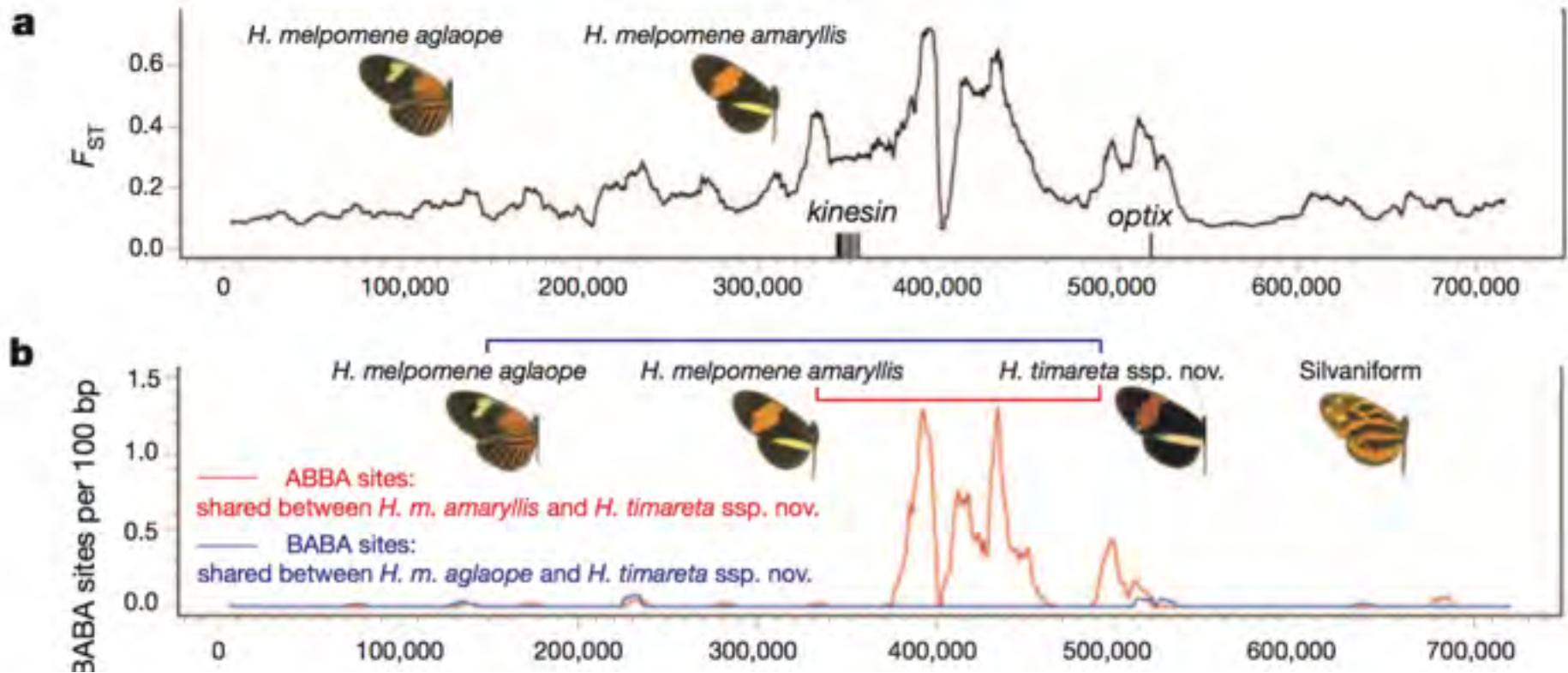


# Cross-species sharing

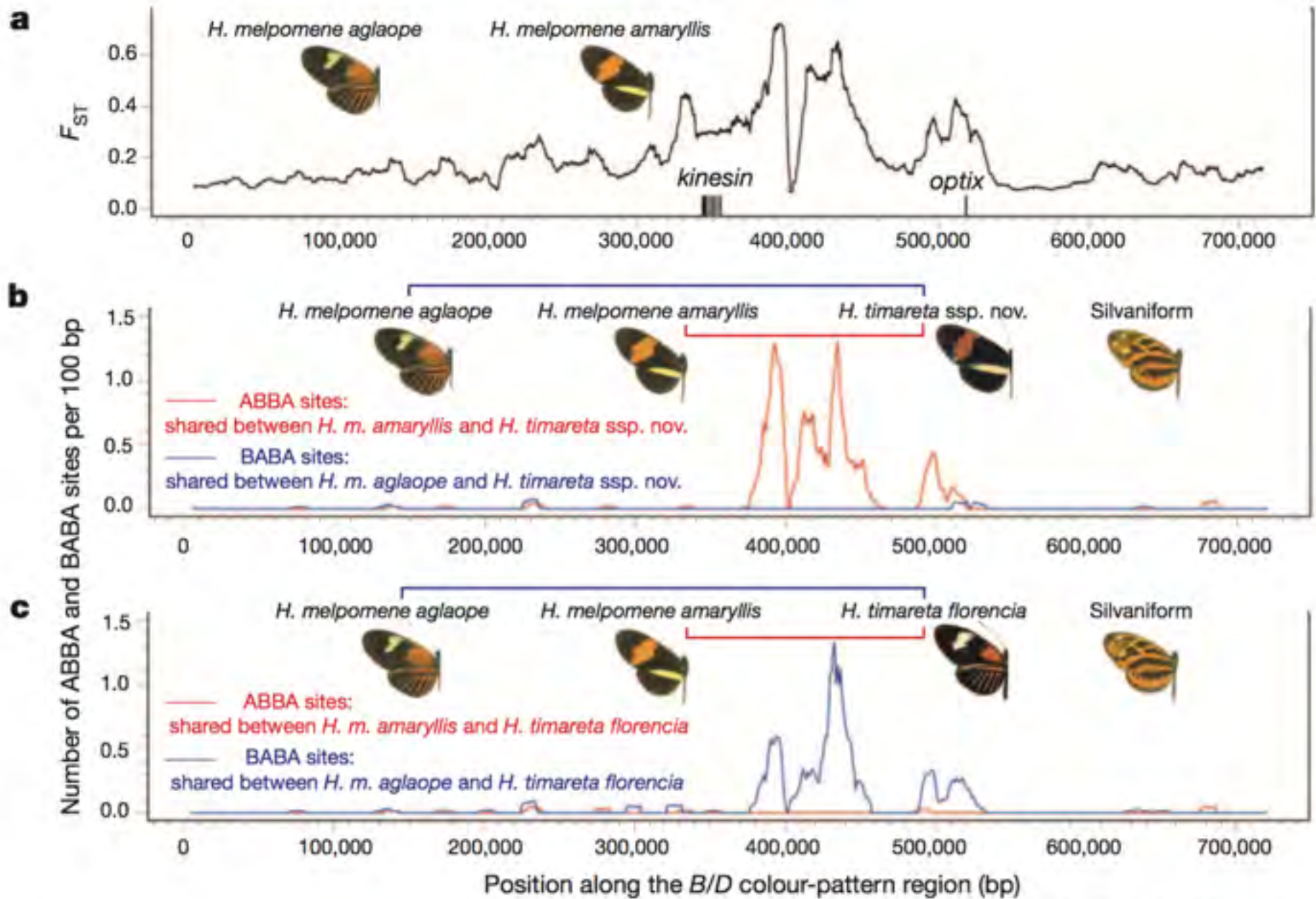




# Cross-species sharing

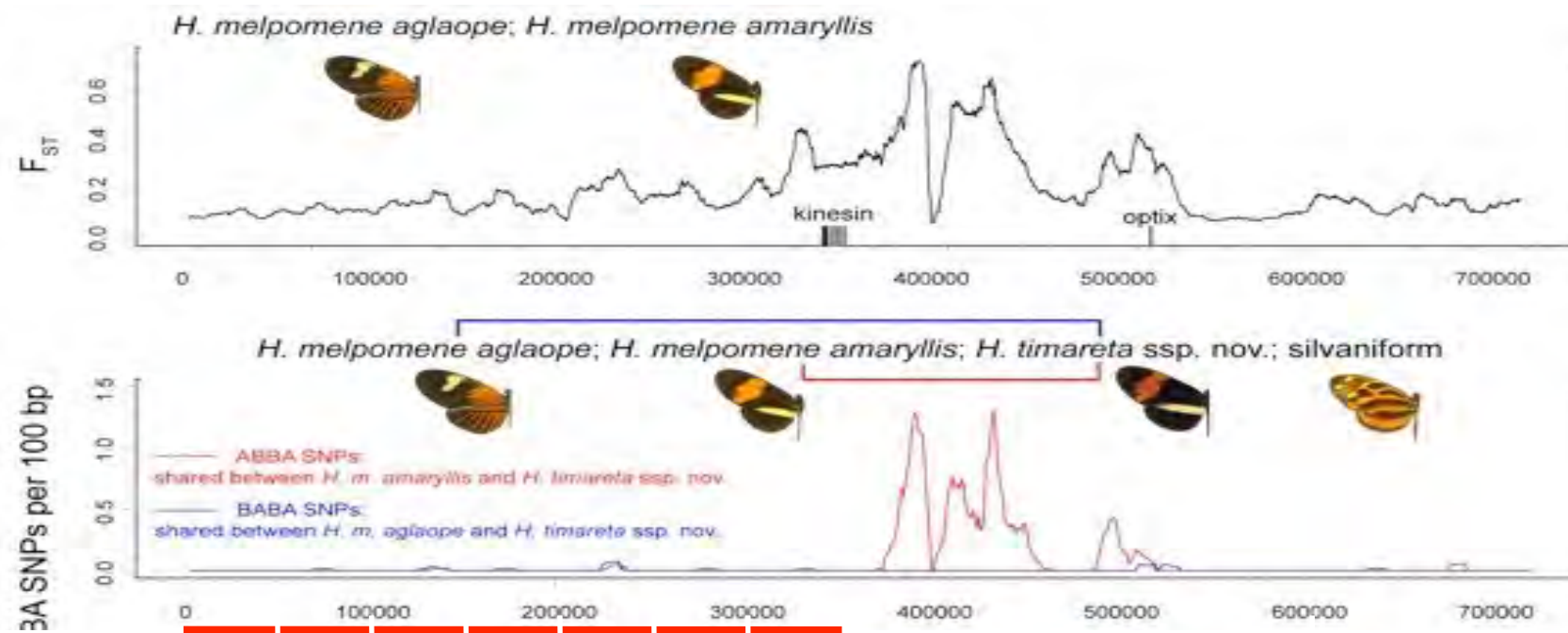


# Cross-species sharing



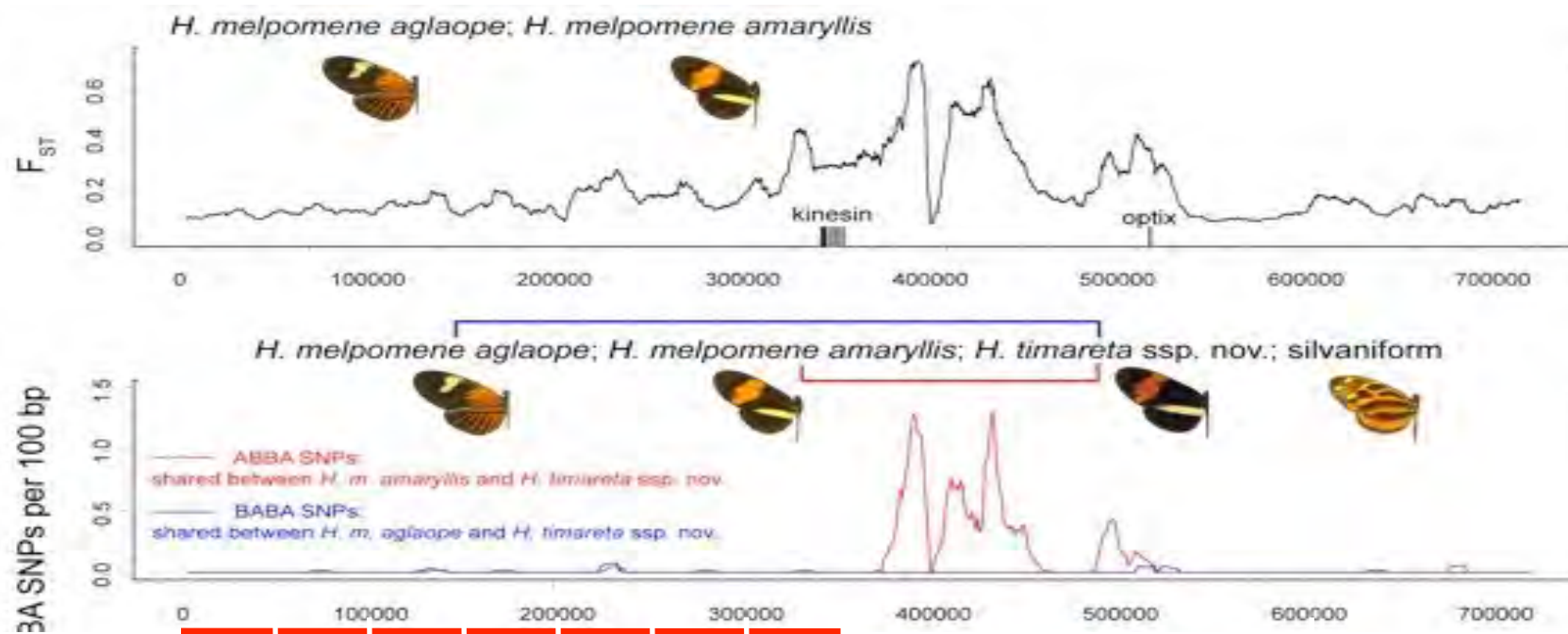
50-kb windows along *B/D*

# Phylogenies across *B/D*

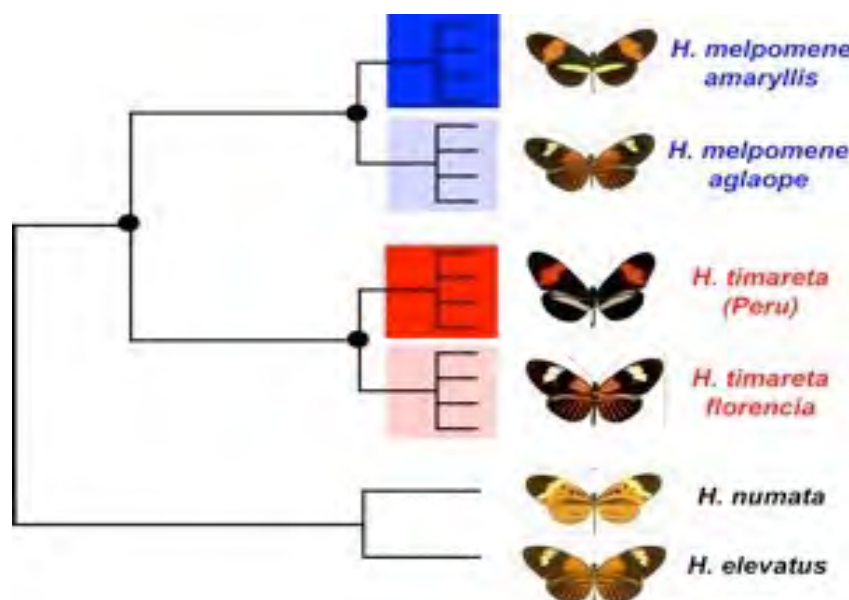


ML tree based  
on  
50,000 bp

# Phylogenies across *B/D*

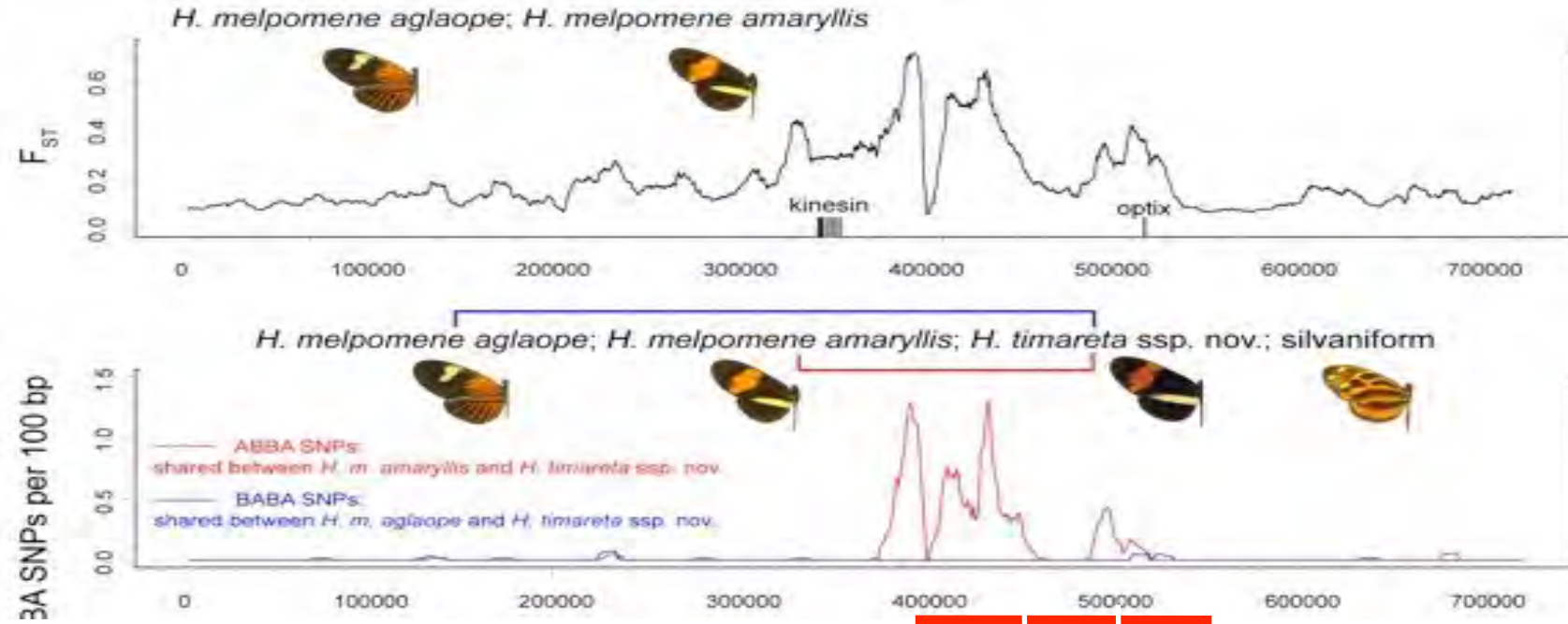


ML tree based  
on  
50,000 bp



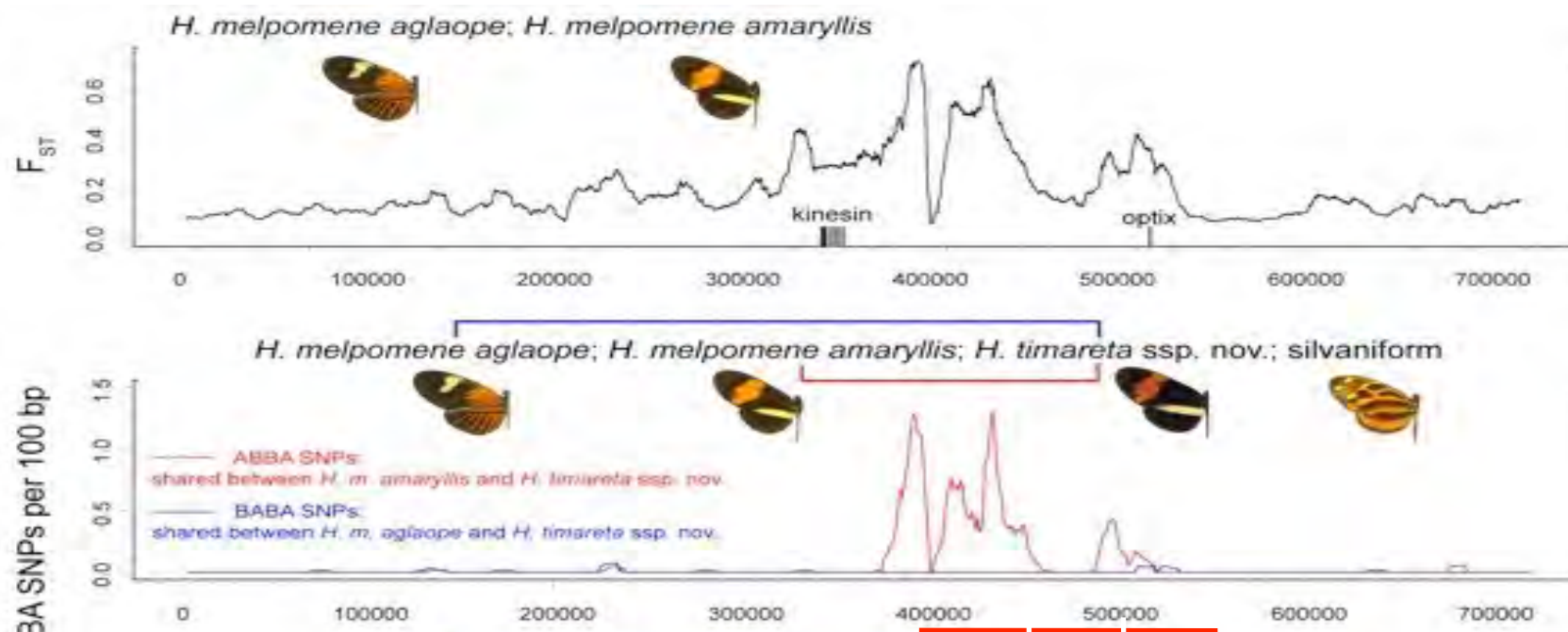


# Phylogenies across *B/D*

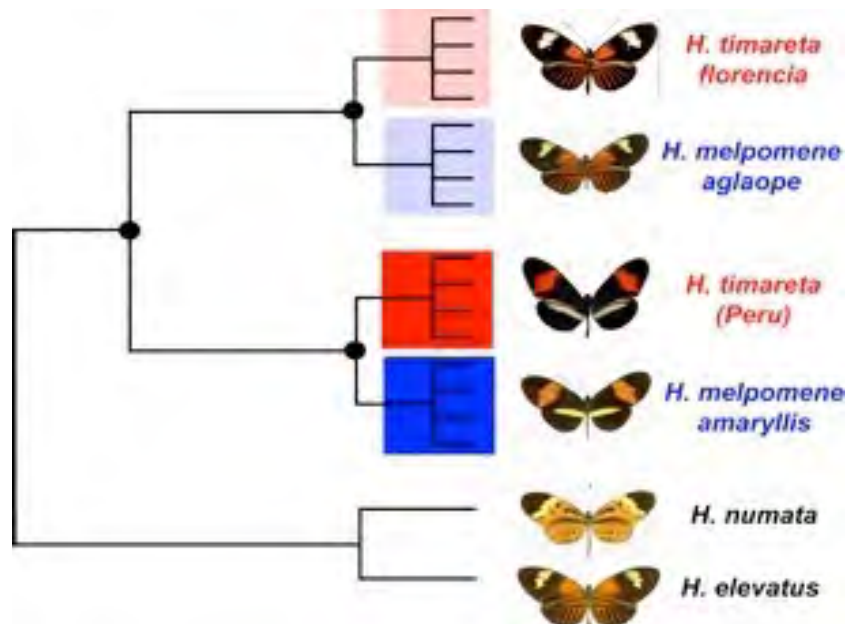


ML tree based  
on  
50,000 bp

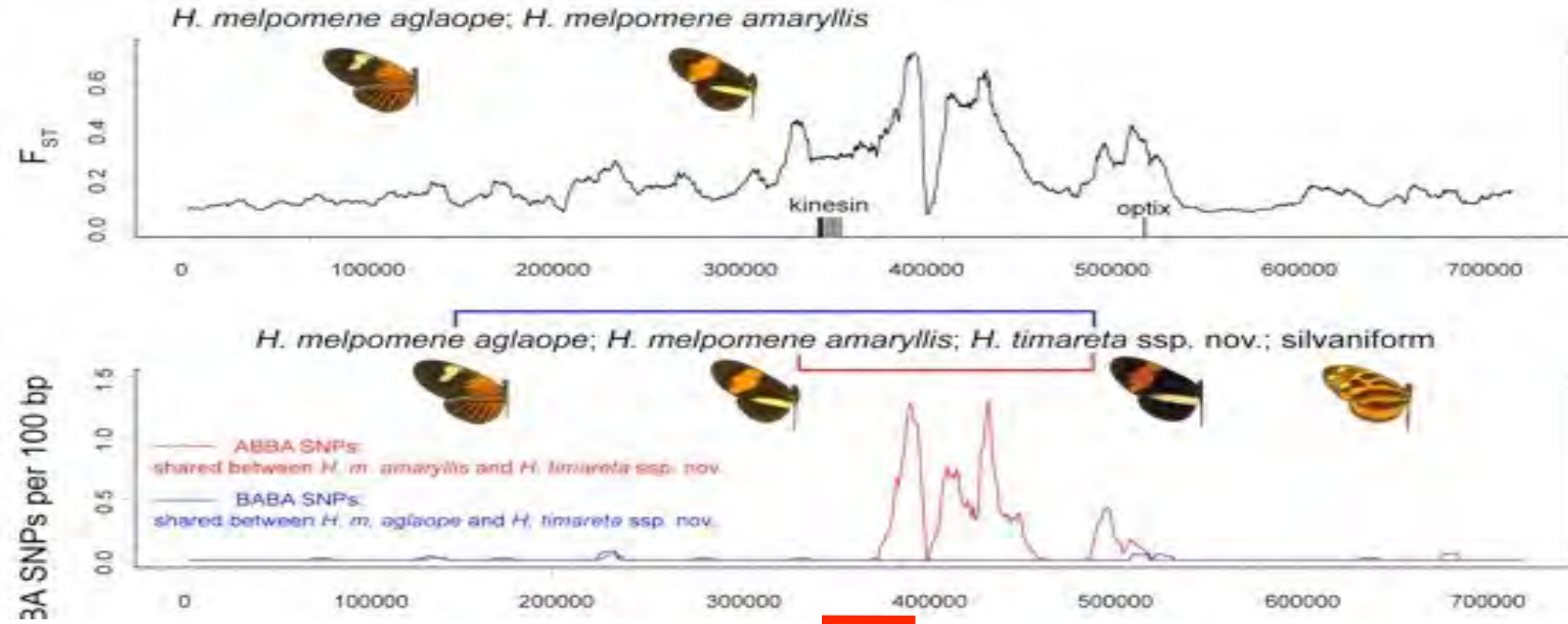
# Phylogenies across *B/D*



ML tree based  
on  
50,000 bp

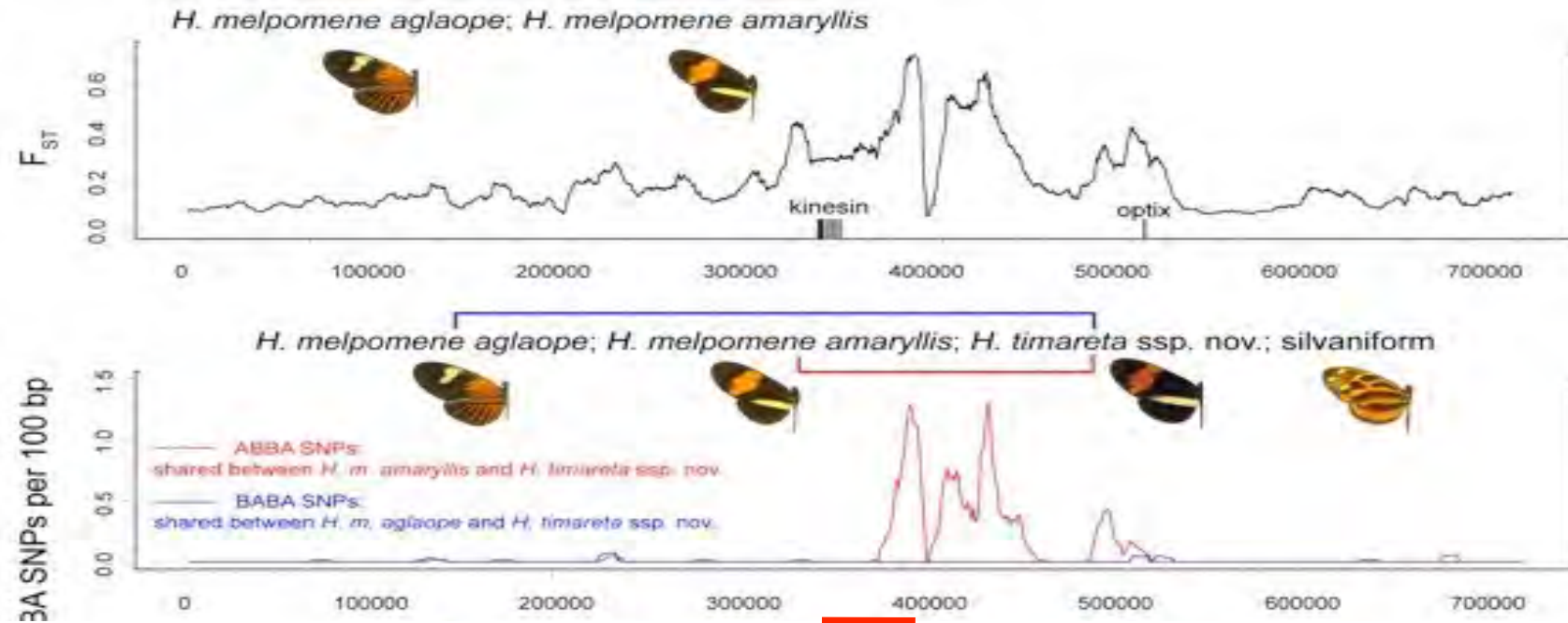


# Phylogenies across *B/D*

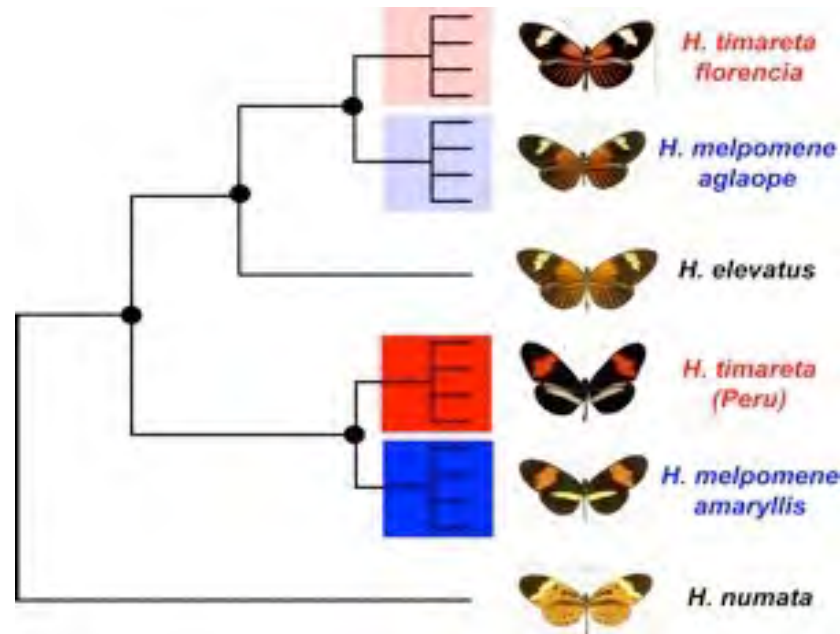


ML tree based  
on  
50,000 bp

# Phylogenies across *B/D*

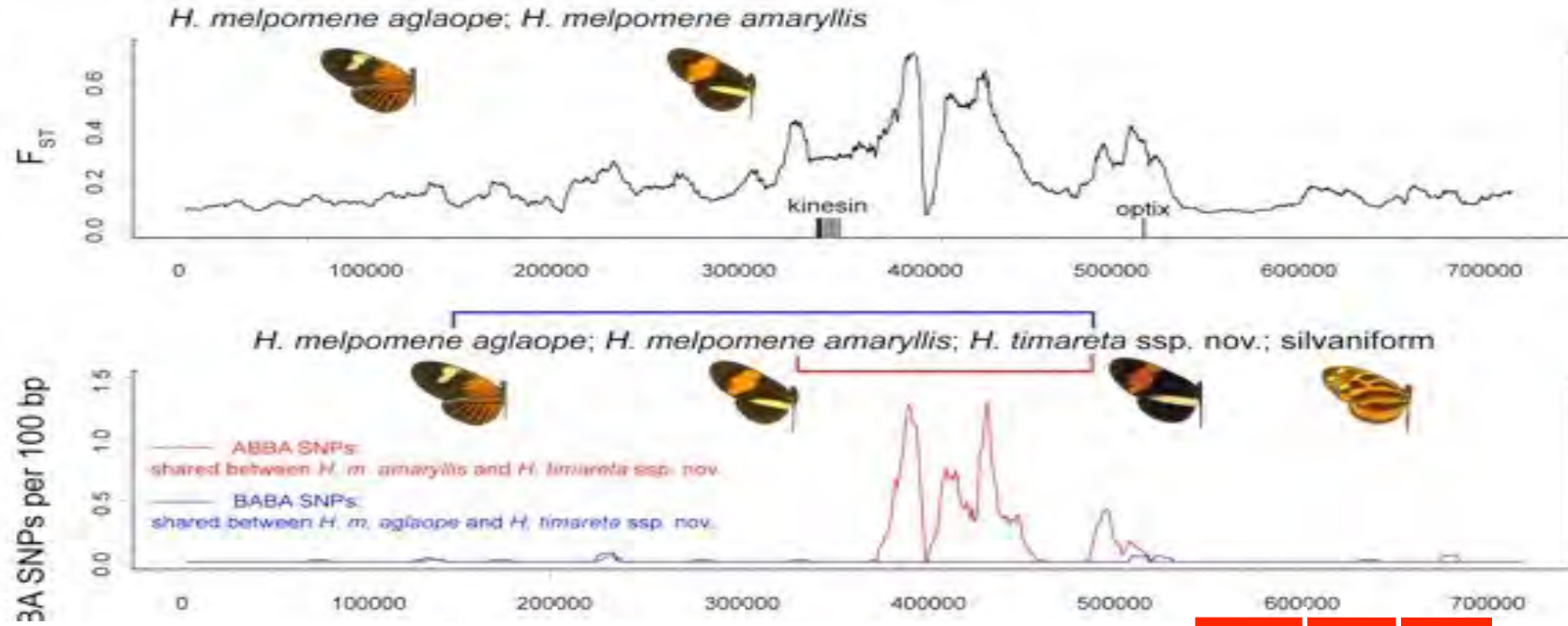


ML tree based  
on  
50,000 bp



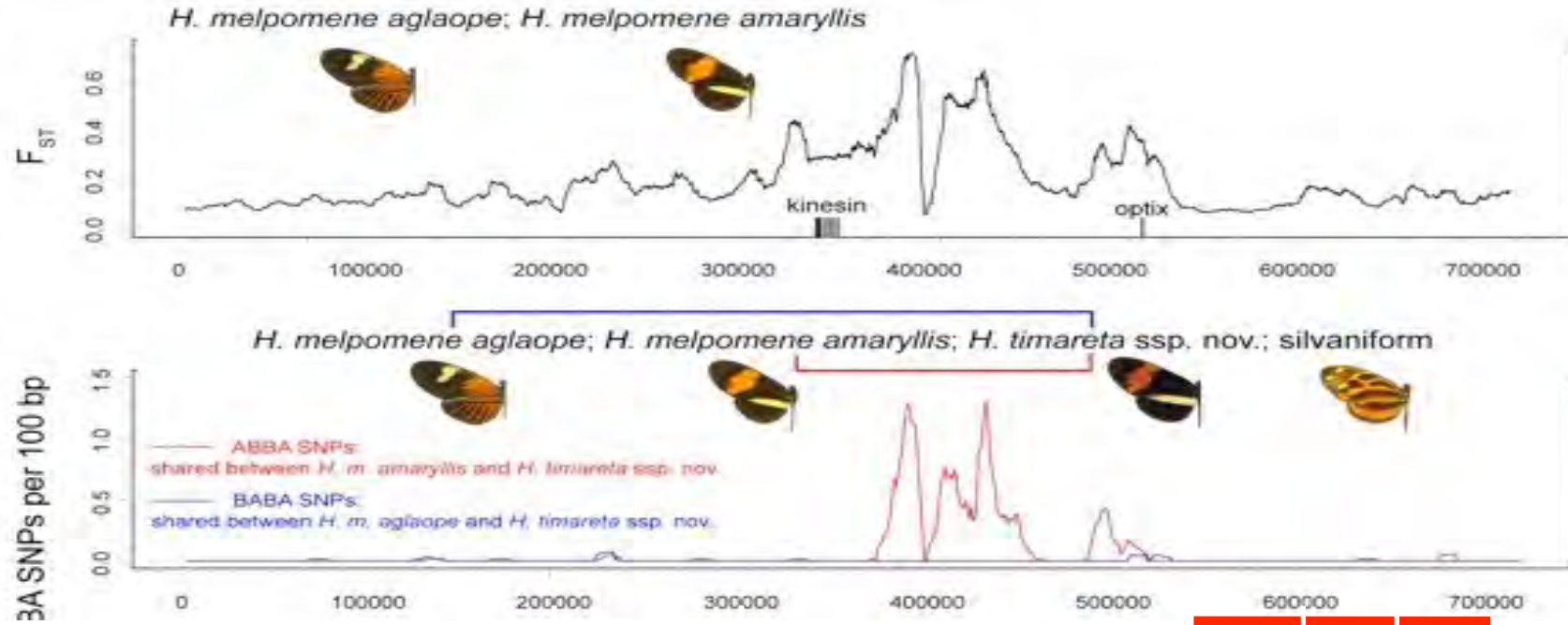


# Phylogenies across *B/D*

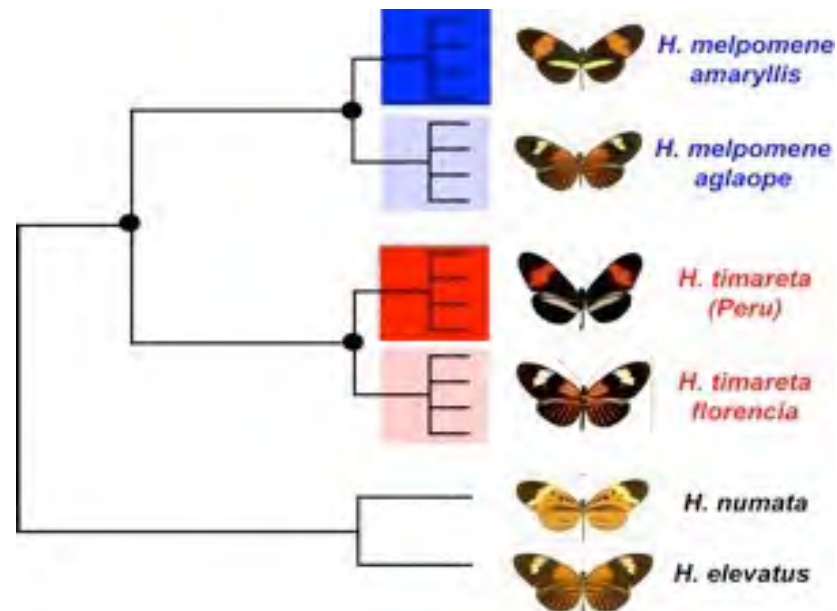


ML tree based  
on  
50,000 bp

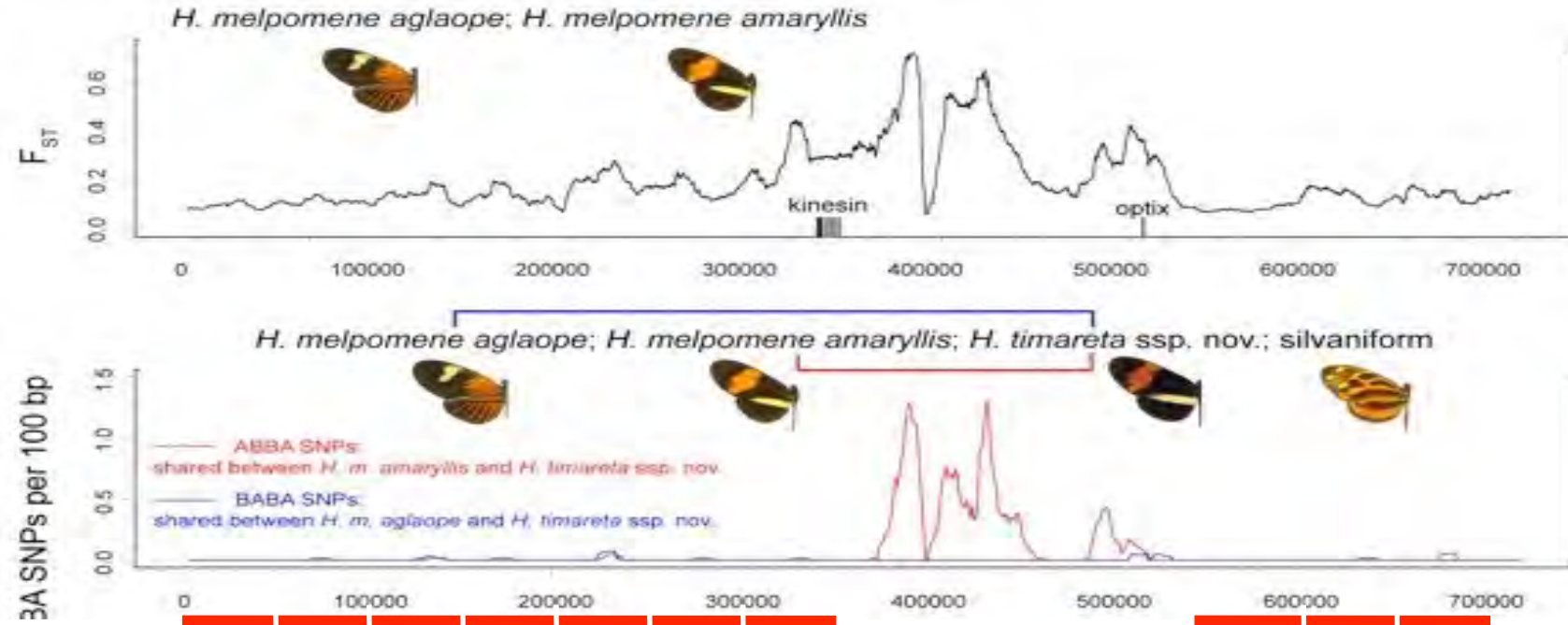
# Phylogenies across *B/D*



ML tree based  
on  
50,000 bp

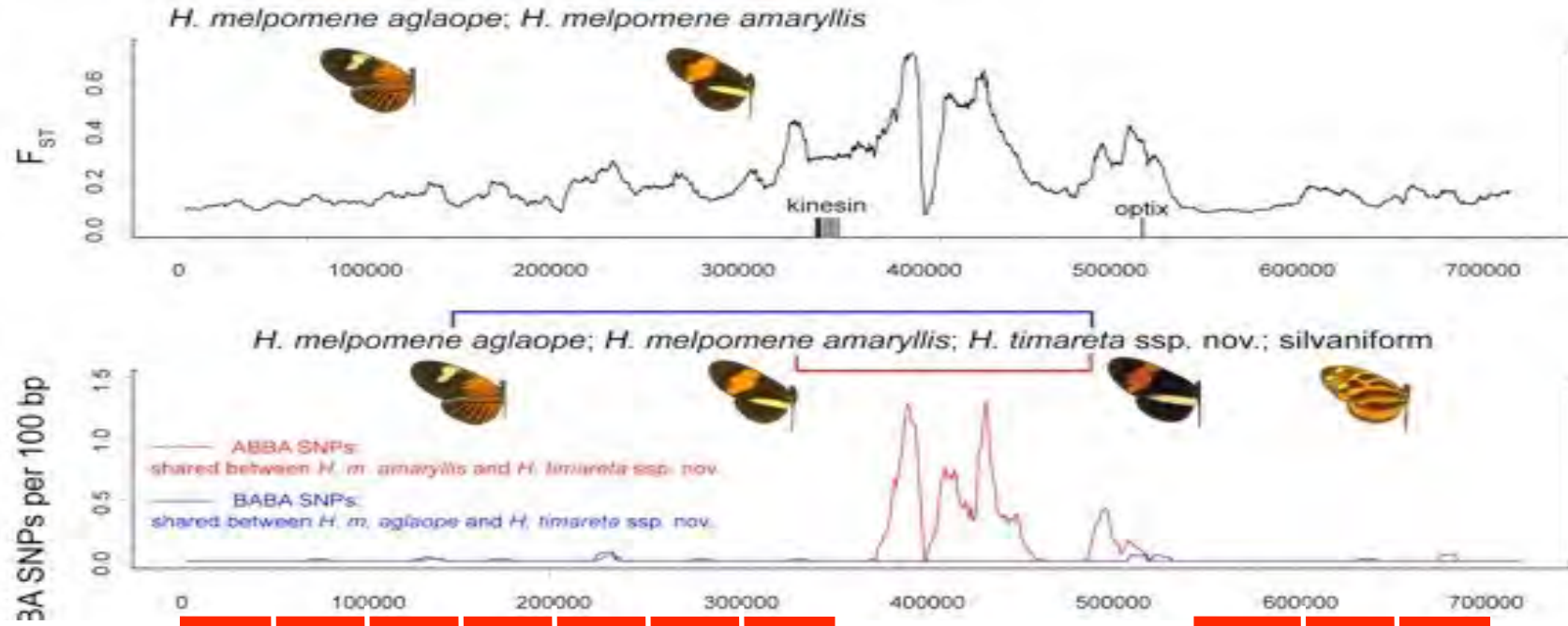


# Phylogenies across *B/D*

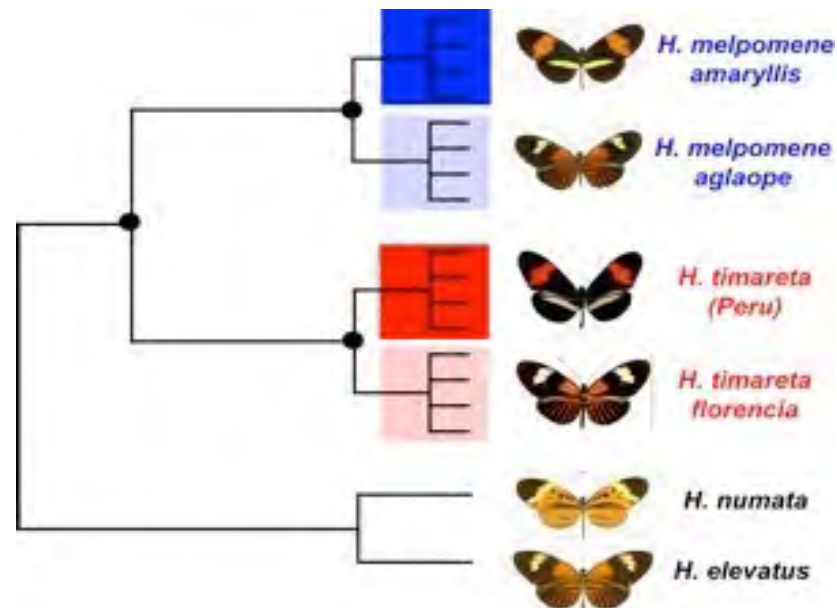


ML tree based  
on  
50,000 bp

# Phylogenies across *B/D*



ML tree based  
on  
50,000 bp





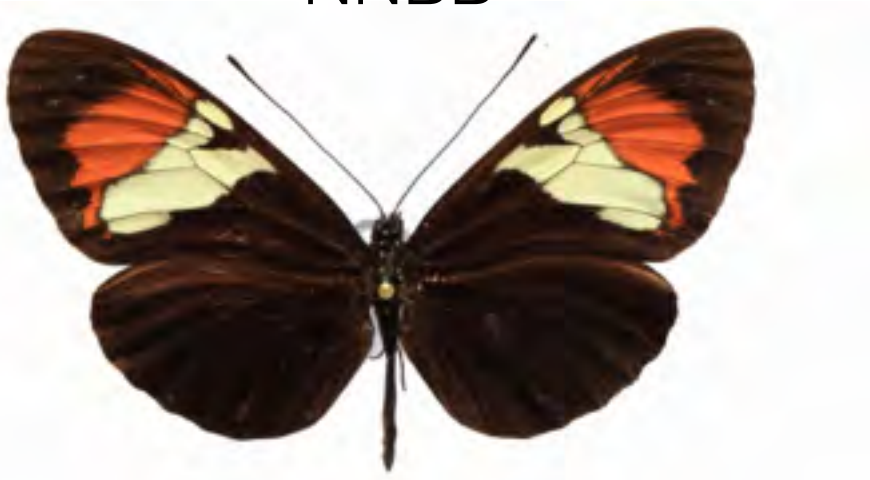
Okay, so introgression causes  
mimicry

Okay, so introgression causes mimicry

But mimicry is weird, right?

# Novelty can arise through introgression and recombination

NNBB



*Heliconius heurippa*



Camilo Salazar

Mavarez et al., Nature 2006

# Novelty can arise through introgression and recombination

NNBB



*Heliconius heurippa*

NNbb



*Heliconius cydno cordula*

nnBB



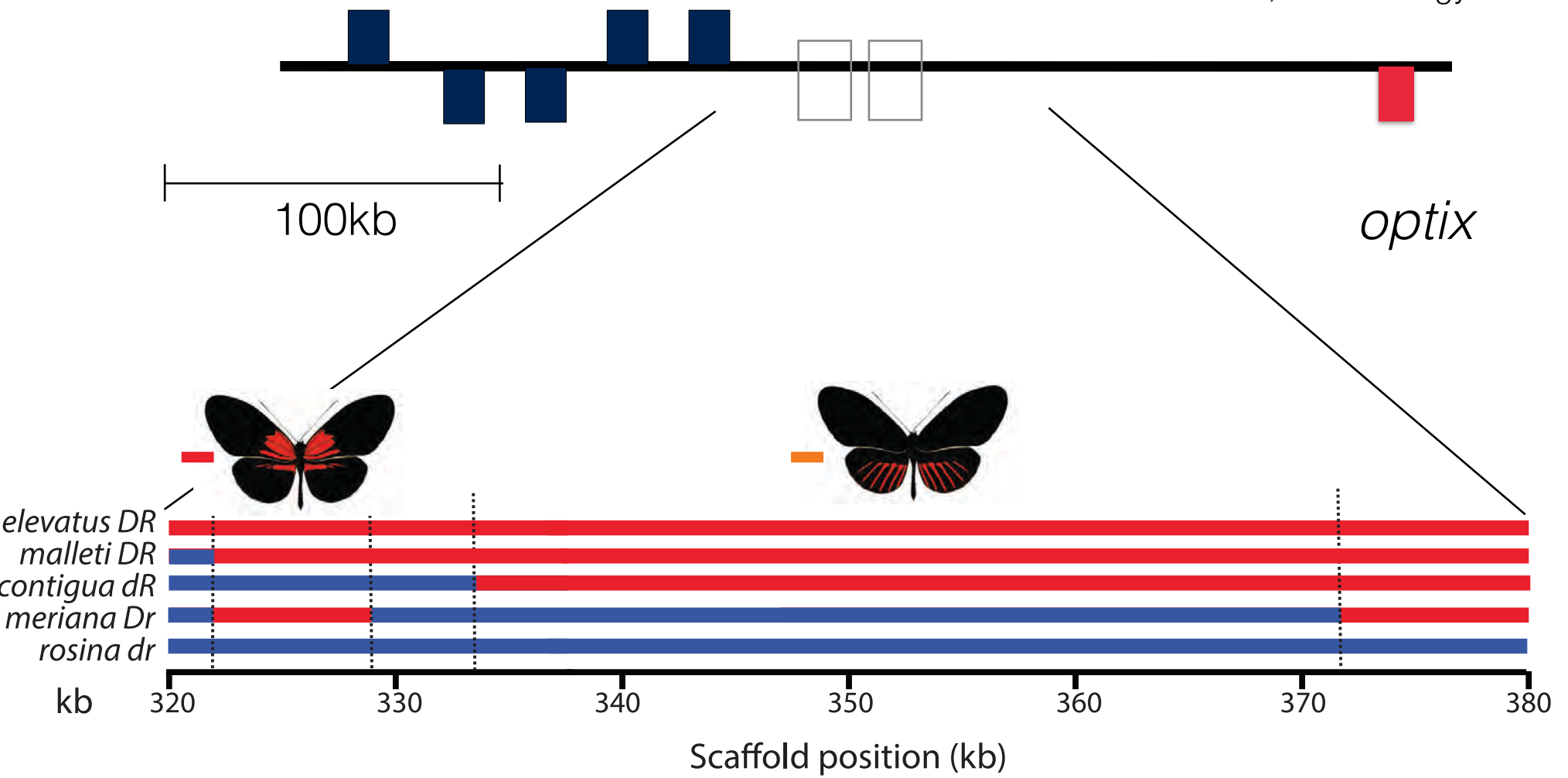
*Heliconius melpomene melpomene*

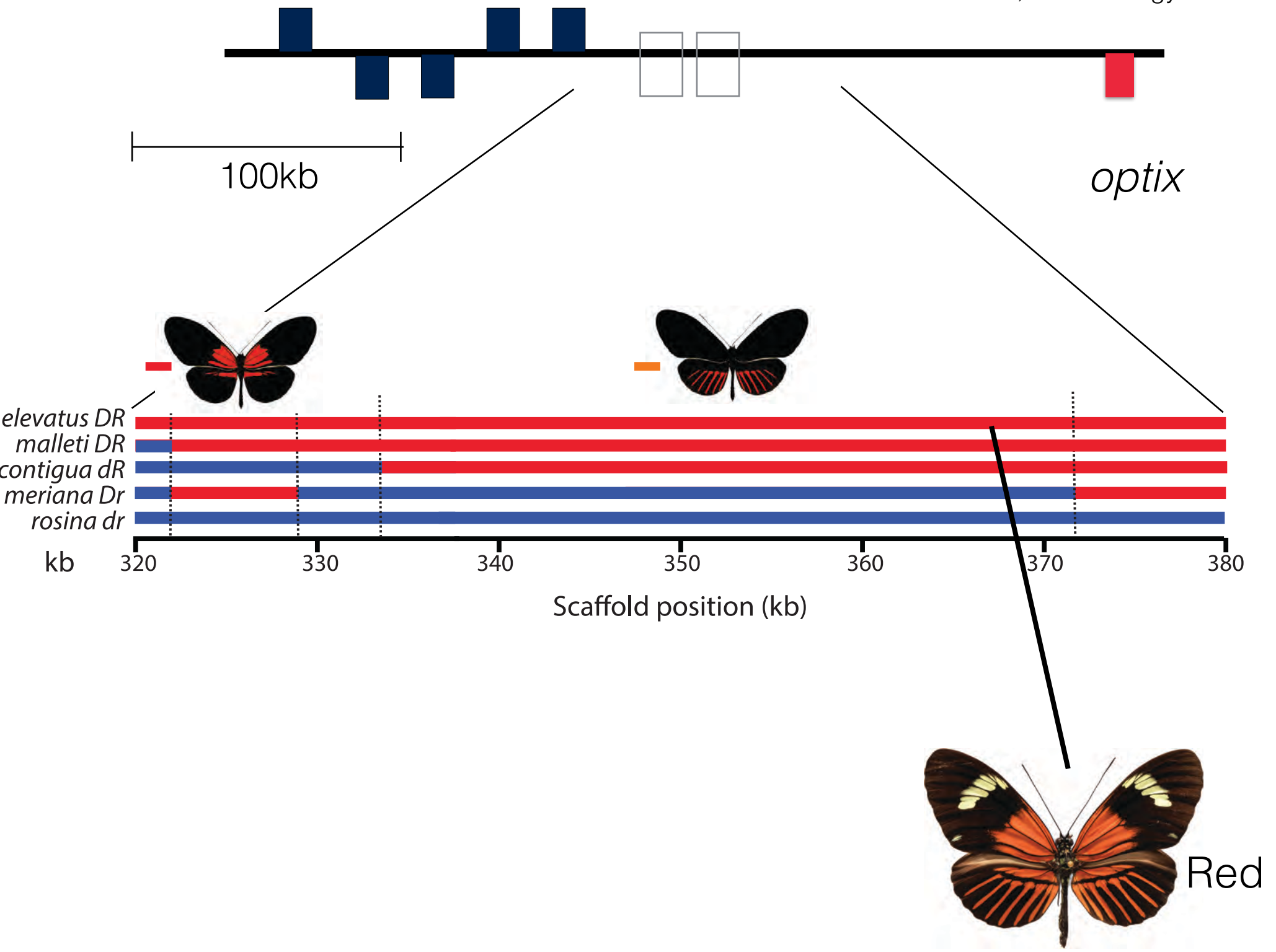


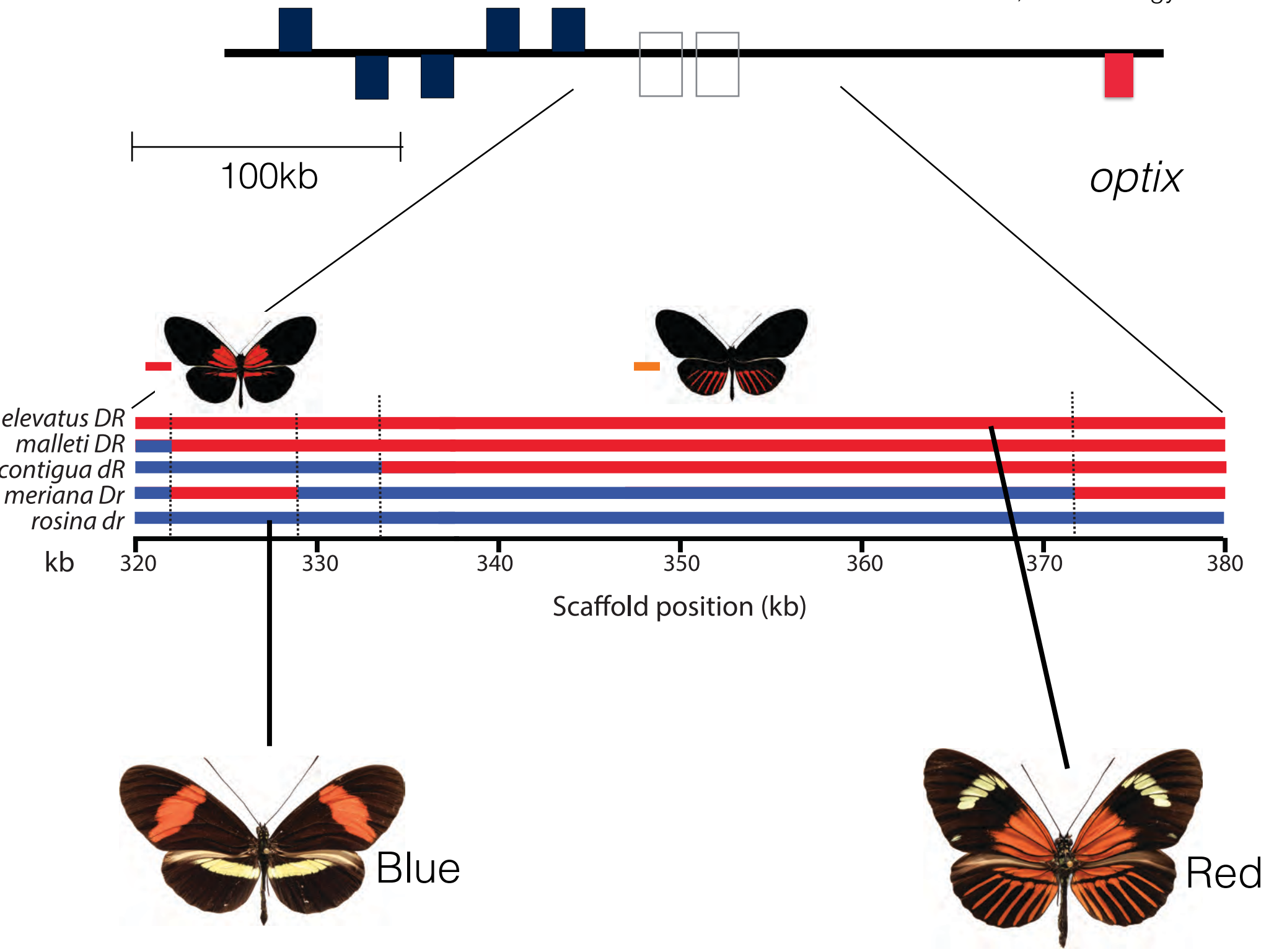
Camilo Salazar

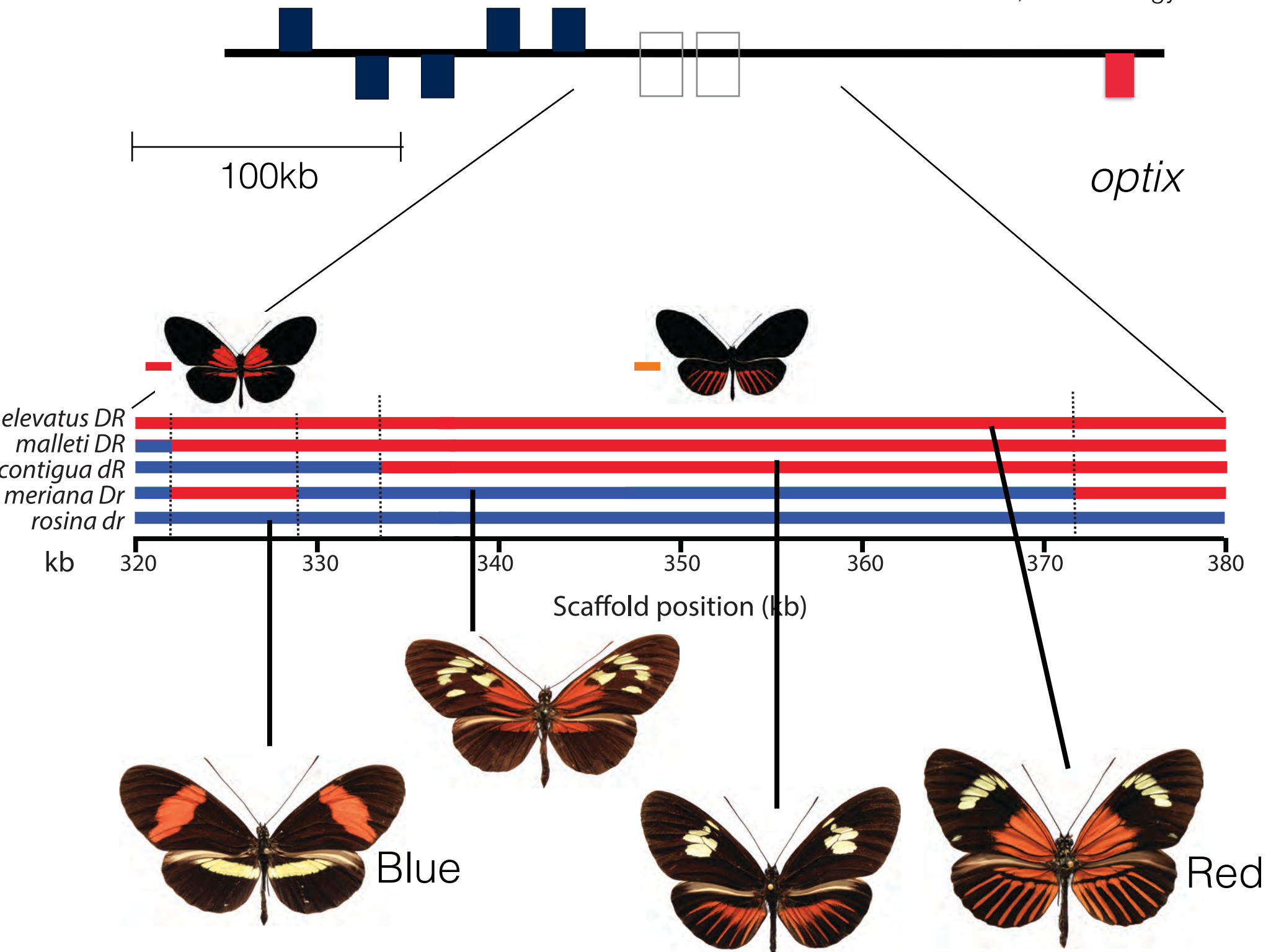
Mavarez et al., Nature 2006



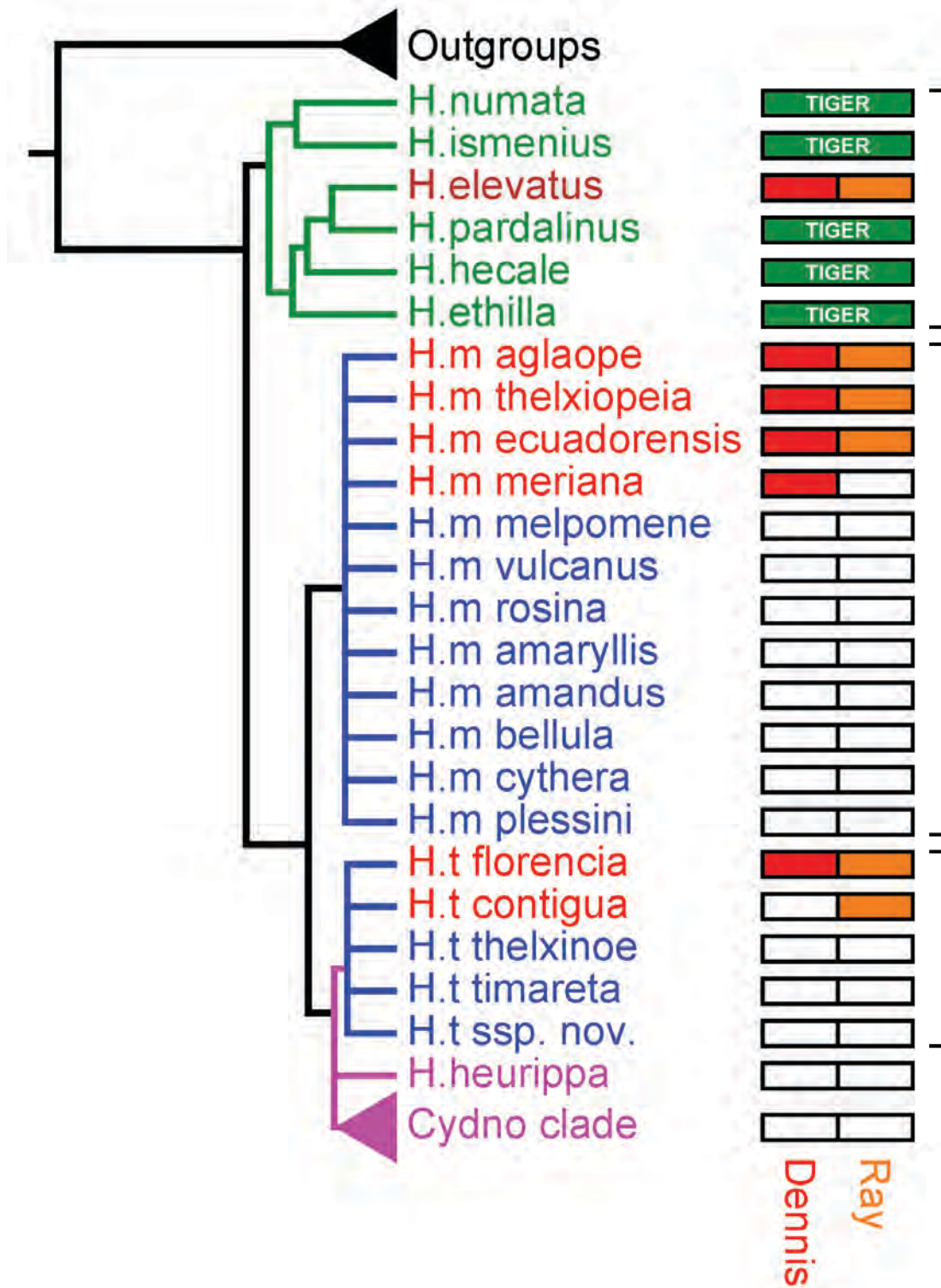


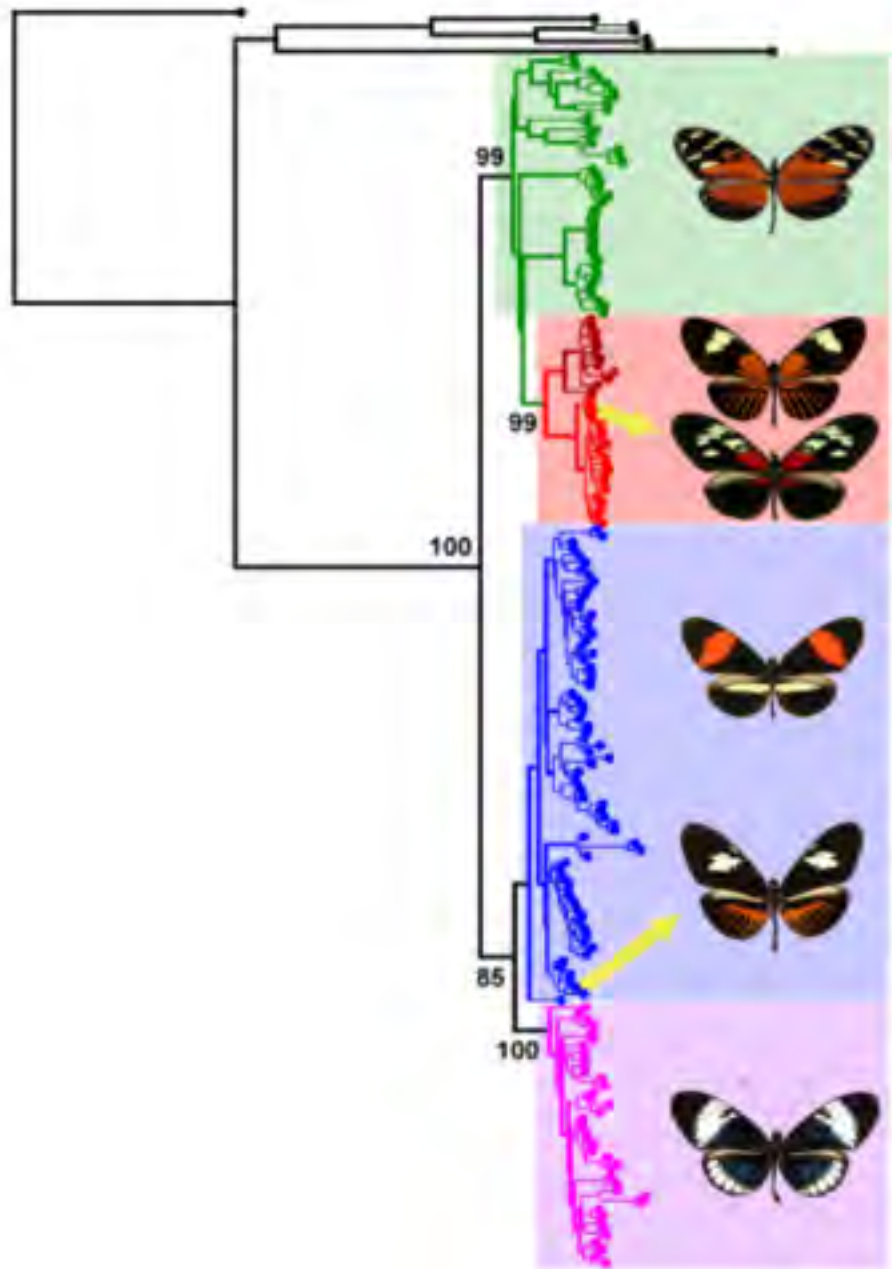




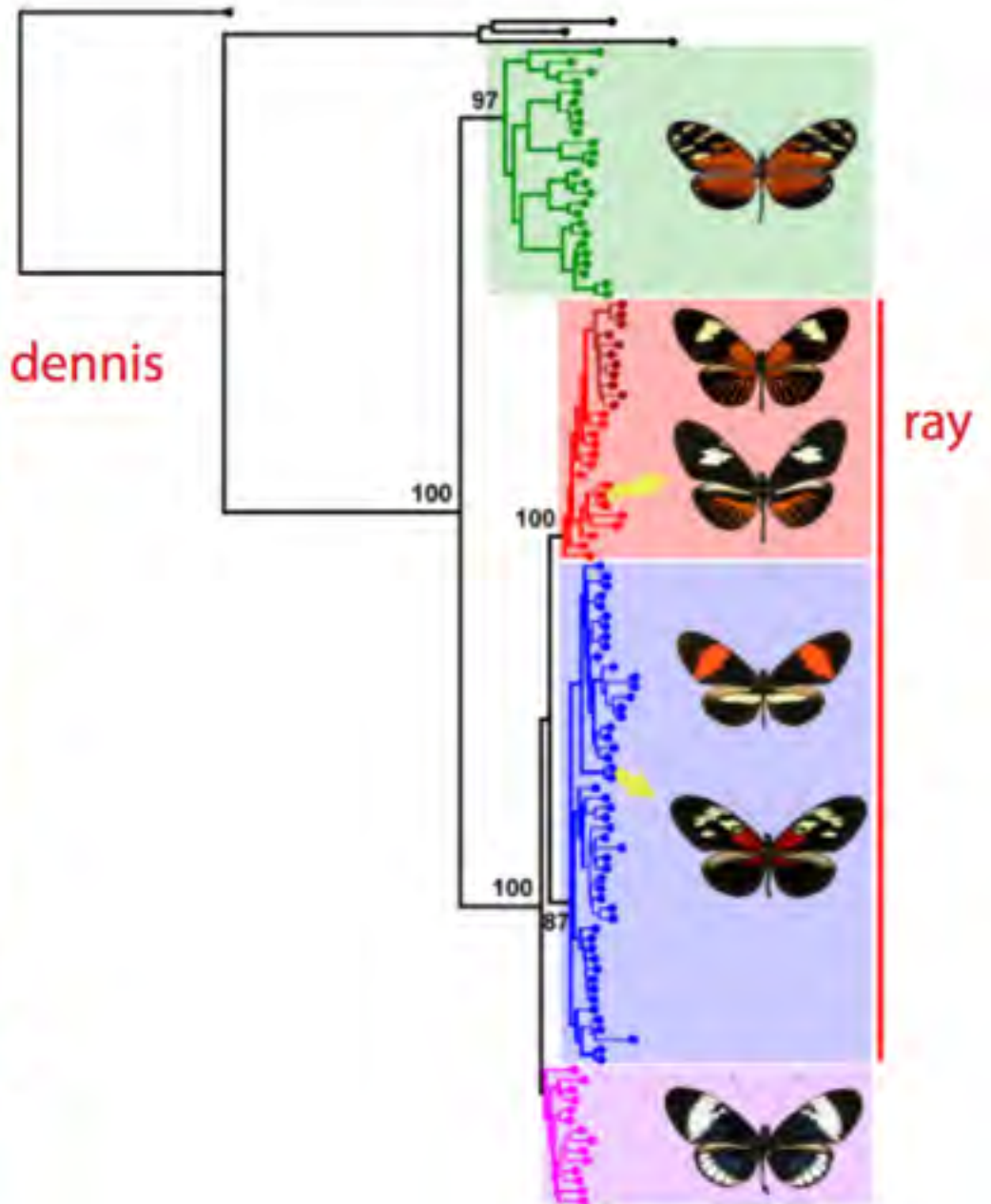






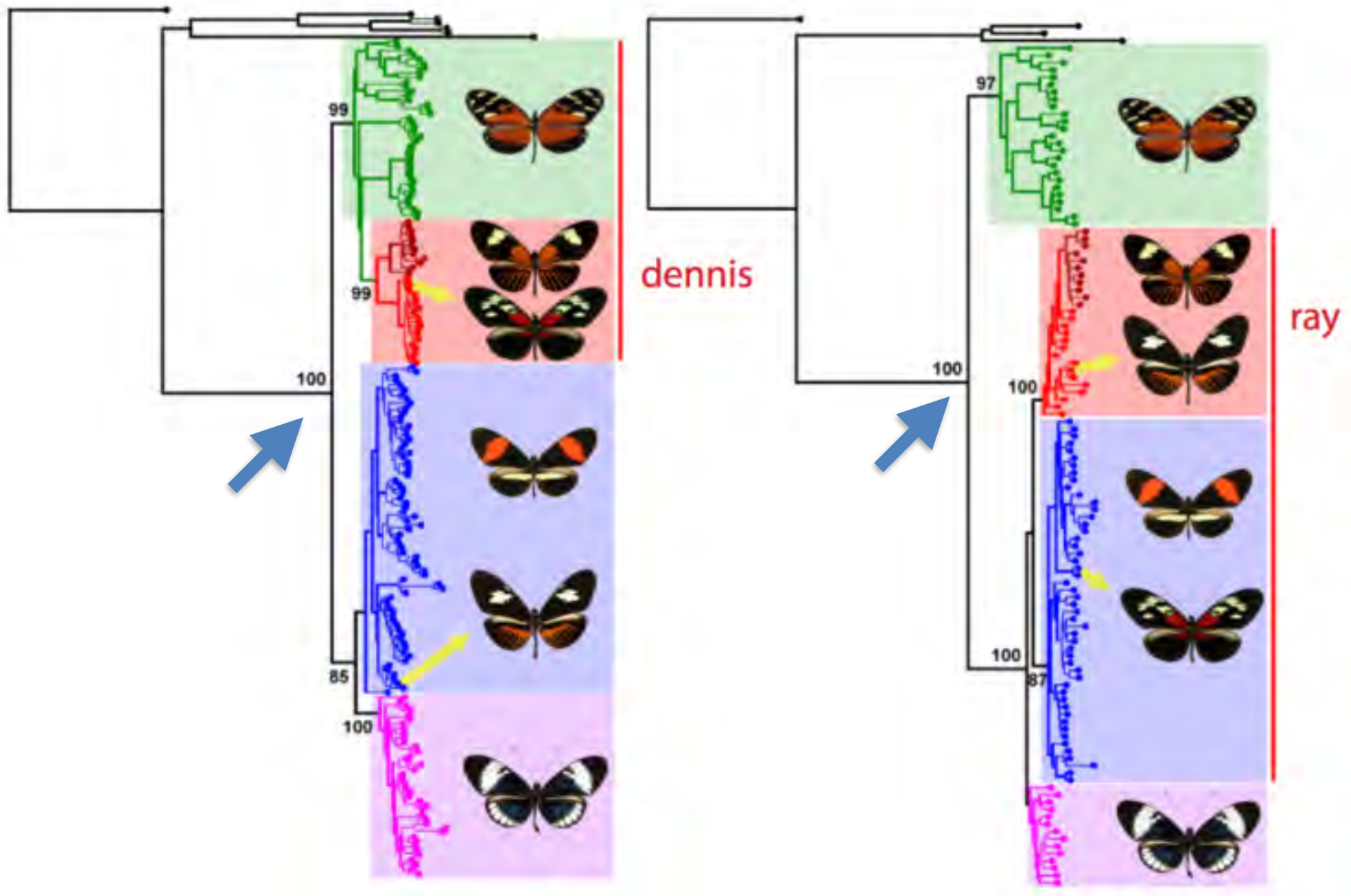


dennis

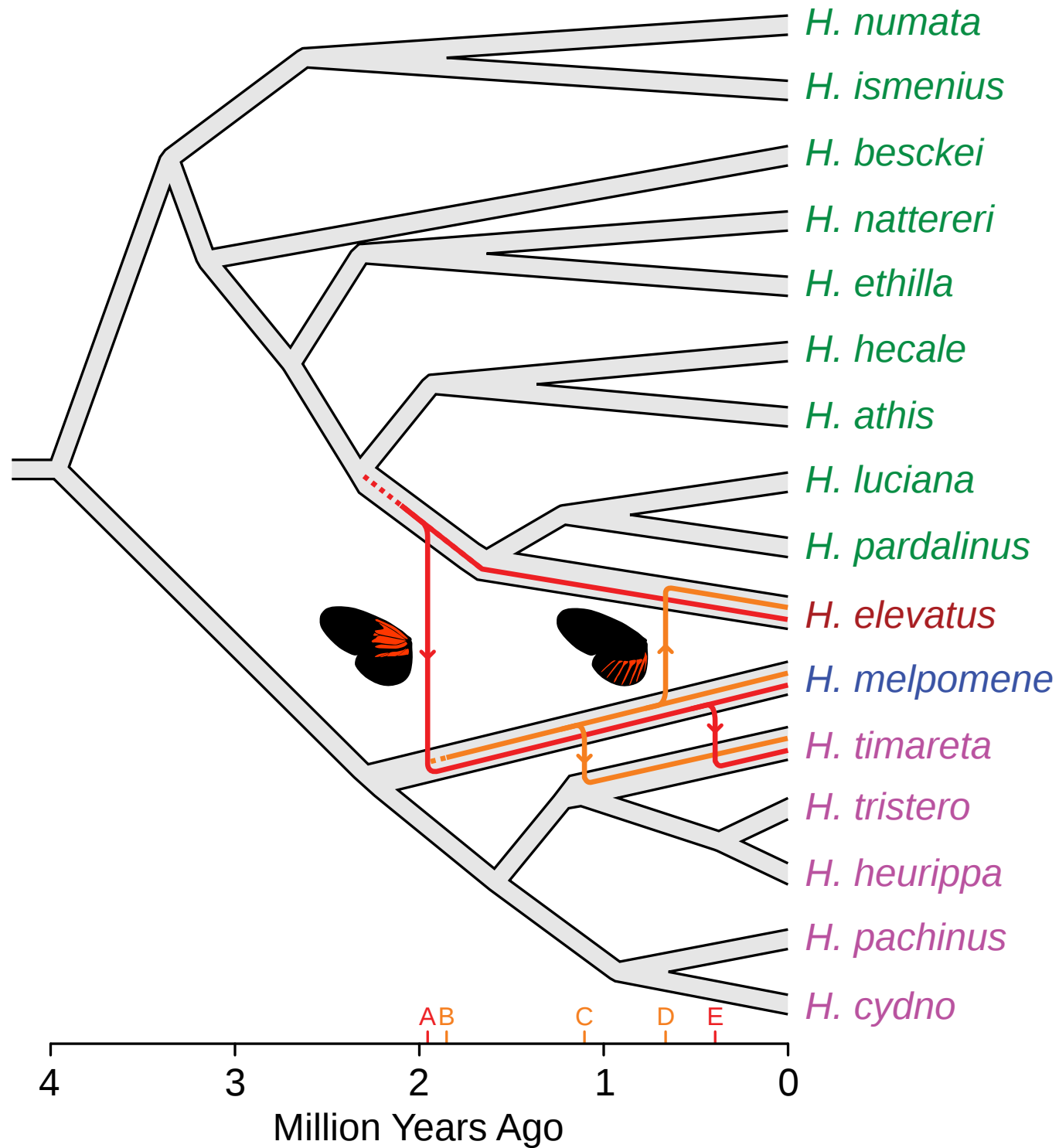


ray





Generate dated trees using this node as a reference point





# What about behaviour?

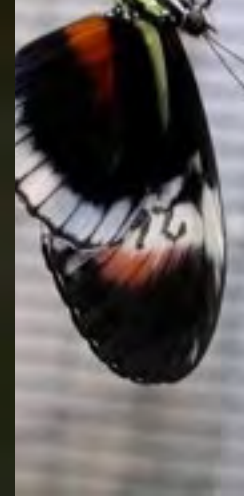
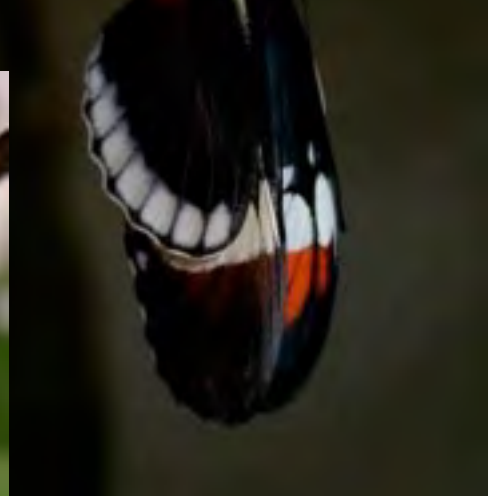
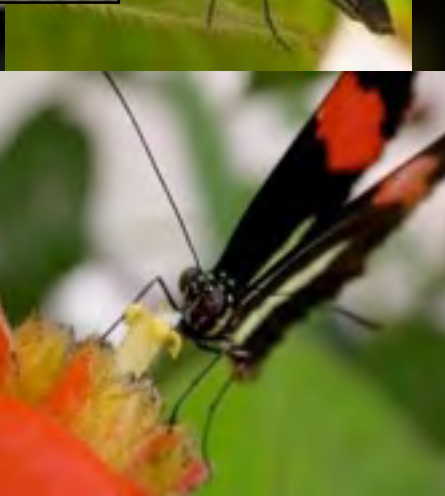
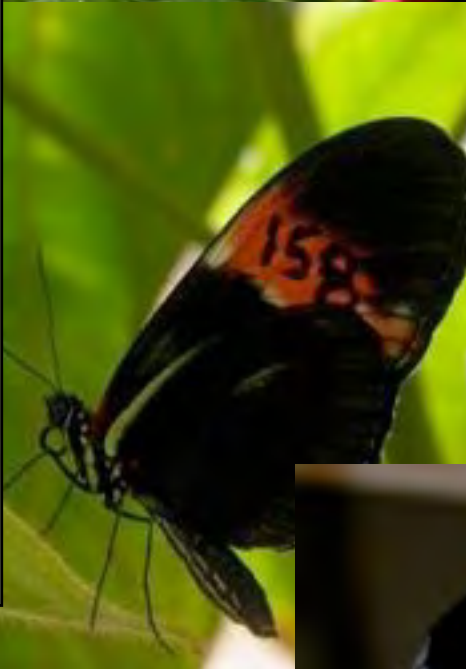
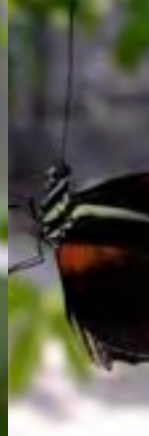
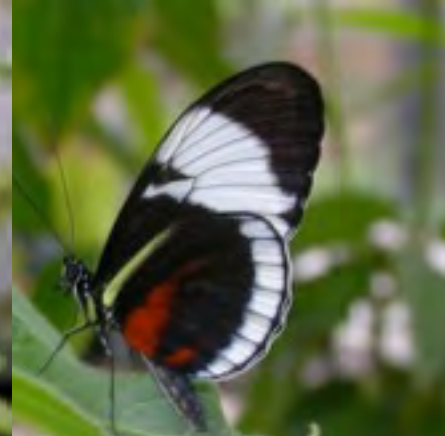


# What about behaviour?

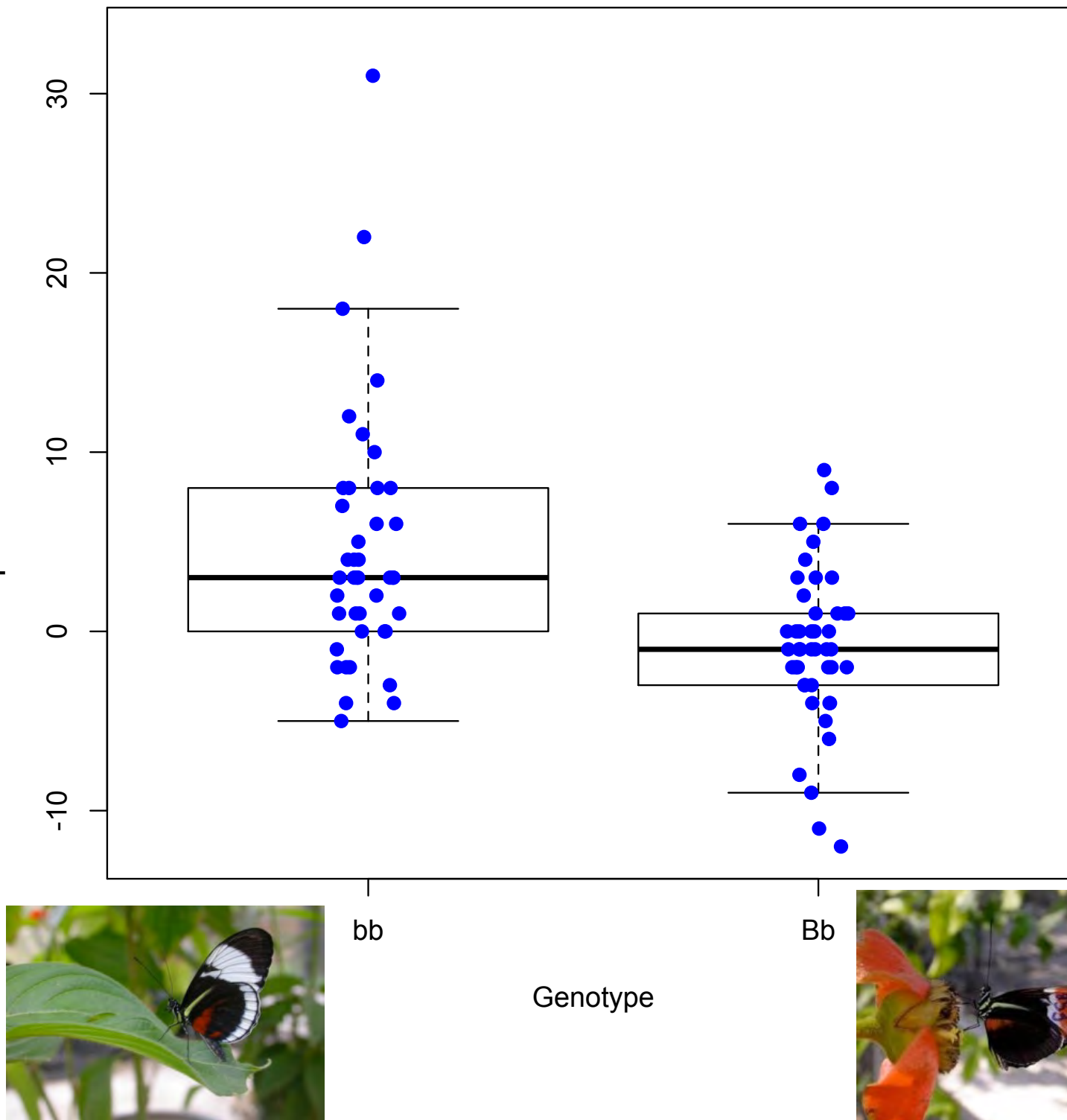




H. melpomene X H. cydno



Difference between  
approaches to cydno and  
melpomene



*bb*



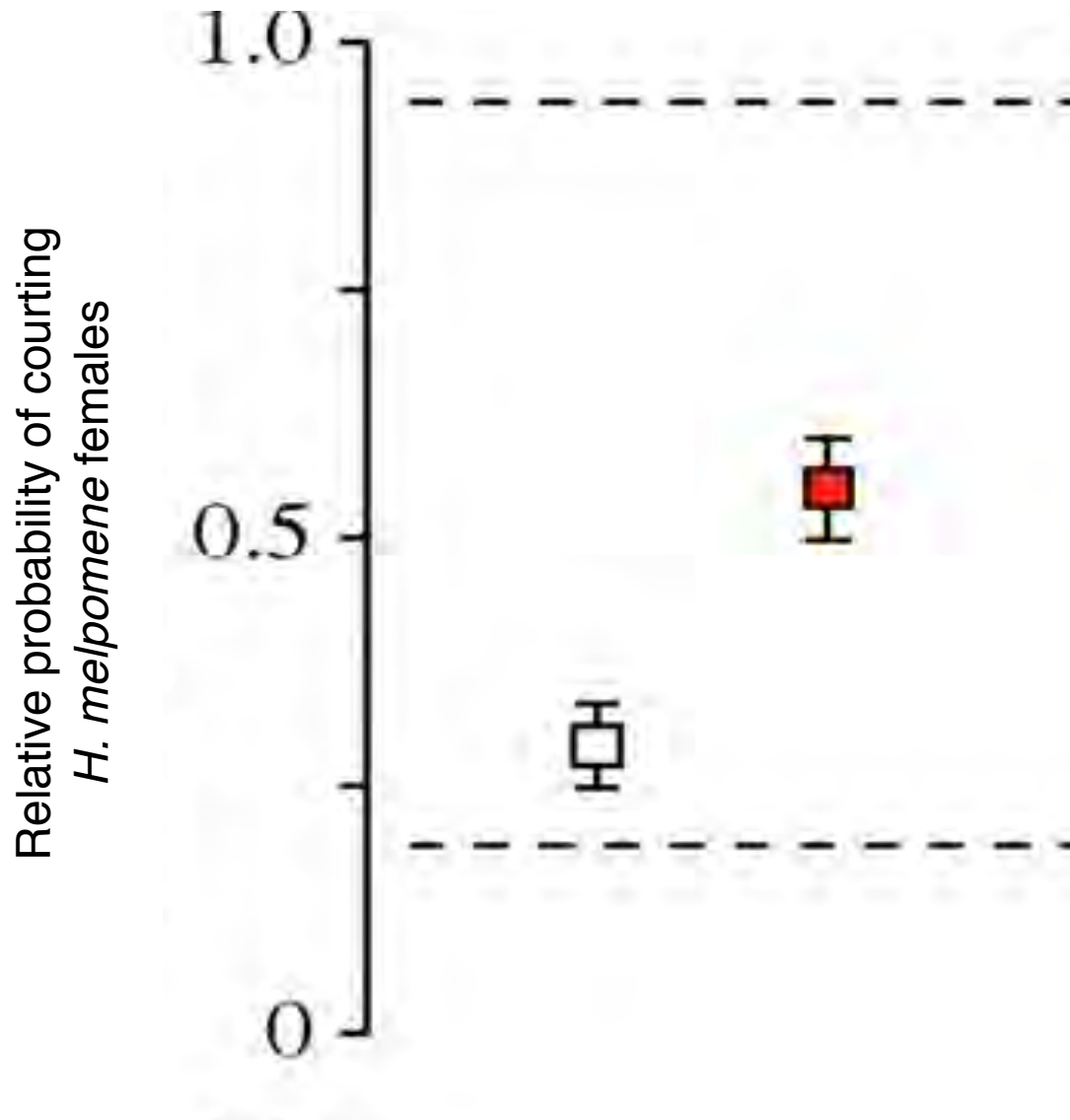
*Bb*



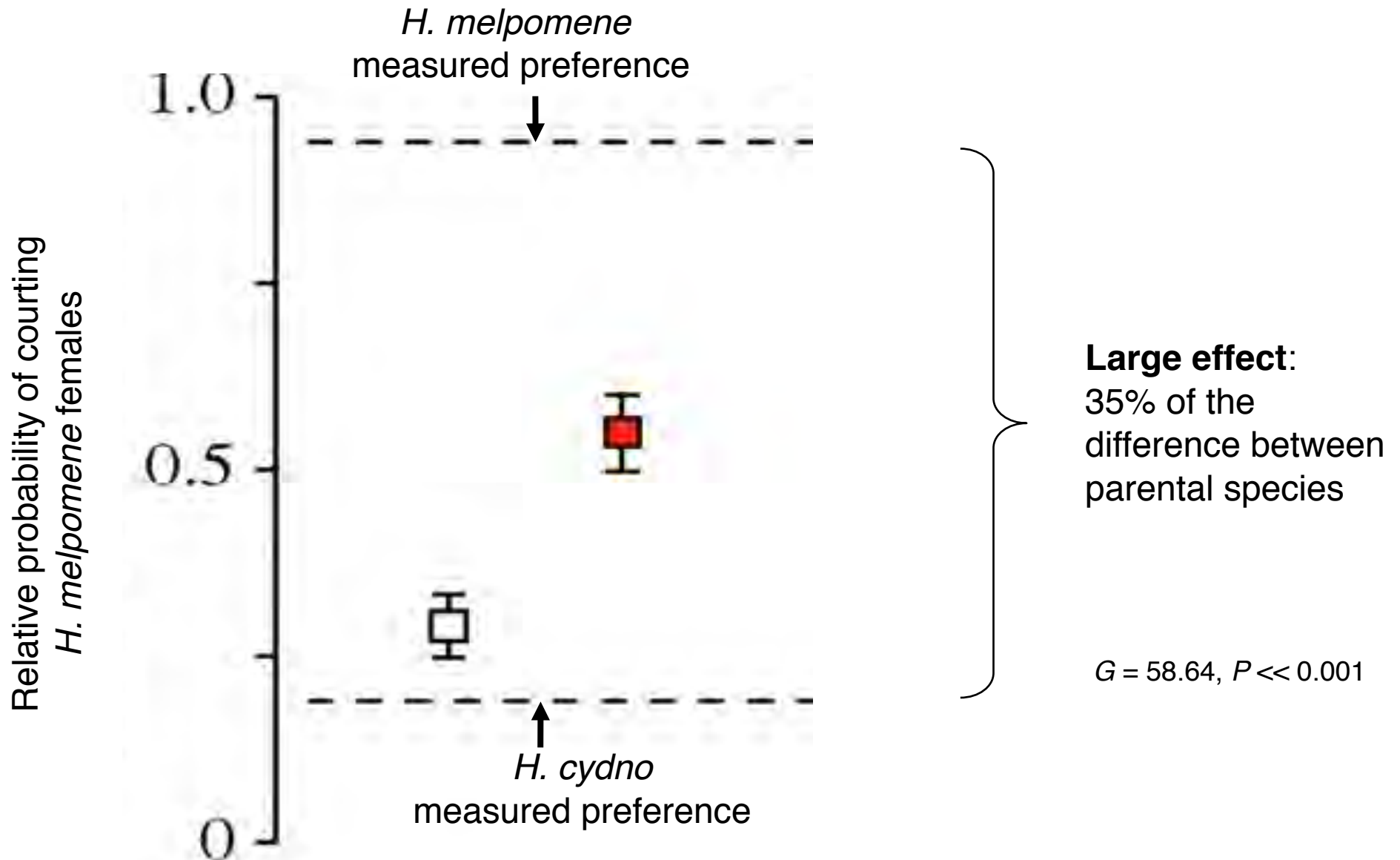
Richard  
Merrill



## Mate preference segregates with forewing colour in backcross hybrids



# Mate preference segregates with forewing colour in backcross hybrids





NNbb

*Heliconius cydno cordula*



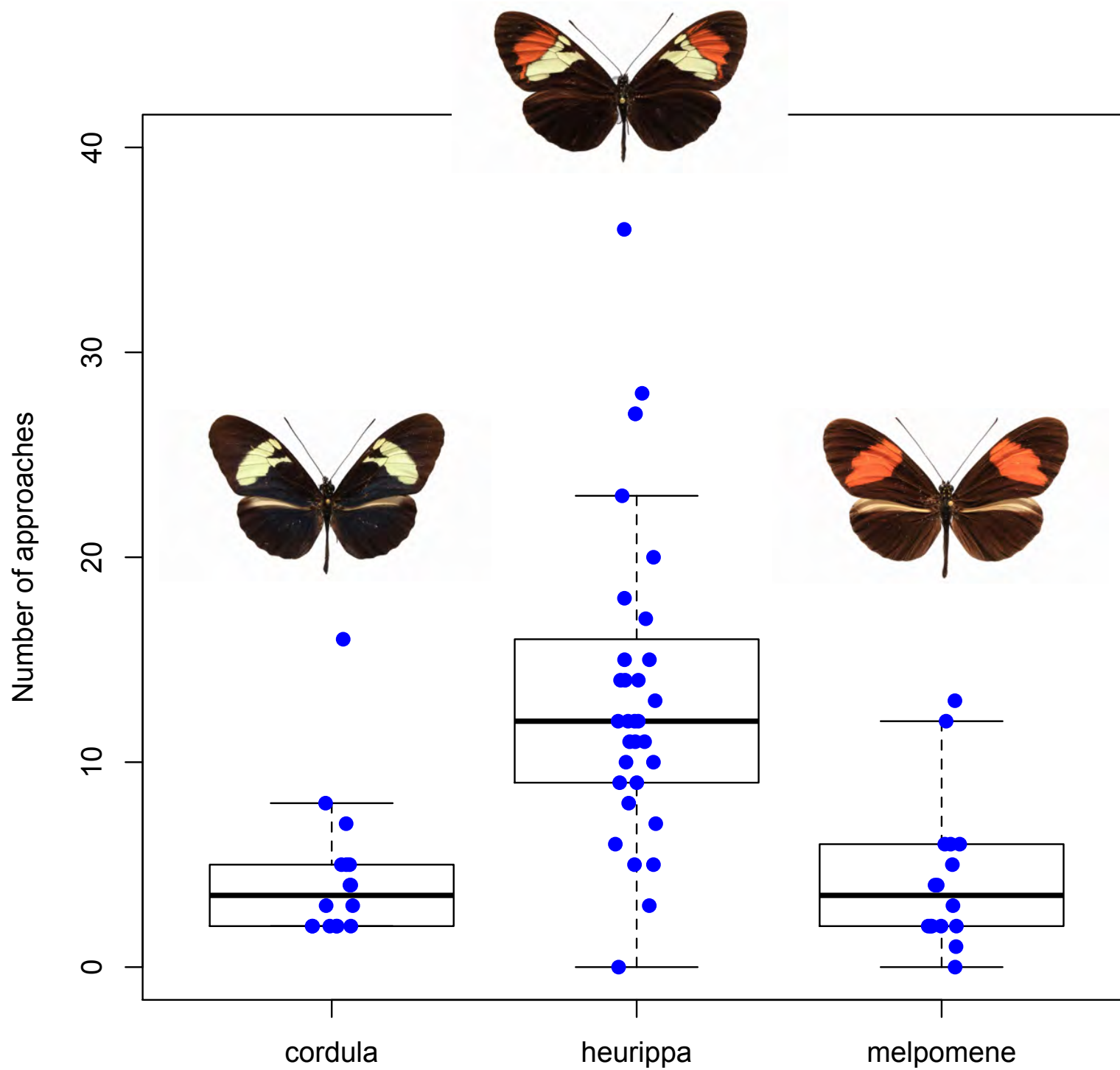
nnBB

*Heliconius melpomene melpomene*



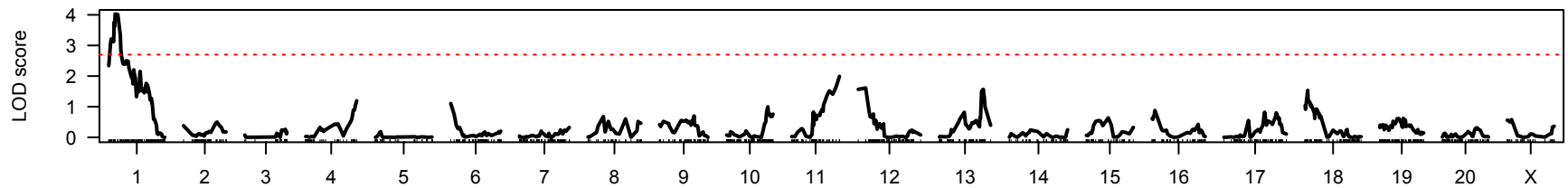
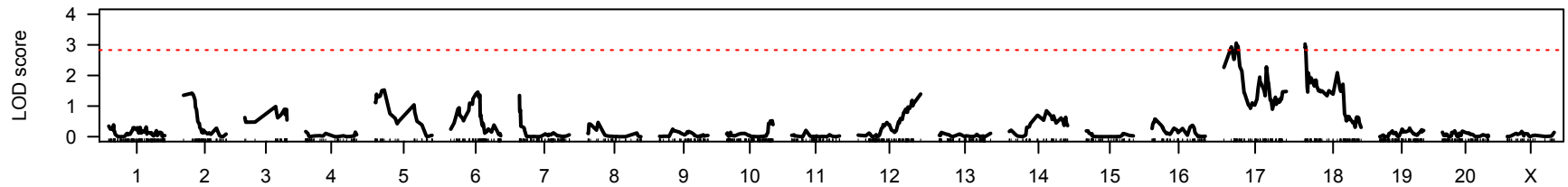
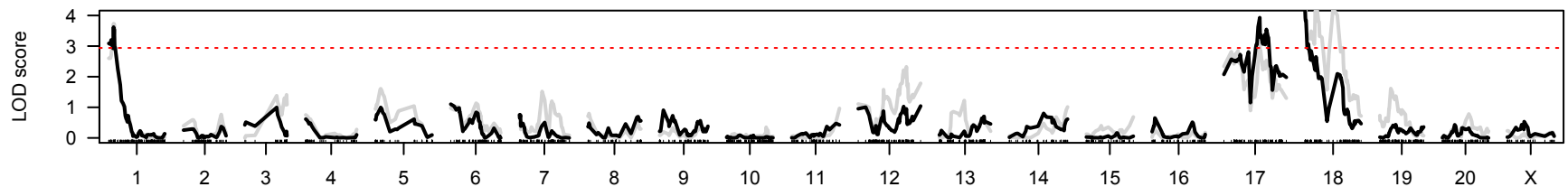
NNBB

*Heliconius heurippa*



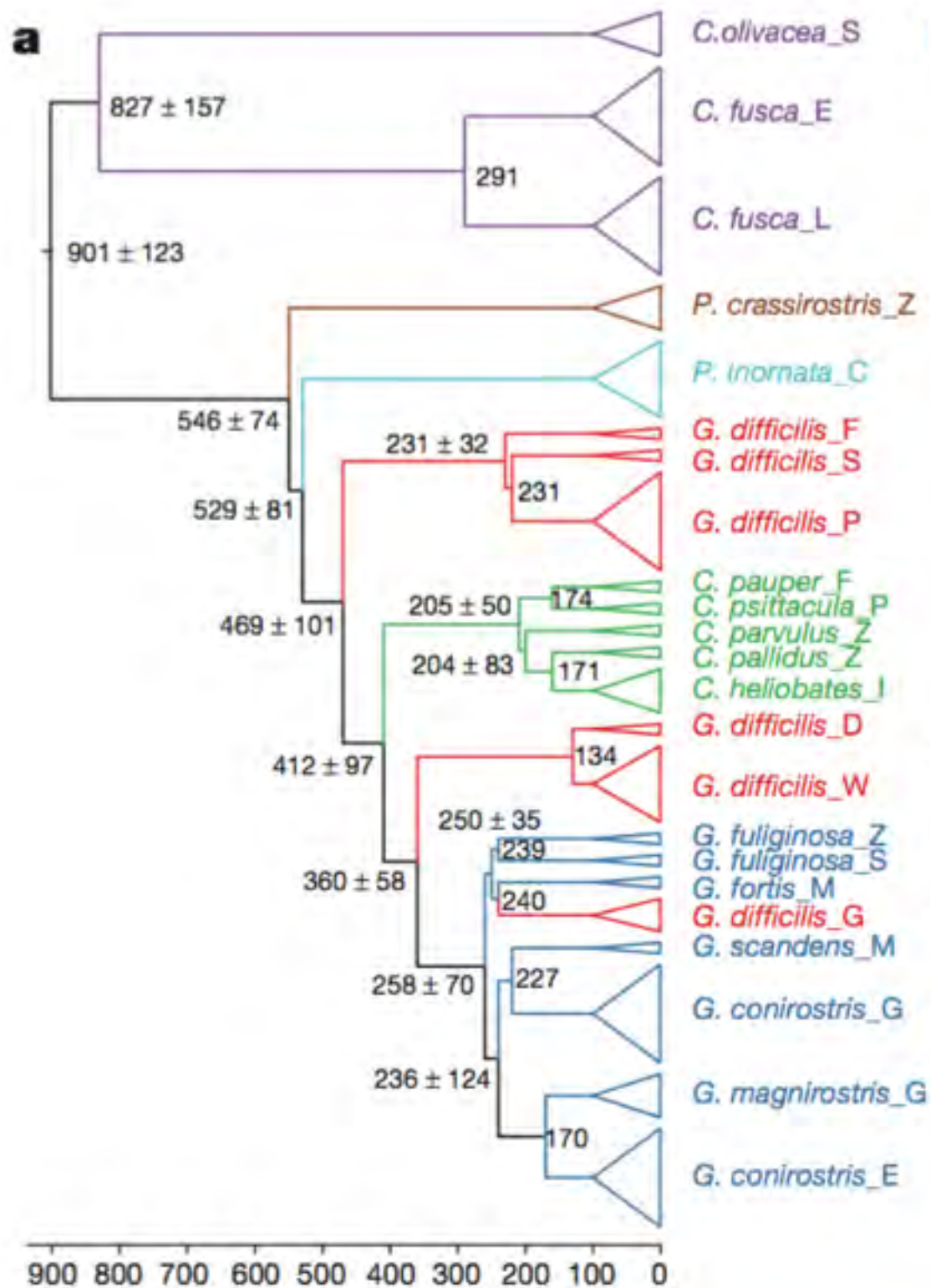


# Now working on QTL maps of species differences in behaviour:

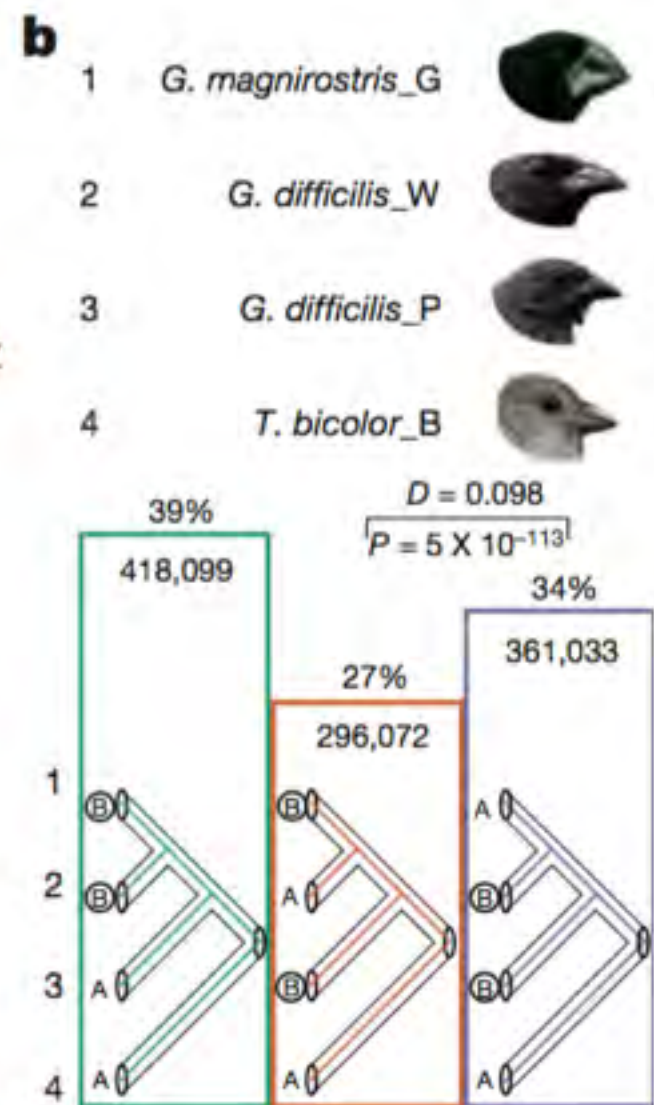
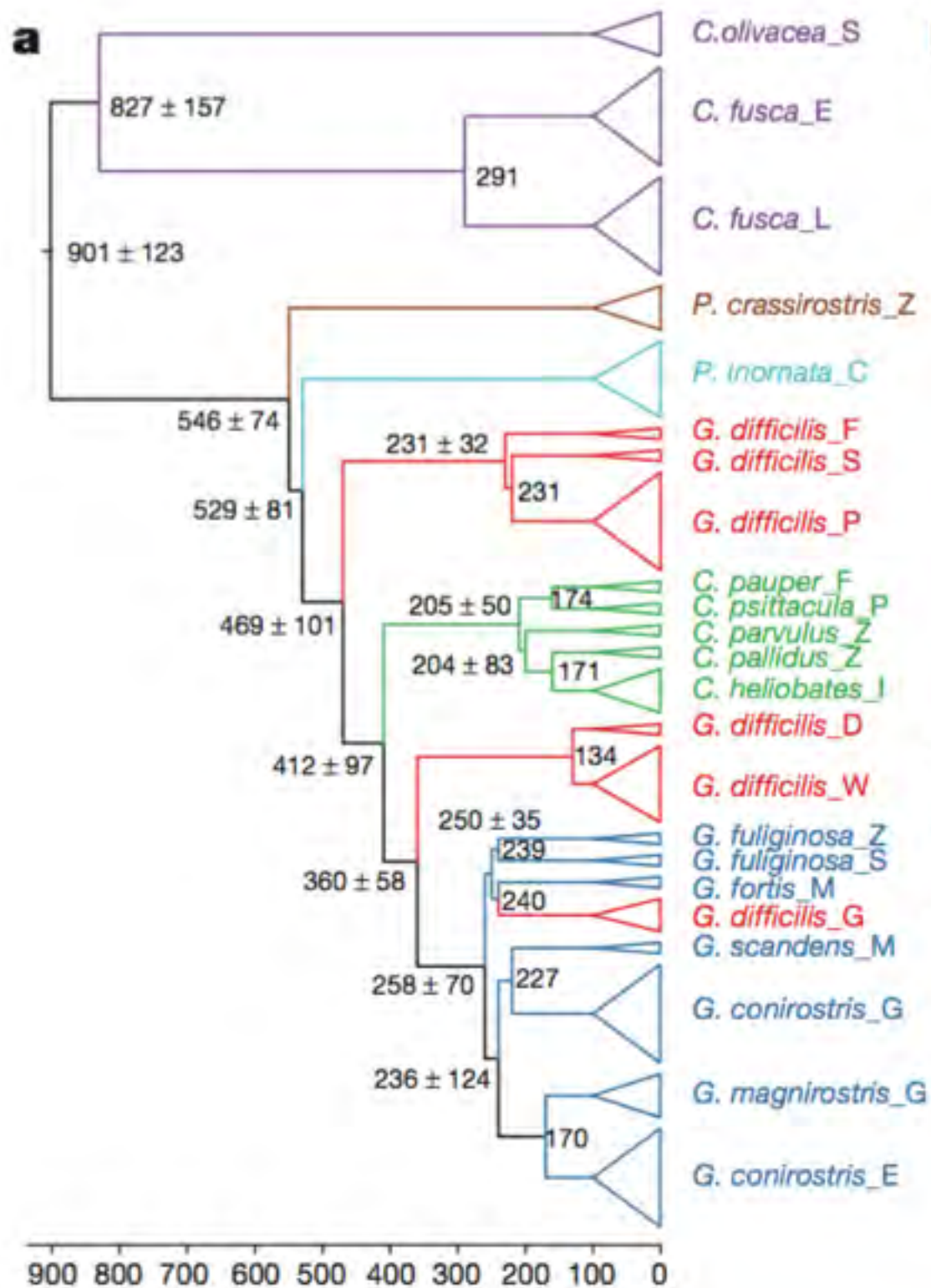


Chromosome



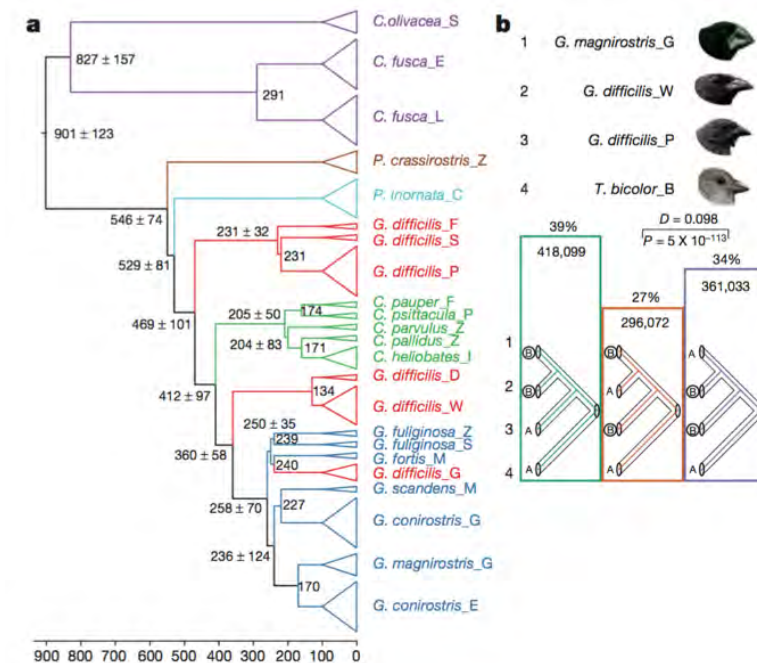
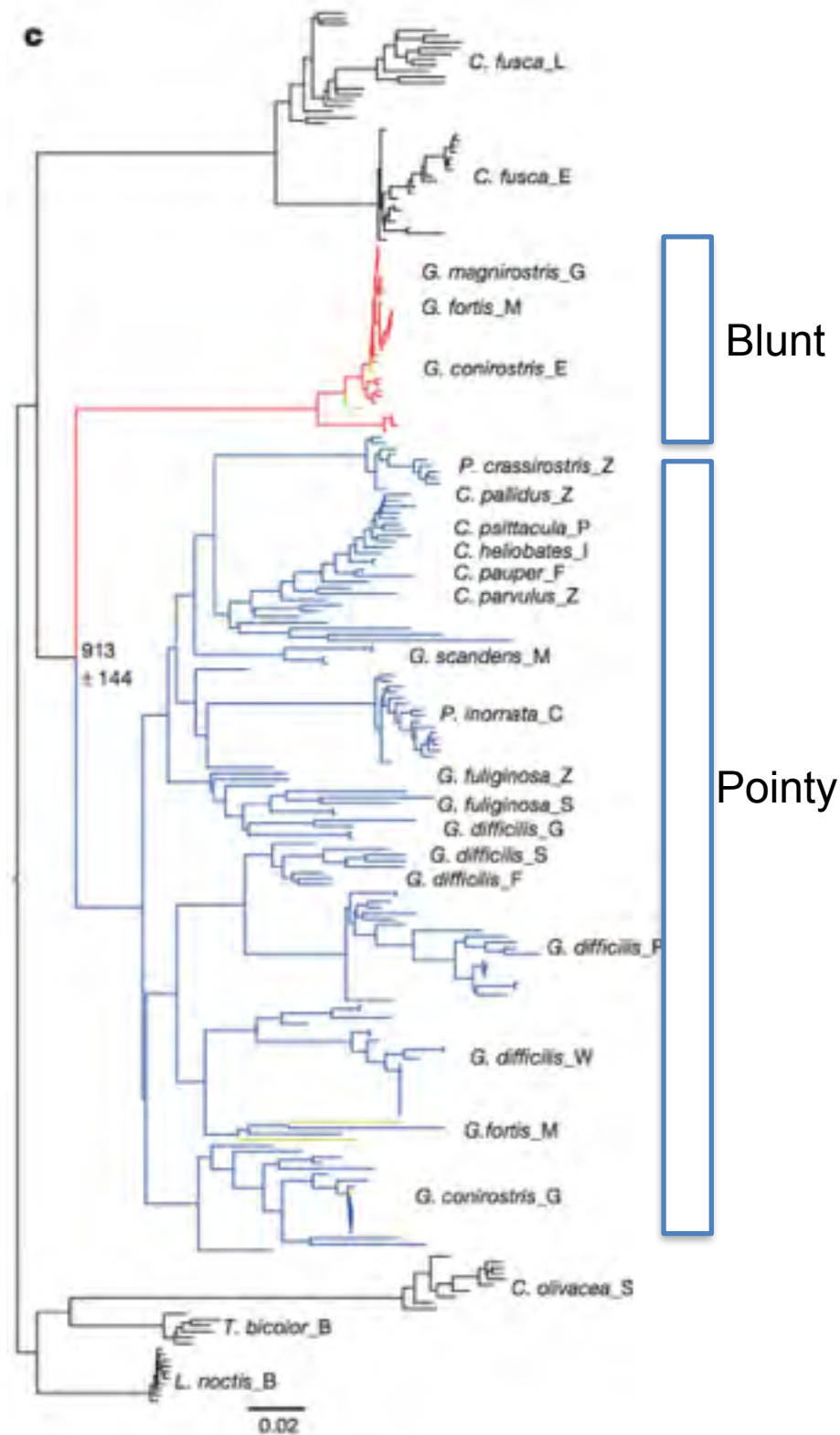








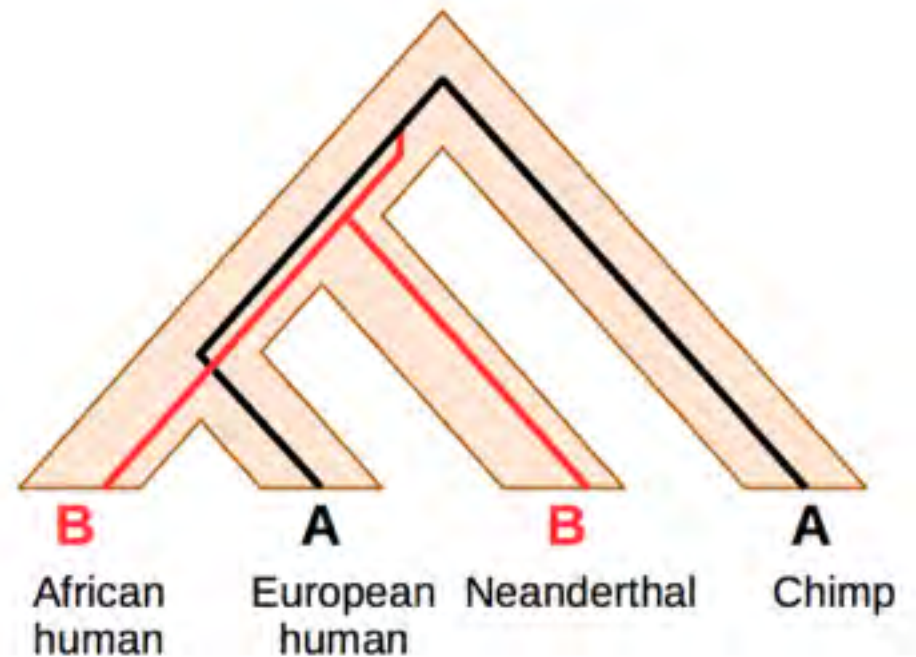
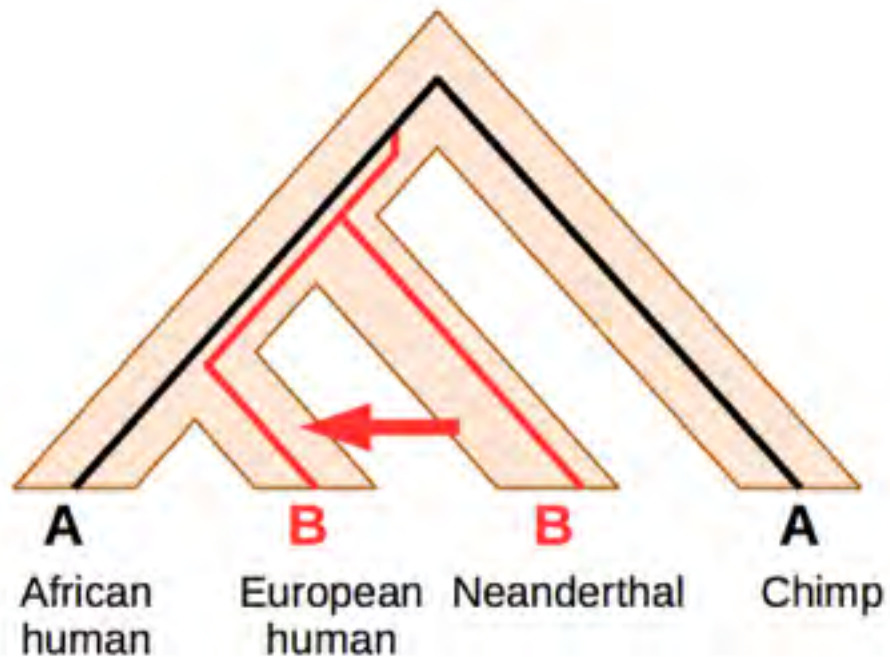
# ALX1 associated with beak shape



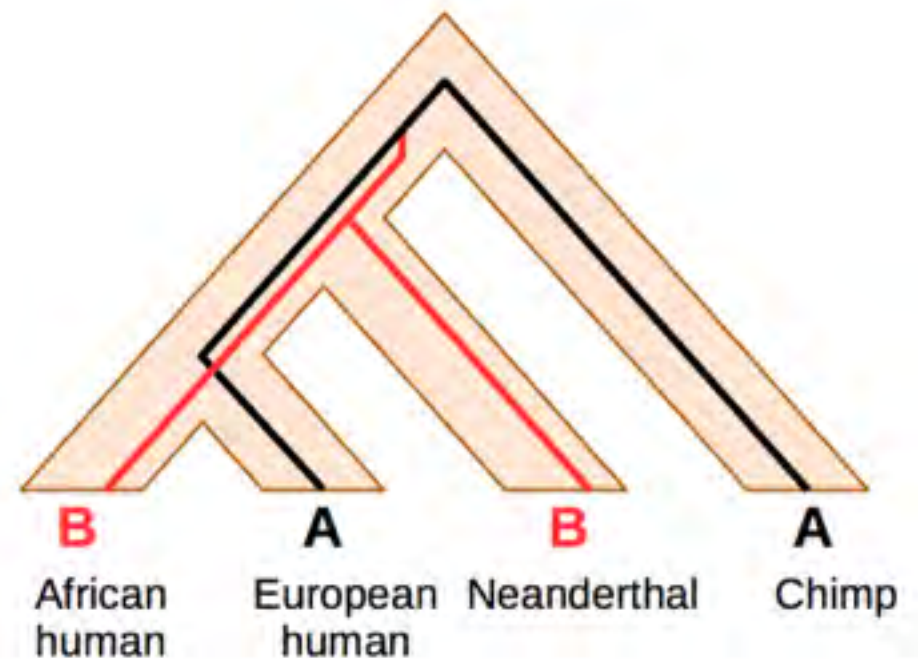
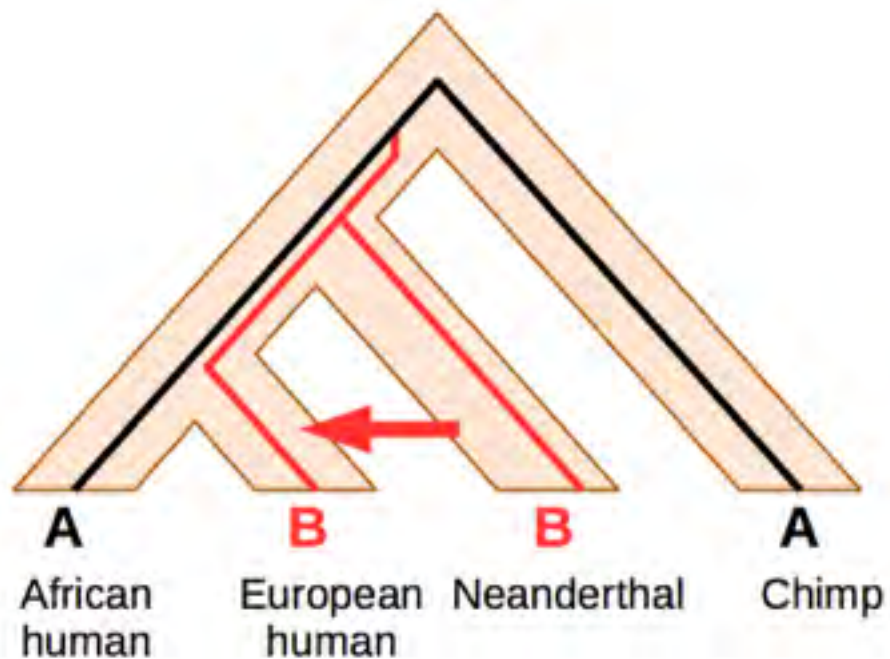
Most of these studies use phenotype associations to identify introgressed loci

But can we identify them a priori using the ABBA-BABA method?

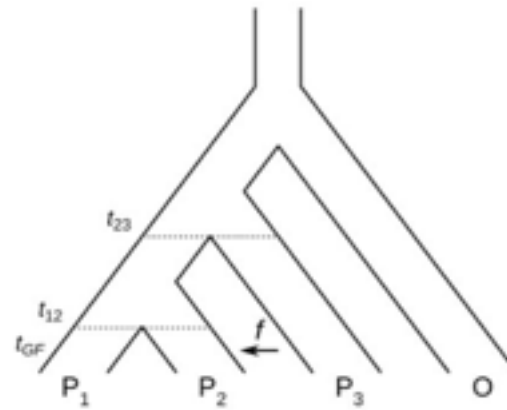
# Explicit tests for gene flow: ABBA-BABA test



# Explicit tests for gene flow: ABBA-BABA test

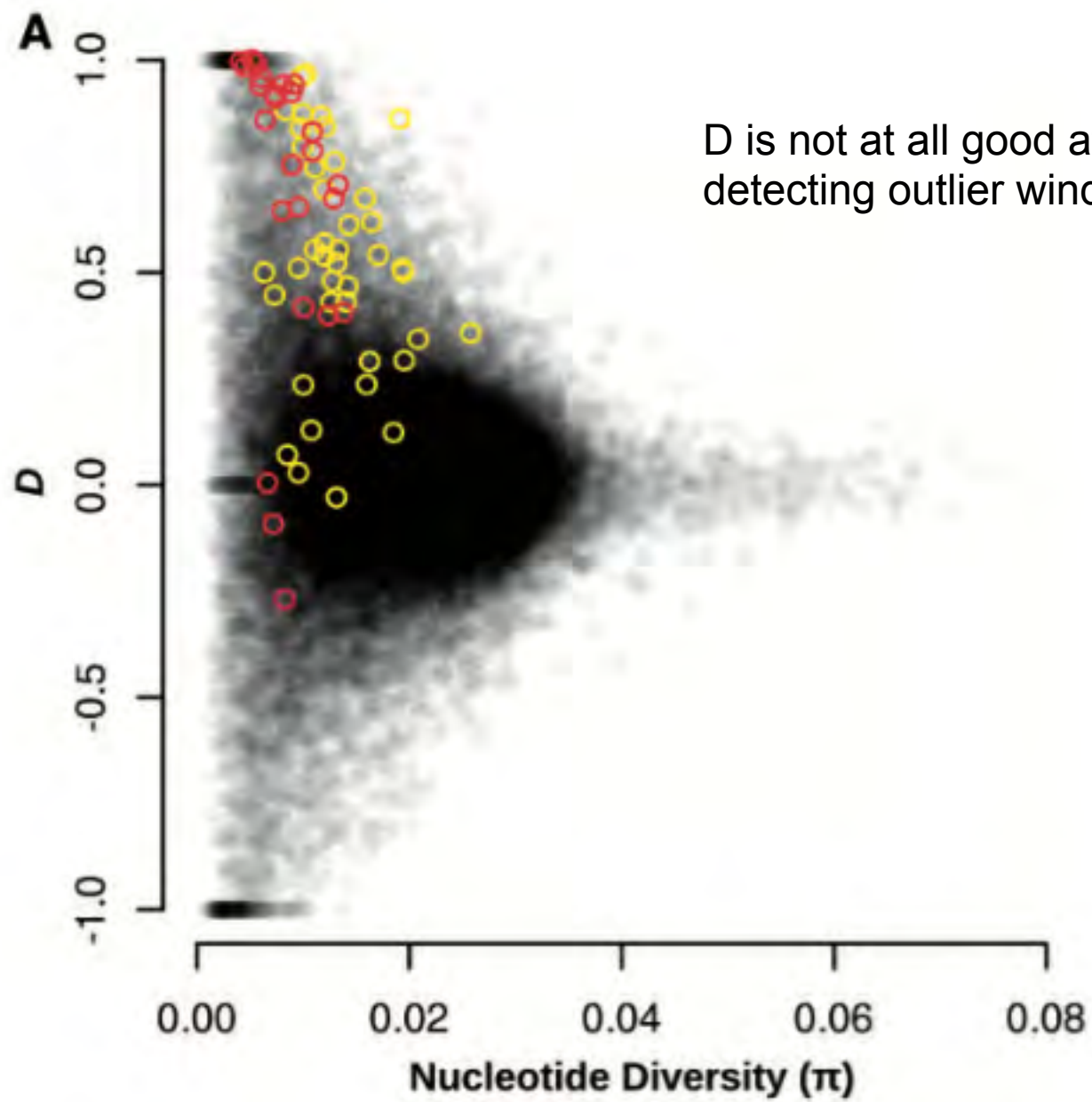




**A**

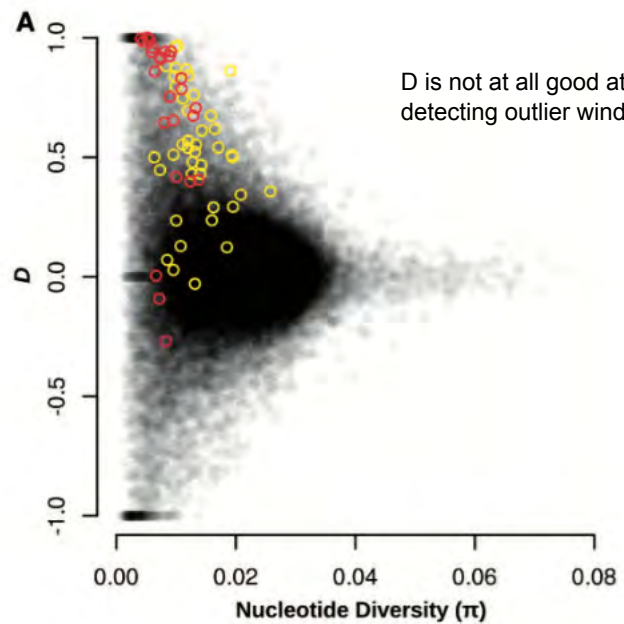
$$D(P_1, P_2, P_3, O) = \frac{\sum C_{ABBA}(i) - C_{BABA}(i)}{\sum C_{ABBA}(i) + C_{BABA}(i)} \quad (1)$$

D is quite dependent on the number of informative sites (the denominator)

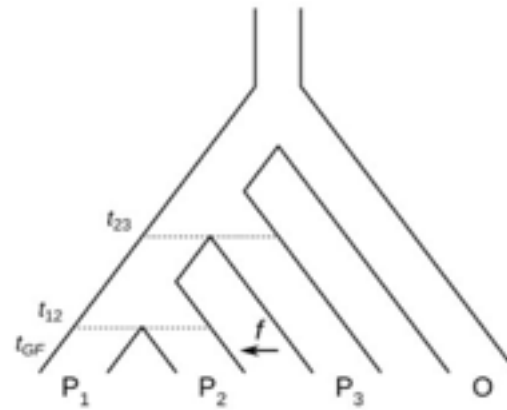


D is not at all good at  
detecting outlier windows

ites (the denominator)

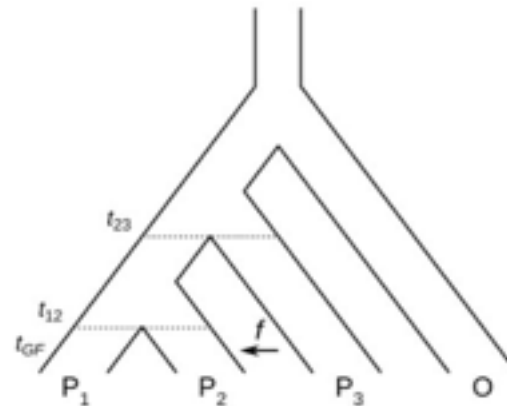
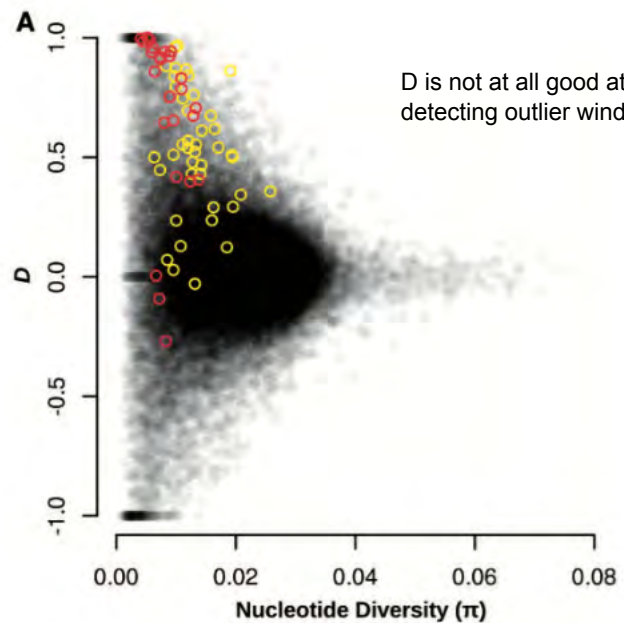


D is not at all good at  
detecting outlier windows



$$D(P_1, P_2, P_3, O) = \frac{\sum C_{ABBA}(i) - C_{BABA}(i)}{\sum C_{ABBA}(i) + C_{BABA}(i)} \quad (1)$$

D is quite dependent on the number of informative sites (the denominator)



$$D(P_1, P_2, P_3, O) = \frac{\sum C_{ABBA}(i) - C_{BABA}(i)}{\sum C_{ABBA}(i) + C_{BABA}(i)} \quad (1)$$

D is quite dependent on the number of informative sites (the denominator)

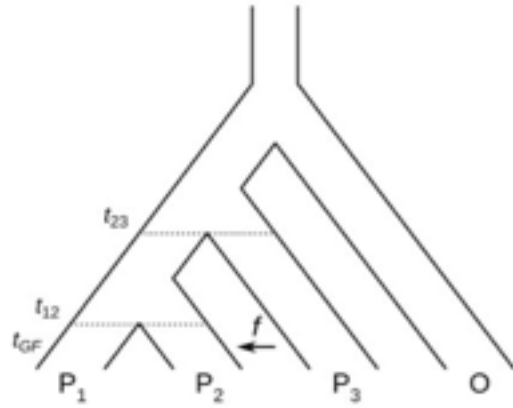
$$\hat{f}_G = \frac{S(P_1, P_2, P_3, O)}{S(P_1, P_{3a}, P_{3b}, O)}$$

Where s is numerator from the D equation

f is the fraction of introgression compared to maximum possible



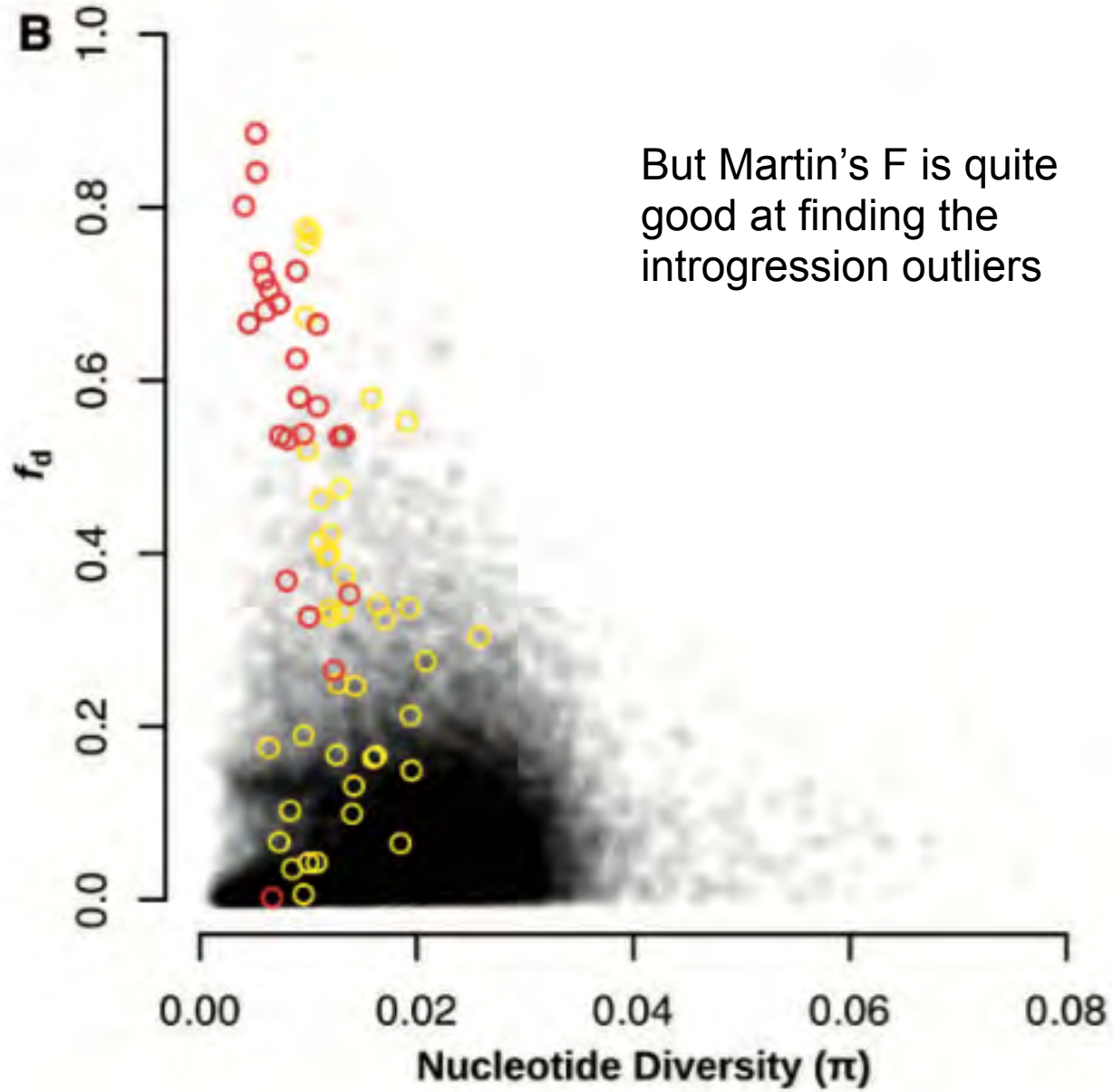
**A**



## Martin's F

$$\hat{f}_d = \frac{S(P_1, P_2, P_3, O)}{S(P_1, P_D, P_D, O)}$$

A





## • **Do *Heliconius* butterfly species exchange mimicry alleles?**

Joel Smith, Marcus R. Kronforst

Published 17 July 2013. DOI: 10.1098/rsbl.2013.0503

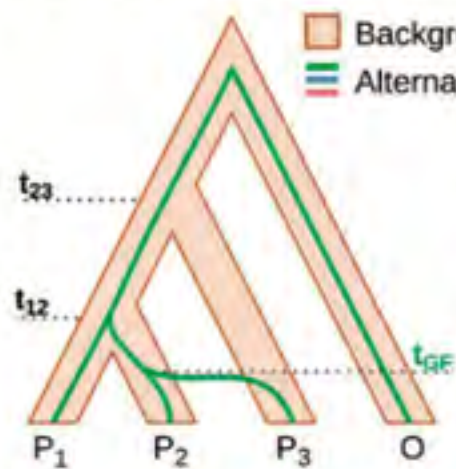
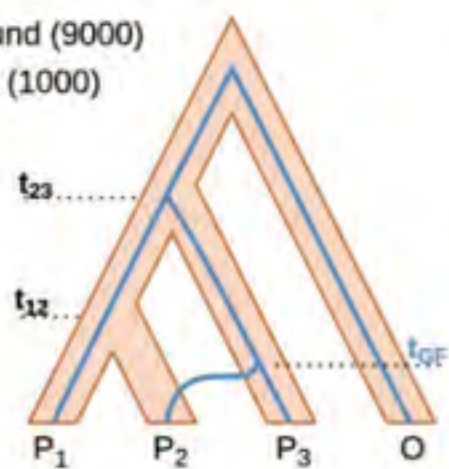
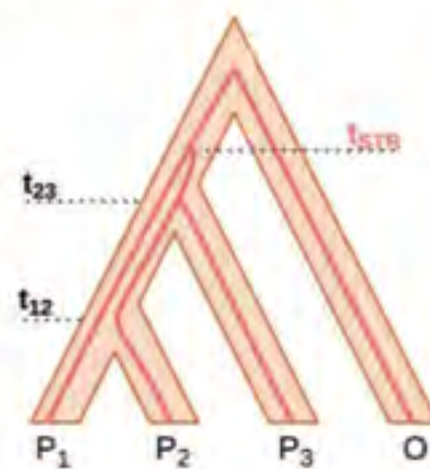
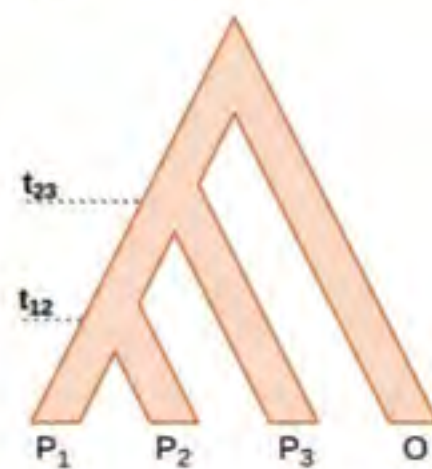
[Article](#) [Figures & Data](#) [Info & Metrics](#) [eLetters](#)

 PDF

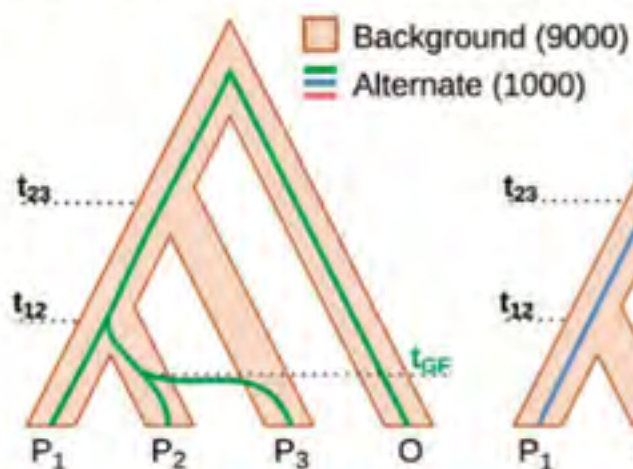
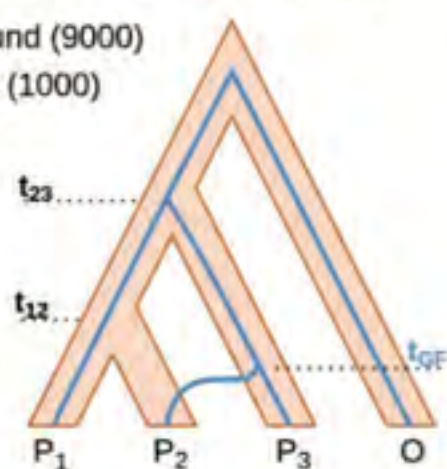
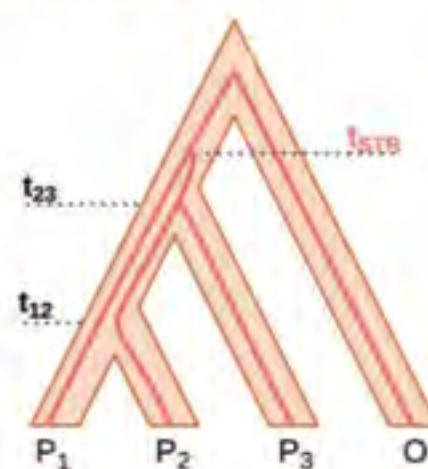
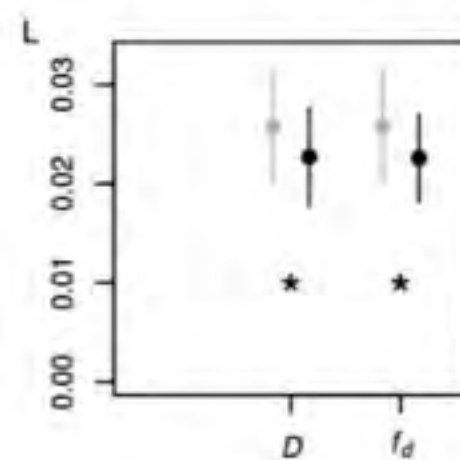
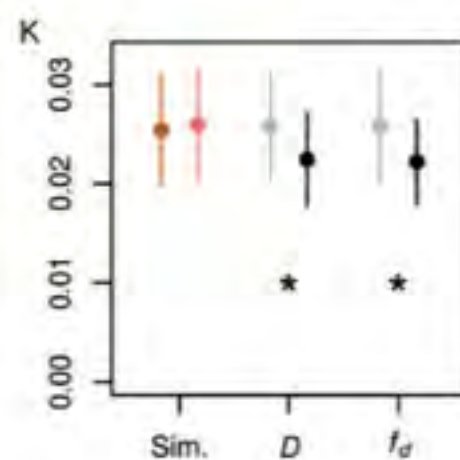
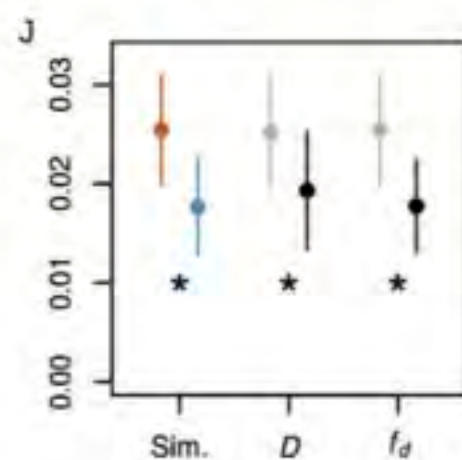
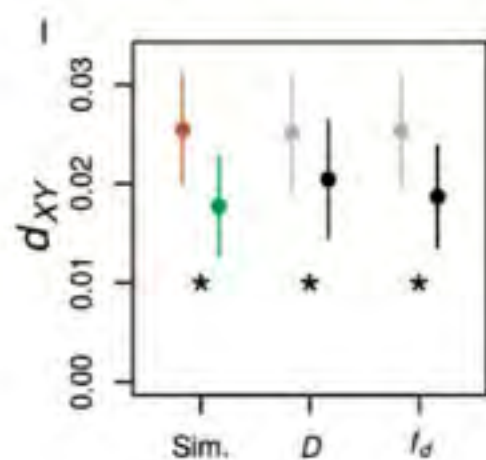
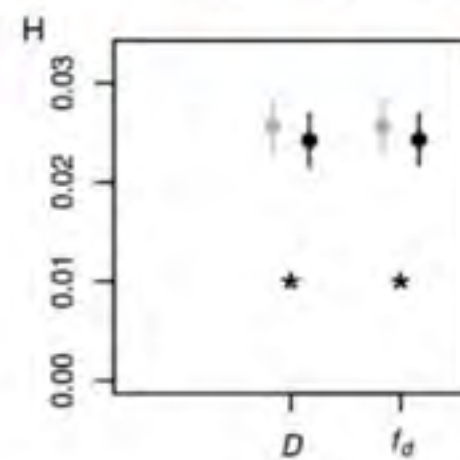
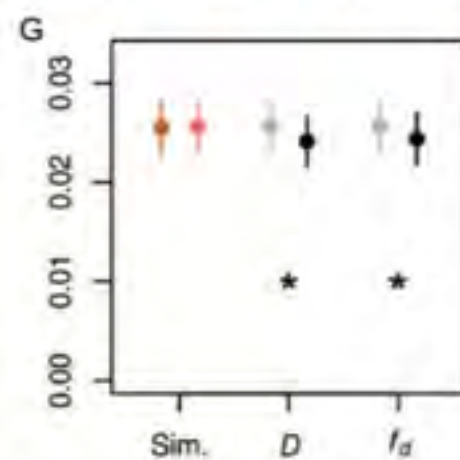
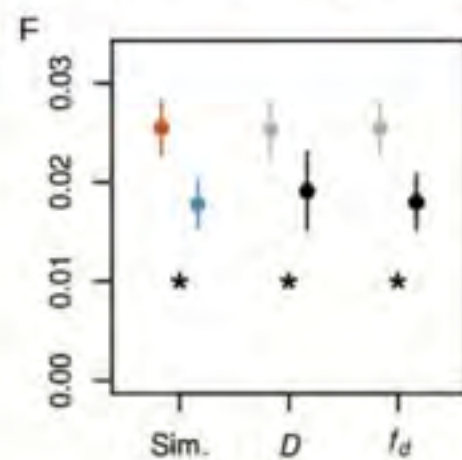
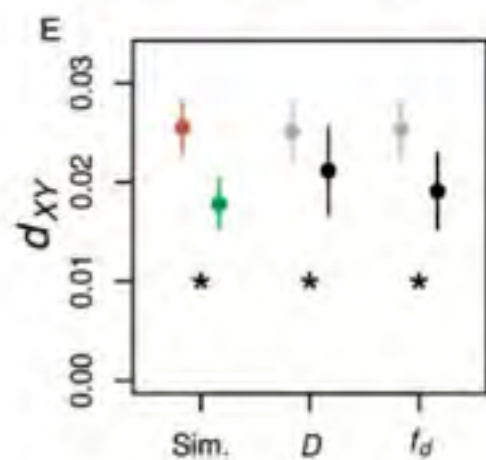
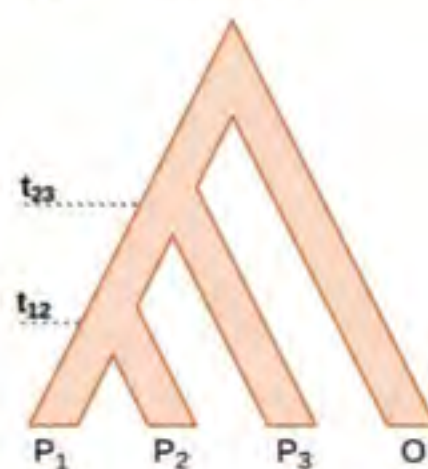
 Previous

Next 

- Smith and Kronforst argued that introgression could be inferred where ABBA-BABA outliers showed lower Dxy compared to genome-wide average

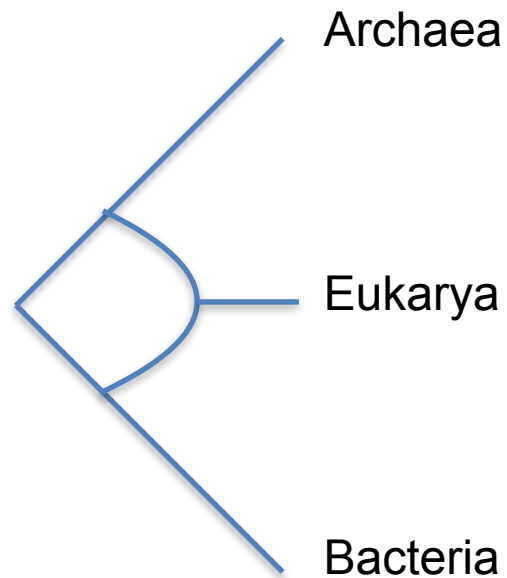
**A** Gene Flow  $P_2$  to  $P_3$ **B** Gene Flow  $P_3$  to  $P_2$ **C** Ancestral Structure**D** Null Model



**A** Gene Flow  $P_2$  to  $P_3$ **B** Gene Flow  $P_3$  to  $P_2$ **C** Ancestral Structure**D** Null Model

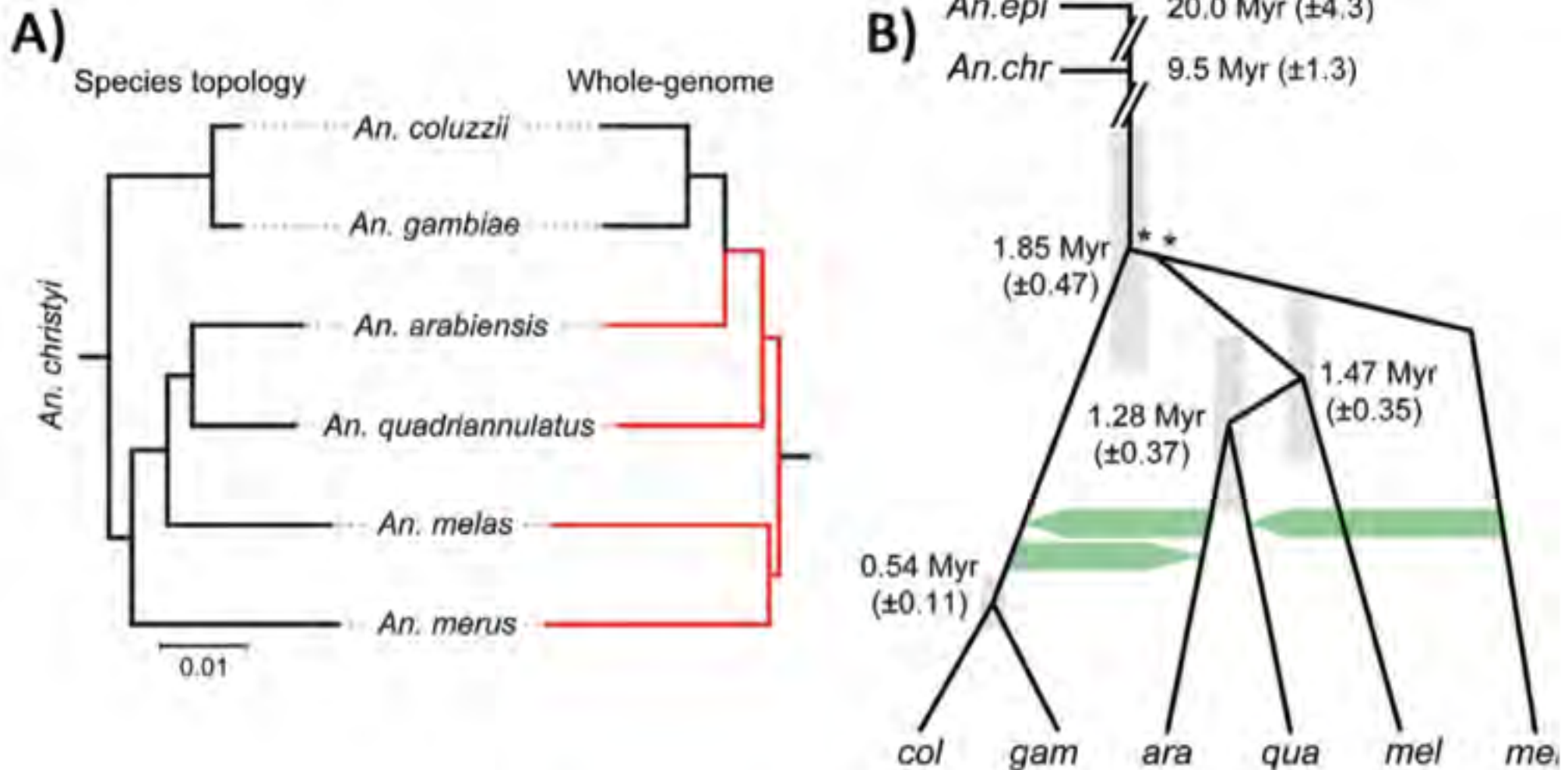
- Be wary of window based D statistics
- F is better than D...
- Sampling design is very important!

# Implications for tree-thinking



**The tree of life is reticulated**

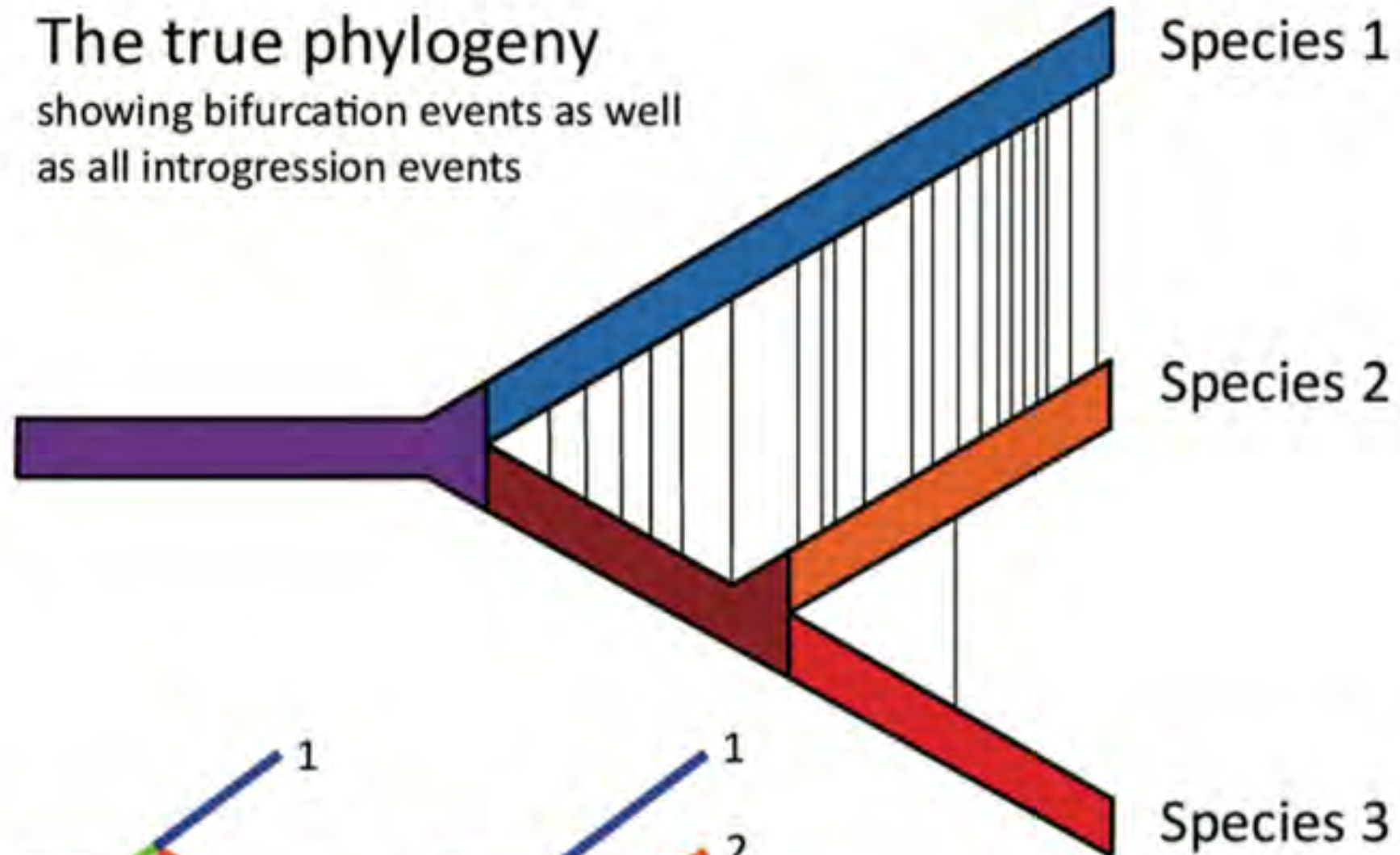
# Implications for tree-thinking





# The true phylogeny

showing bifurcation events as well as all introgression events



'Whole-genome'  
or 'democratic  
majority' tree



'Species tree,' or  
bifurcation history

Okay, so what have we learnt  
and where do we go from here?