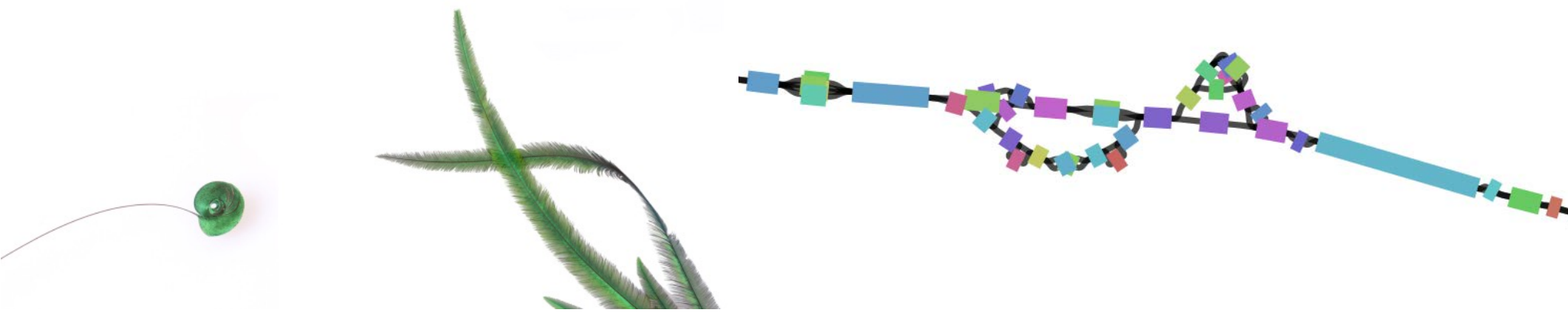




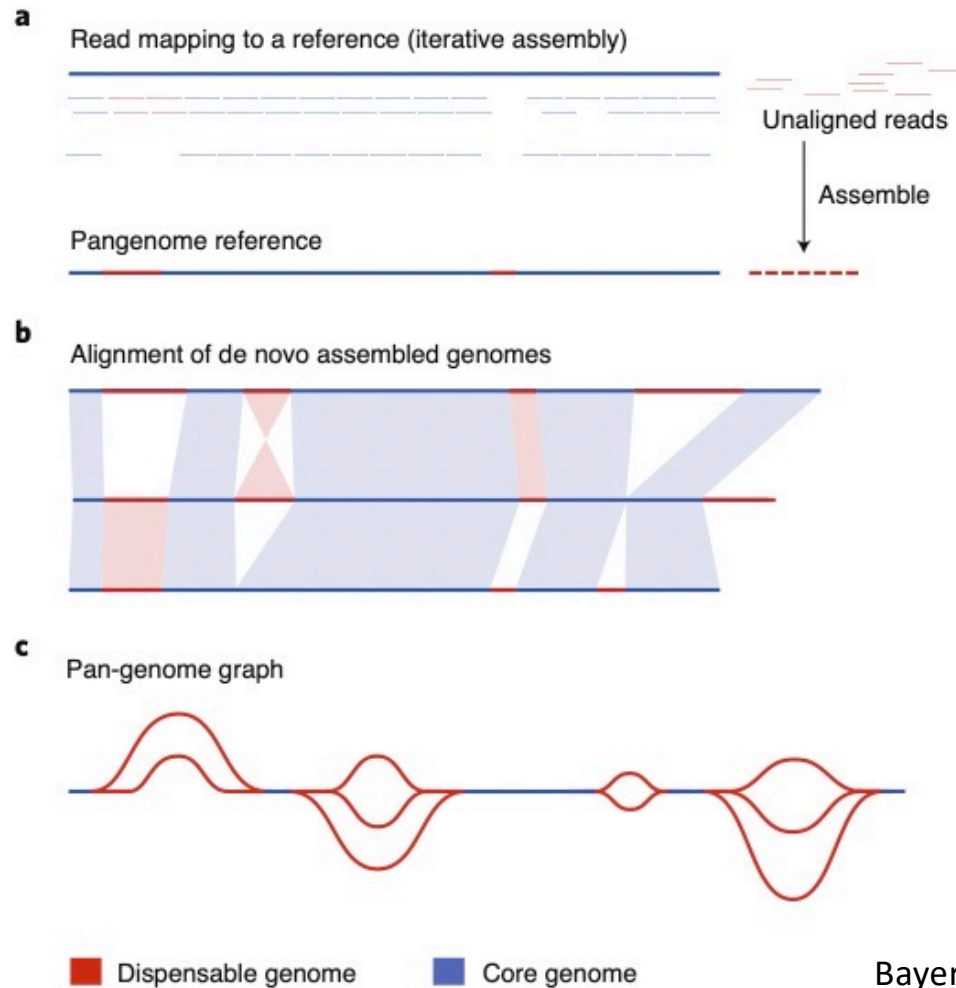
# Pangenomes as a new tool for studying ecology and evolution of natural populations

Scott V. Edwards

Museum of Comparative Zoology, Harvard University, Cambridge, USA



# Pangenomes: moving beyond reference-based genomics



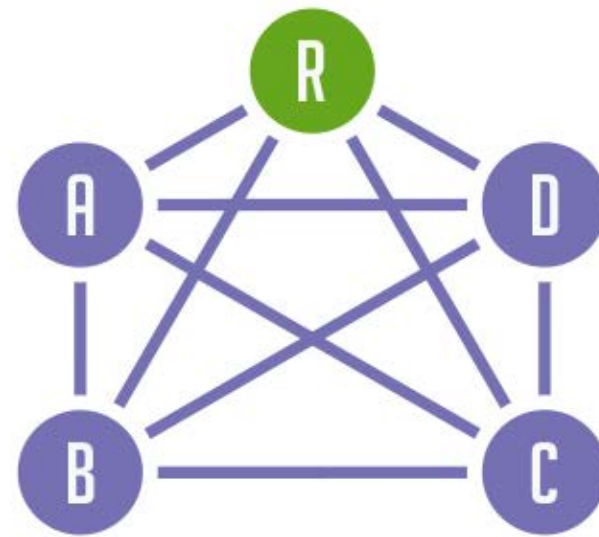
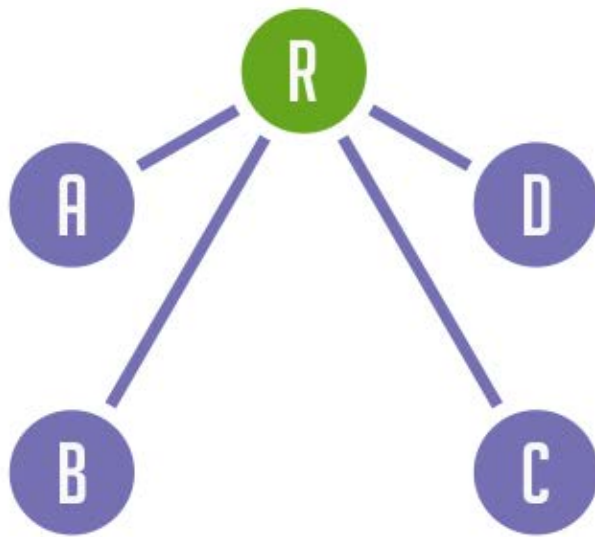
Bayer et al. 2020. *Nature Plants* 6: 914-920.

# Reference-free genomics

Genomic

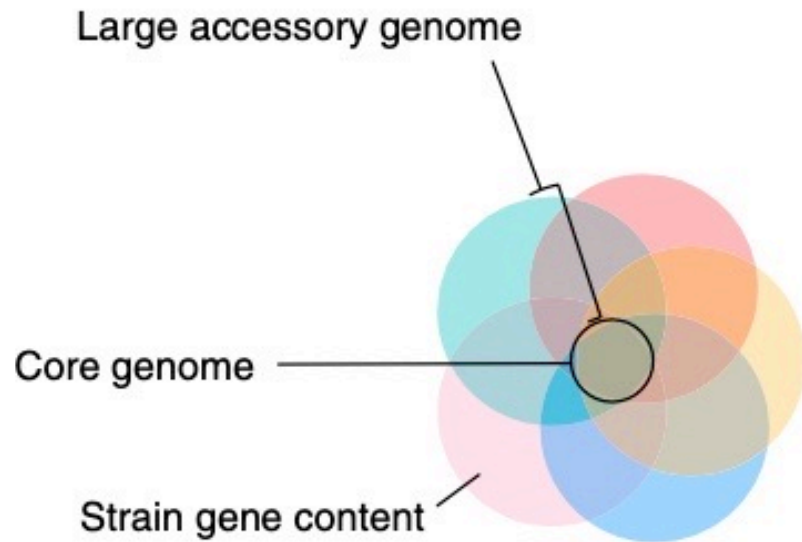
Pangenomic

Reference model

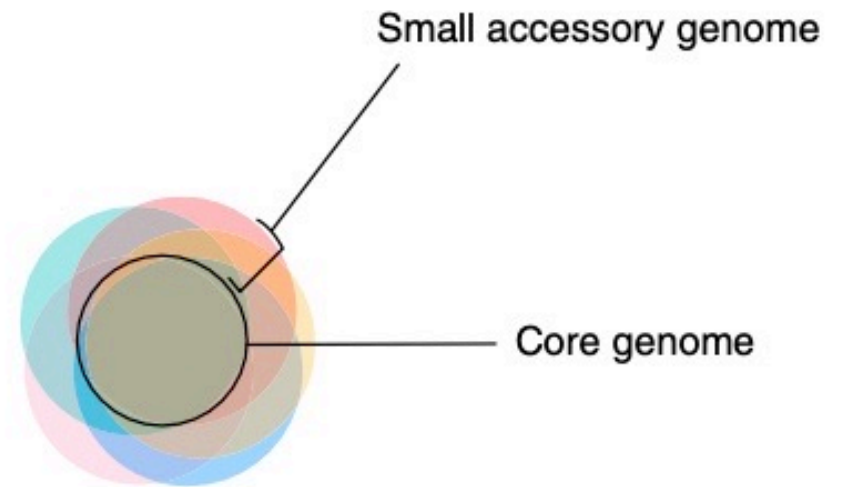


# Open and closed pangenomes

Open pangenomes

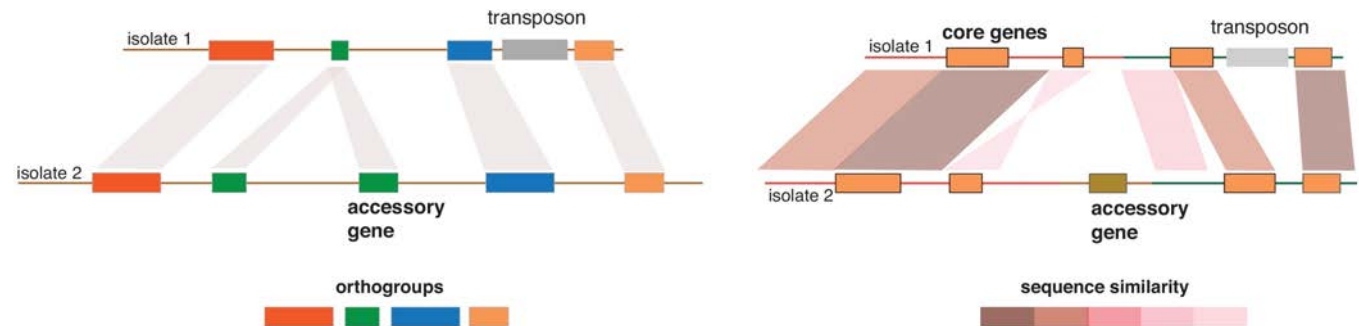
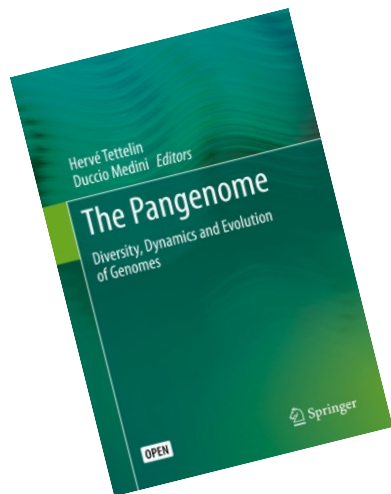


Closed pangenomes



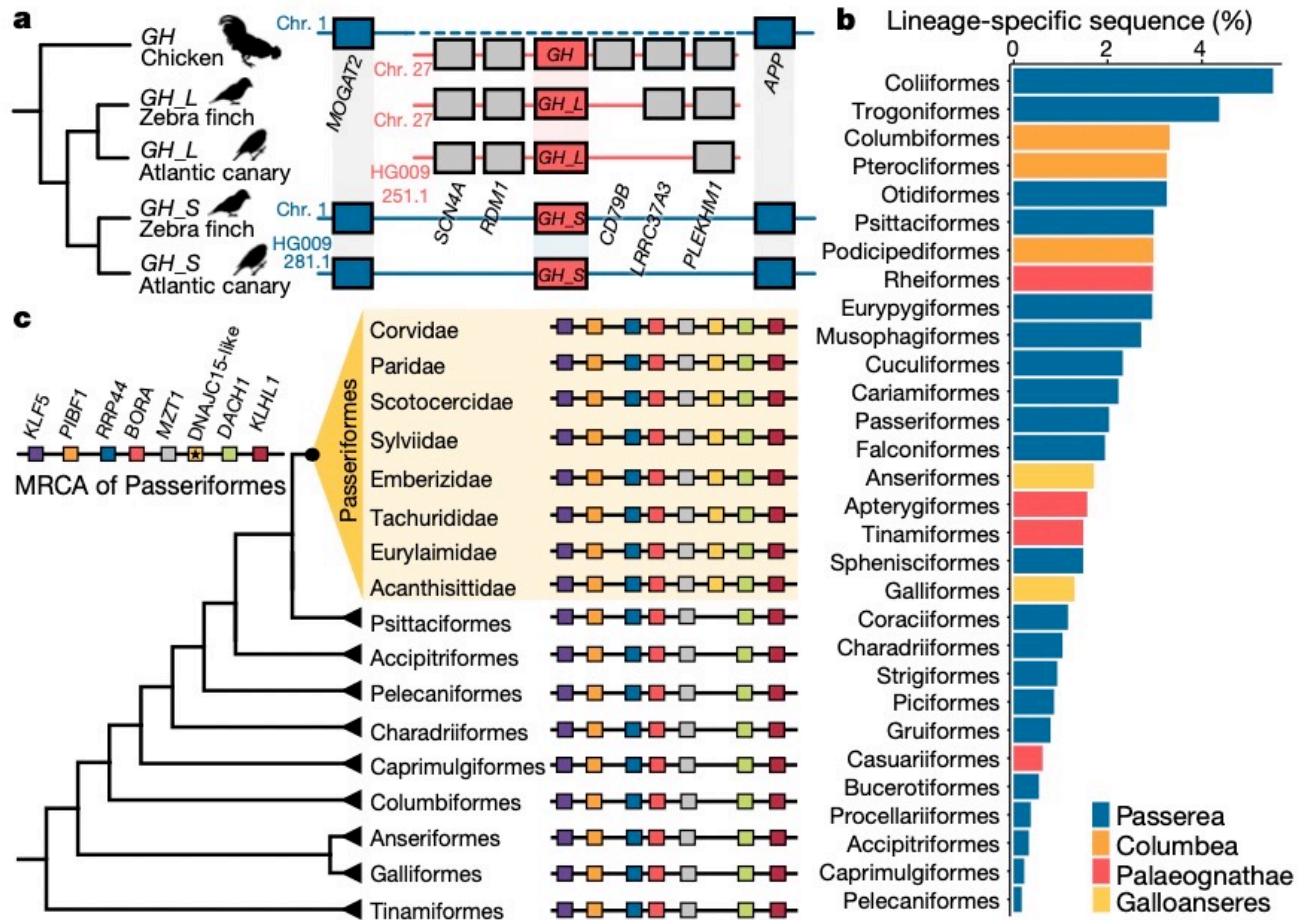
# The eukaryotic pangenome

- “The existence of pangenomes in eukaryotes is debated...Pangenome studies in eukaryotes are challenging due to their more complex genome and architectures and a lack of replete genome-level sampling” (Brockhurst et al. 2019. *Current Biology*)



<https://pathogen-genomics.org/research/>

# Pangenome approach to comparative genomics



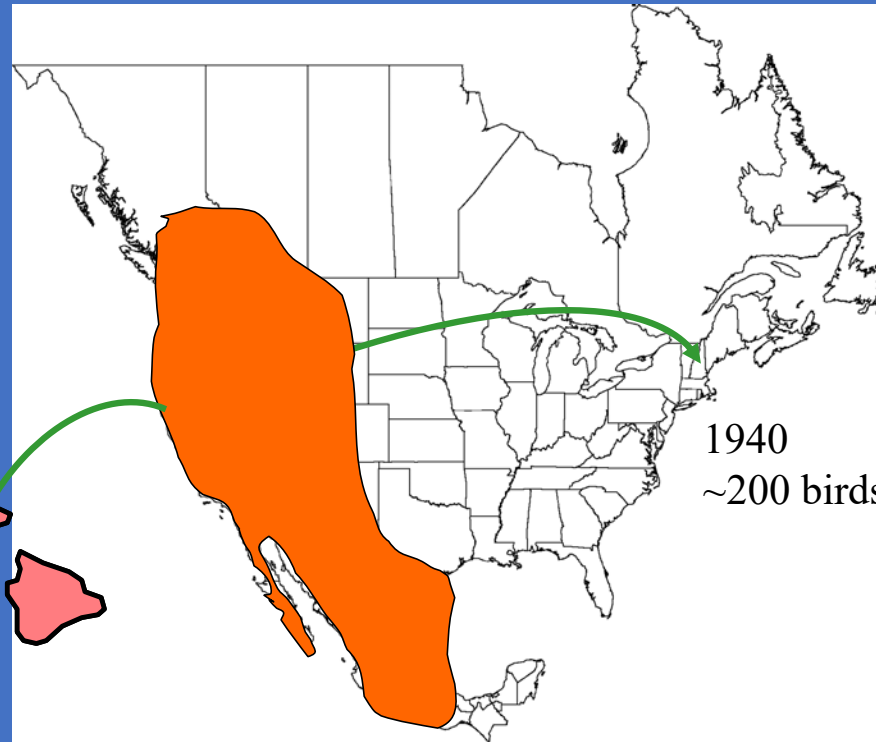
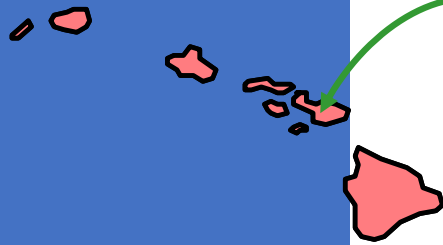
Feng et al. 2020. *Nature* 587:252-257.



# Recent history of House Finch populations

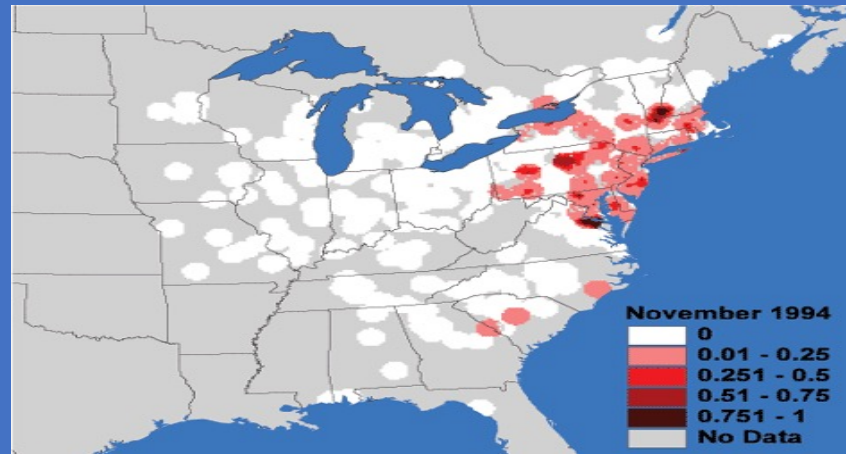
■ historic range

~1870 bottleneck?



1940  
~200 birds

# Rapid spread of *Mycoplasma* in House Finch populations

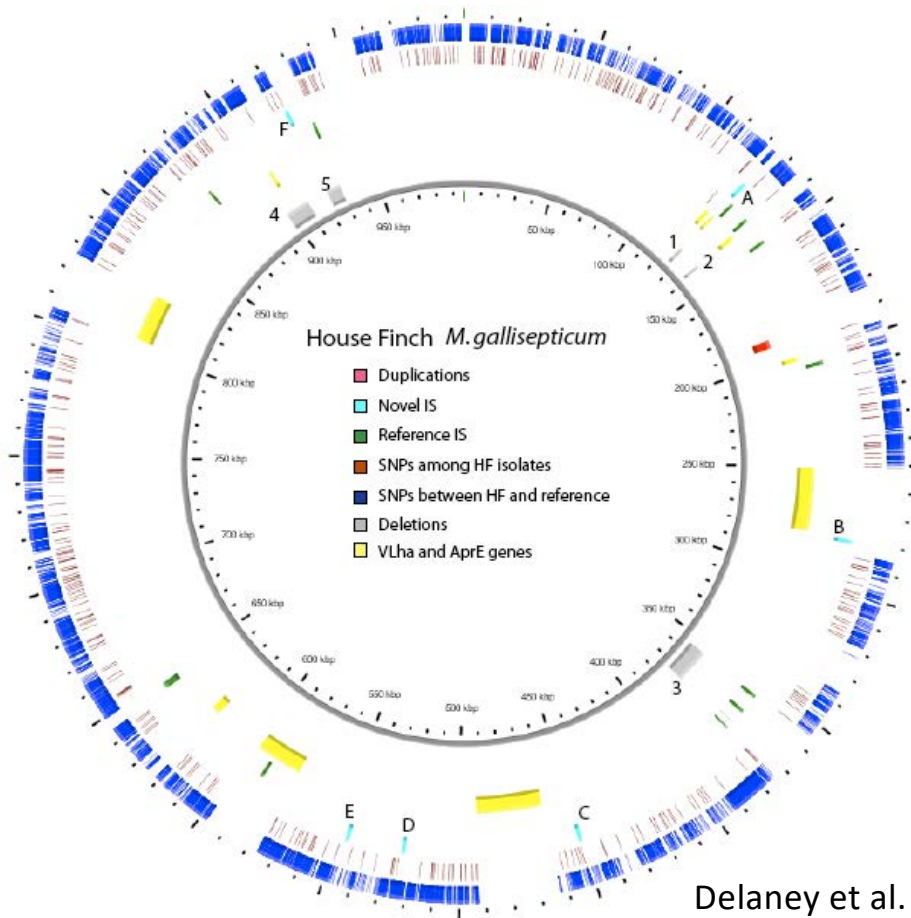


Courtesy Cornell Lab of Ornithology

- *Mycoplasma* is transmitted horizontally, often at bird feeders
- Expanded throughout the eastern US in just five years
- Has now crossed the Rockies and is spreading south through California and the southwest.



# House Finch *Mycoplasma* genome ~1 Mb



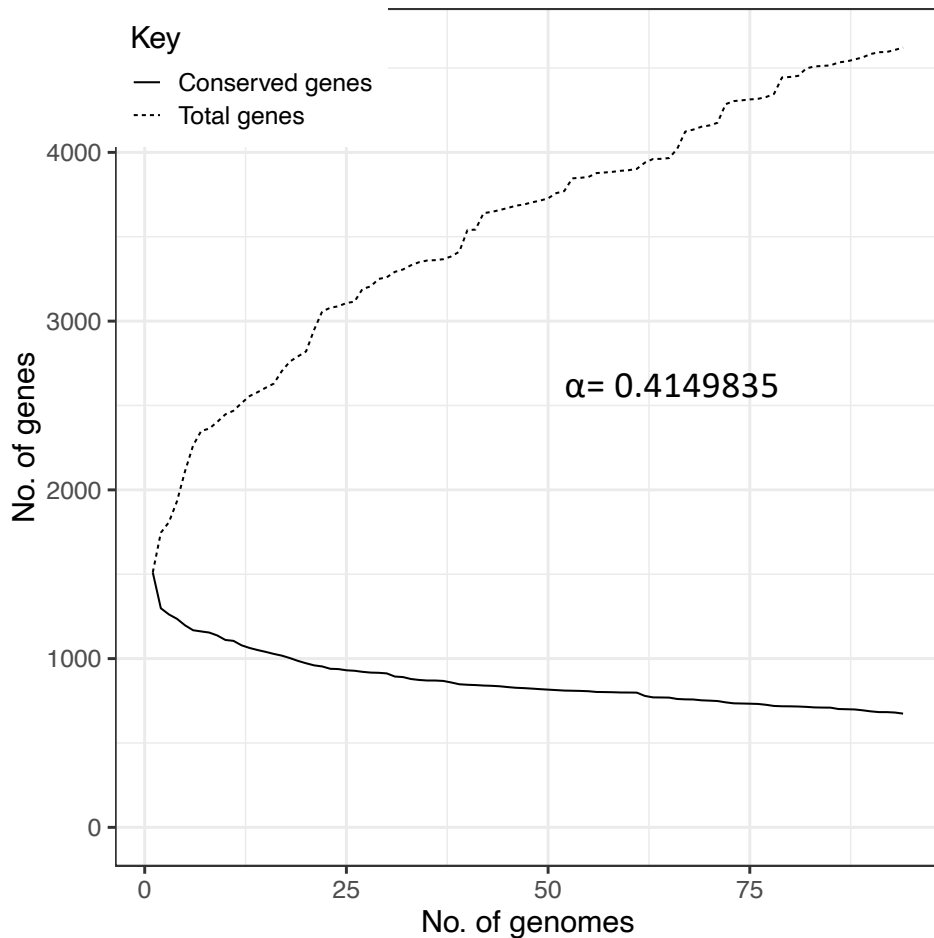
Delaney et al. 2012. *PLoS Genetics*

Analyzed 81 *Mycoplasma* strains from chicken, turkey and house finch, available on NCBI

Added 12 new House Finch *Mycoplasma* strains, sequenced with PacBio

Used

# Pangenome of *Mycoplasma gallisepticum*



The size of the pan-genome was determined using 10,000 permutations by microPan

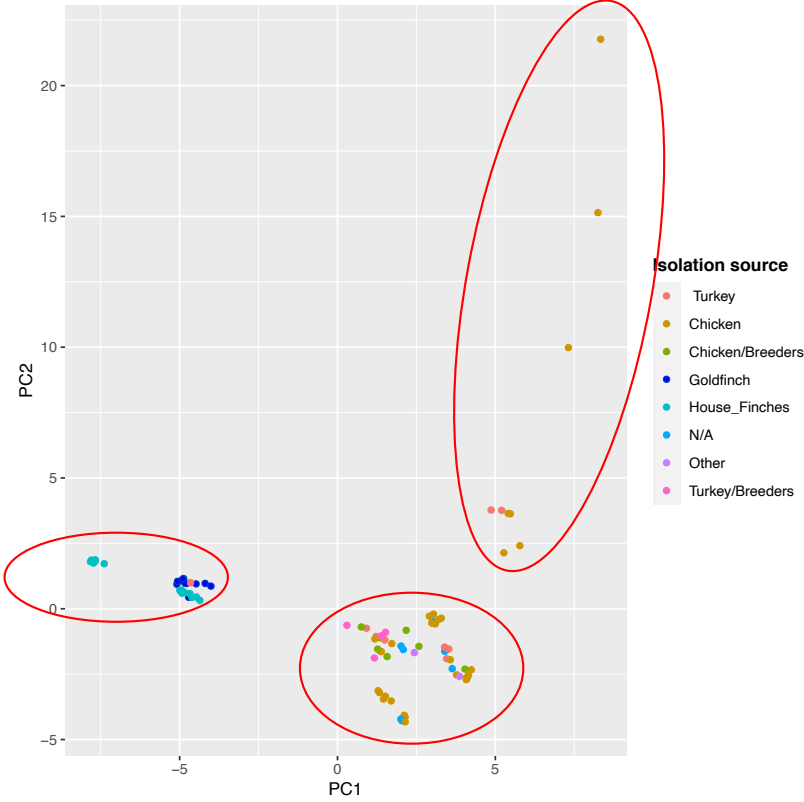
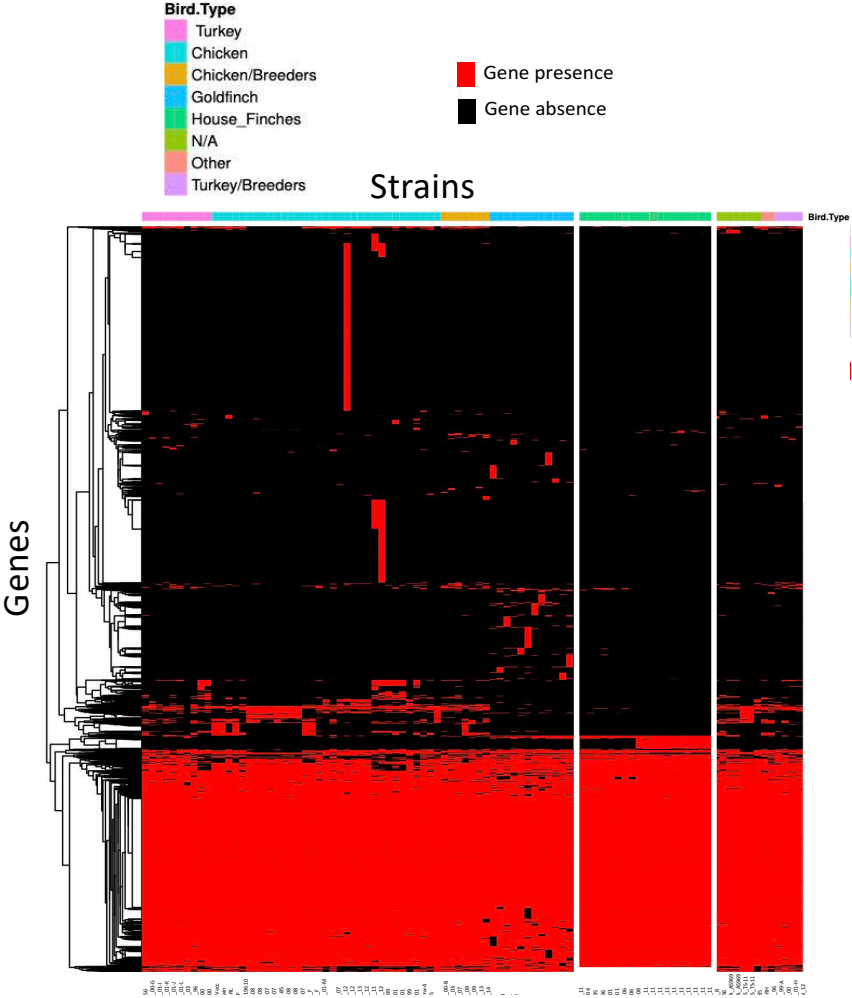
Feature	Info	Number of genes	Percentage
Core genes	(99% <= strains <= 100%)	674	14.586
Soft core genes	(95% <= strains < 99%)	464	10.041
Shell genes	(15% <= strains < 95%)	412	8.916
Cloud genes	(0% <= strains < 15%)	3071	66.457
SGF	one copy in all strains	141	3.051
SGF	without recombination signals	117	2.532
Total genes	(0% <= strains <= 100%)	4621	100

Alpha value: the number of gene clusters we would see if we collected *all* genomes of the species

**New data: Determine the alpha value using MicroPan**

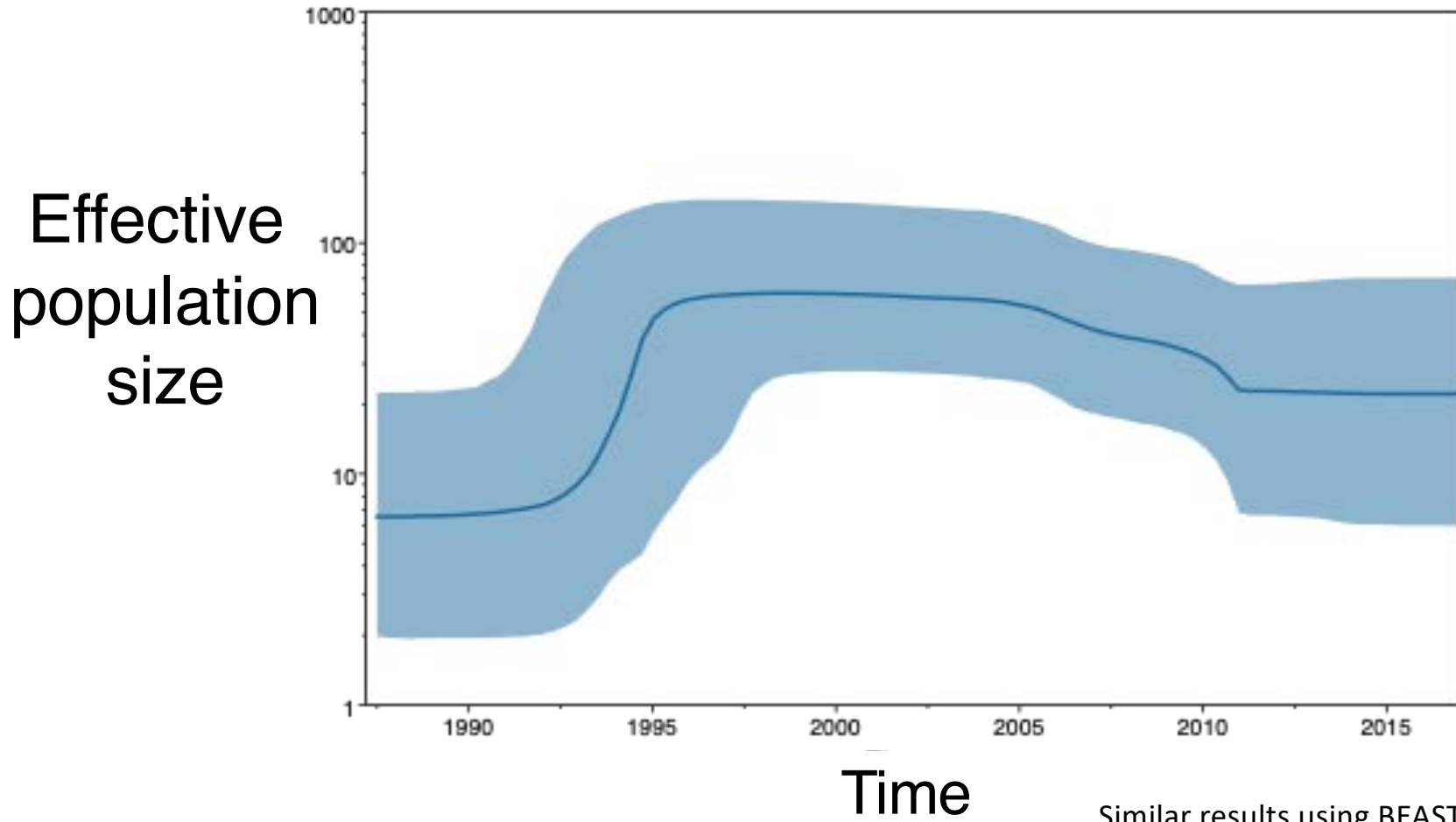
**\*the pan-genome is closed if the estimated alpha is above 1.0**

# Mycoplasma pangenome gene repertoire is highly strain-specific



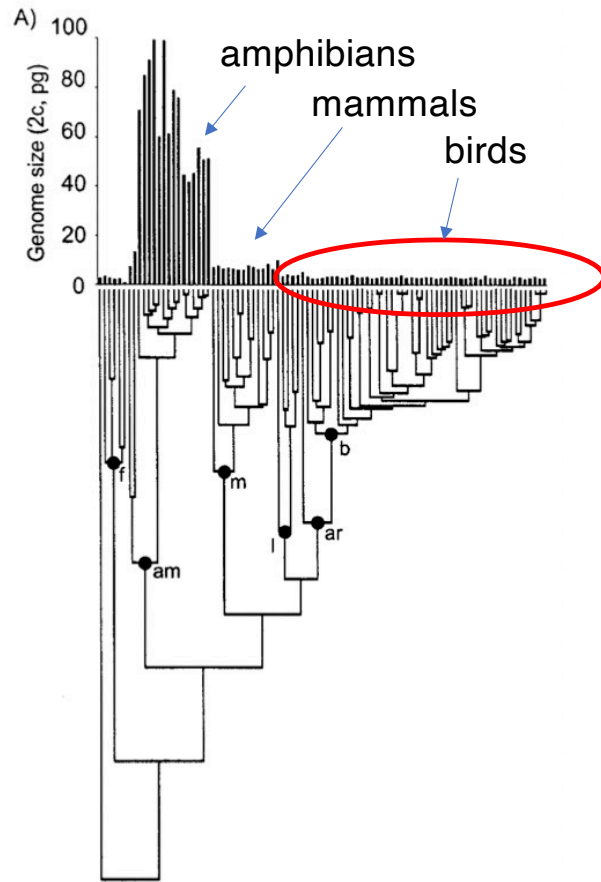


*Mycoplasma* epizootic likely began ~2 years before first detection

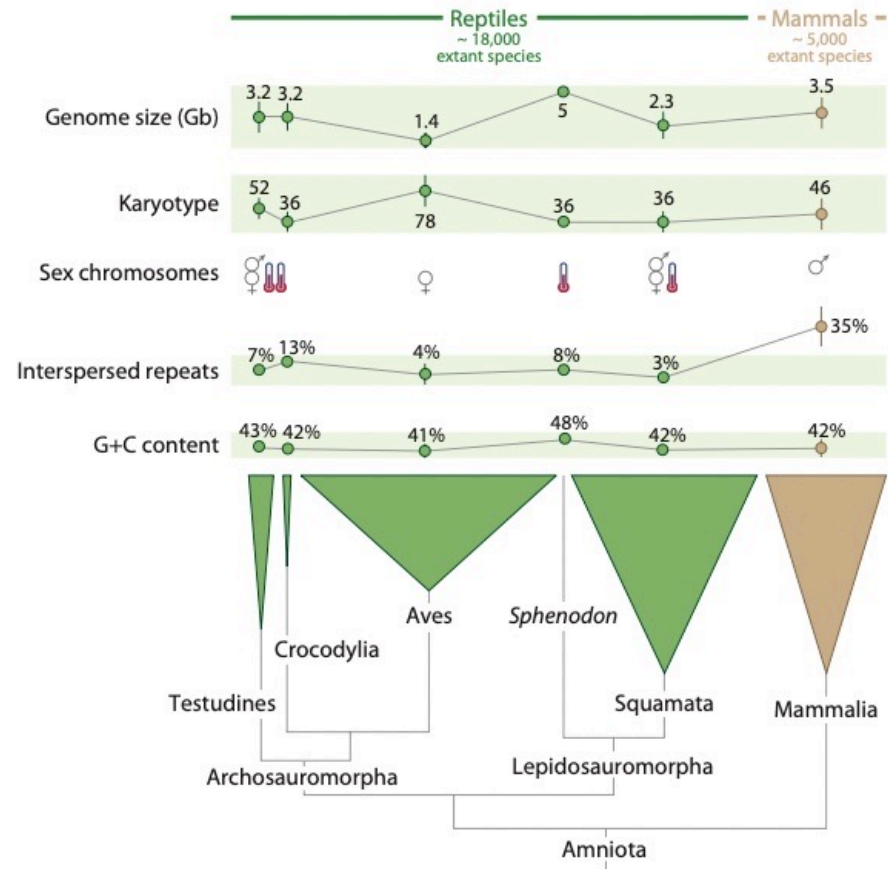


Similar results using BEAST and Stairway plot

# Birds have small, streamlined genomes

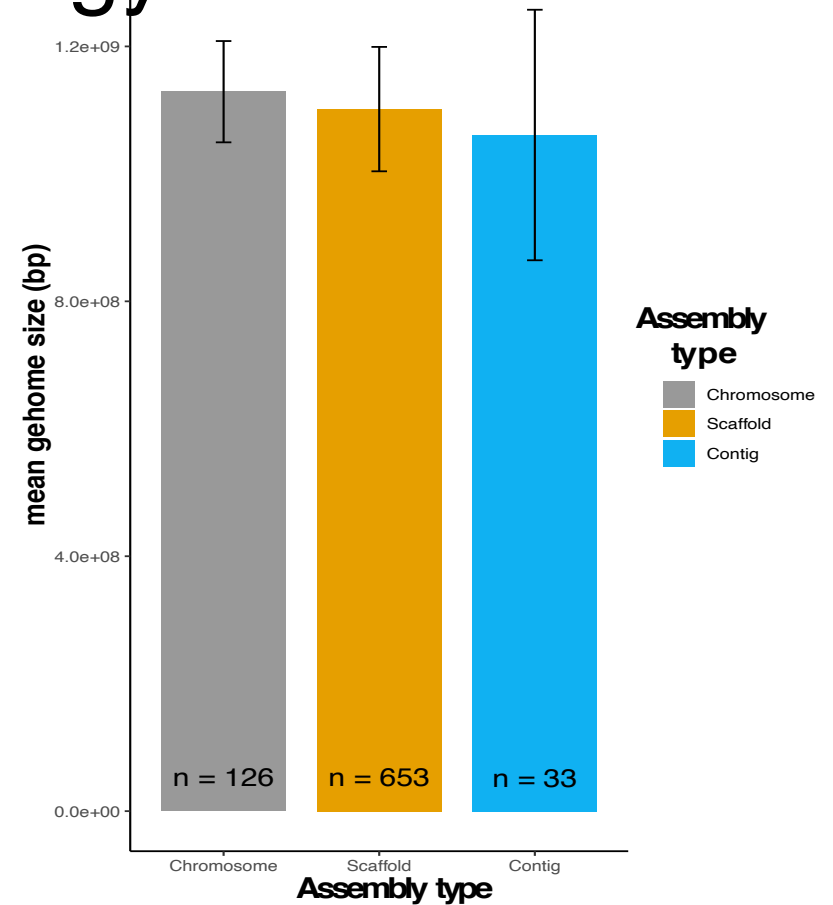
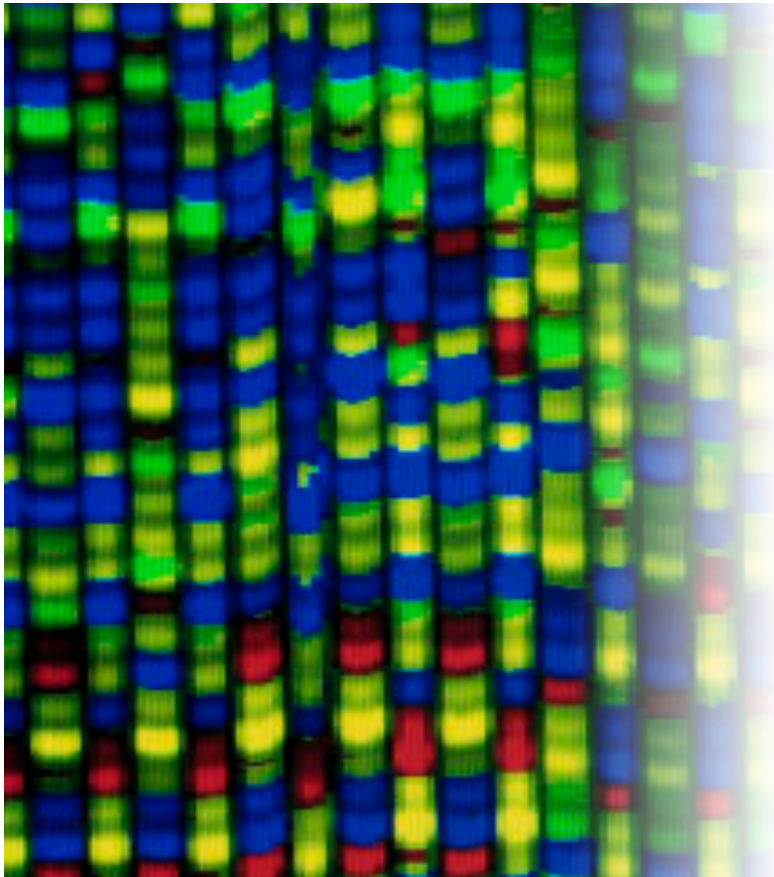


Waltari & Edwards. 2002. *Am. Nat.*



Organ et al. 2010. *Ann. Rev. Genom. Hum. Genet.*

# Avian genomes are growing with each new technology



Data from NCBI, accessed 13 Nov. 2021

# Three scrub-jay (*Aphelocoma*) species in pangenome project



n = 14



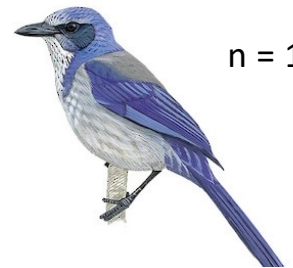
Woodhouse's scrub jay  
*A. woodhouseii*  
weight 76.9 -77.7 g

- Goal: study genome complexity and estimate fitness effects of structural variation



Florida scrub jay  
*A. coerulescens*  
weight 75 – 79.3 g

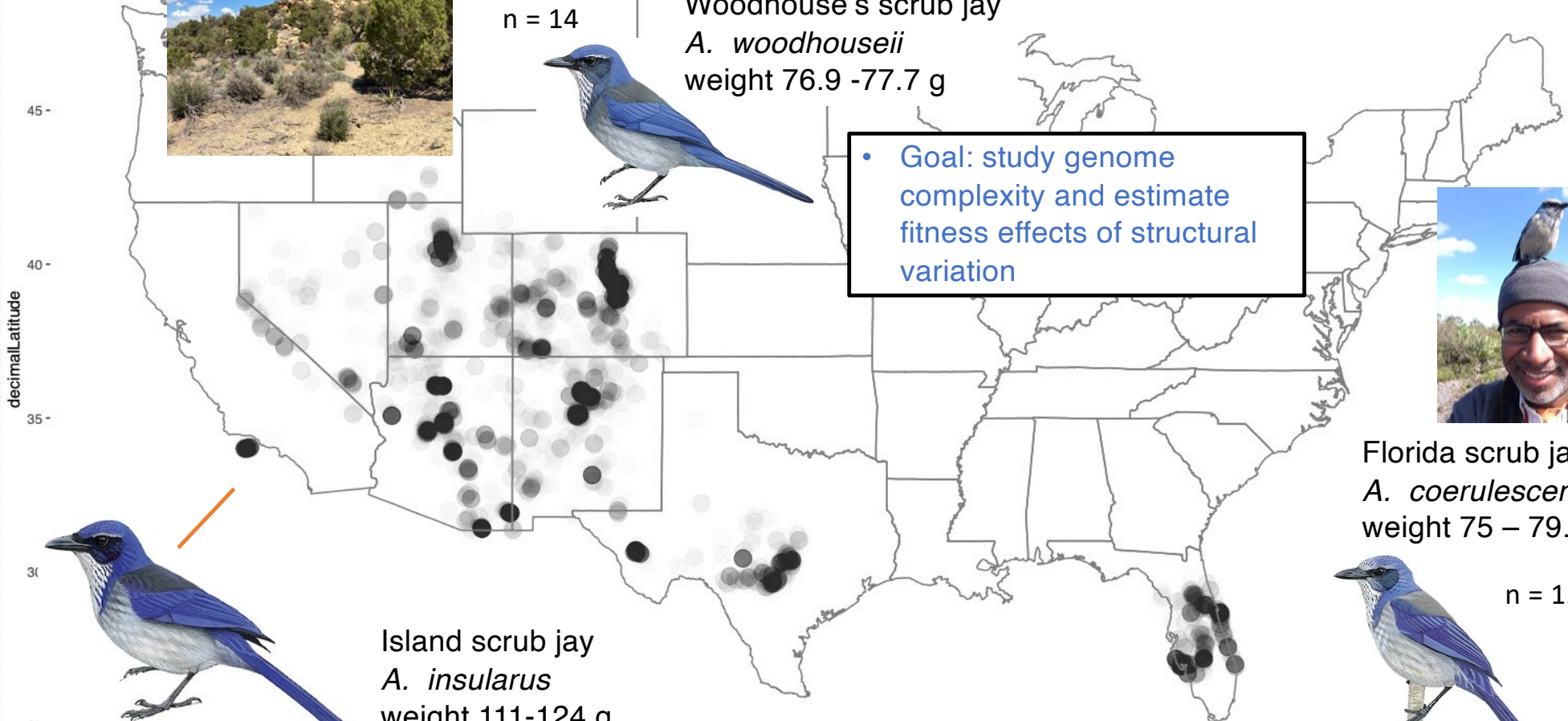
n = 15



Island scrub jay  
*A. insularis*  
weight 111-124 g



n = 15




Datapoints from gbif.org



# The Evolution of Comparative Phylogeography: Putting the Geography (and More) into Comparative Population Genomics

GBE

Scott V. Edwards <sup>1,2,\*</sup>, V. V. Robin<sup>3</sup>, Nuno Ferrand<sup>4</sup>, and Craig Moritz<sup>5</sup>

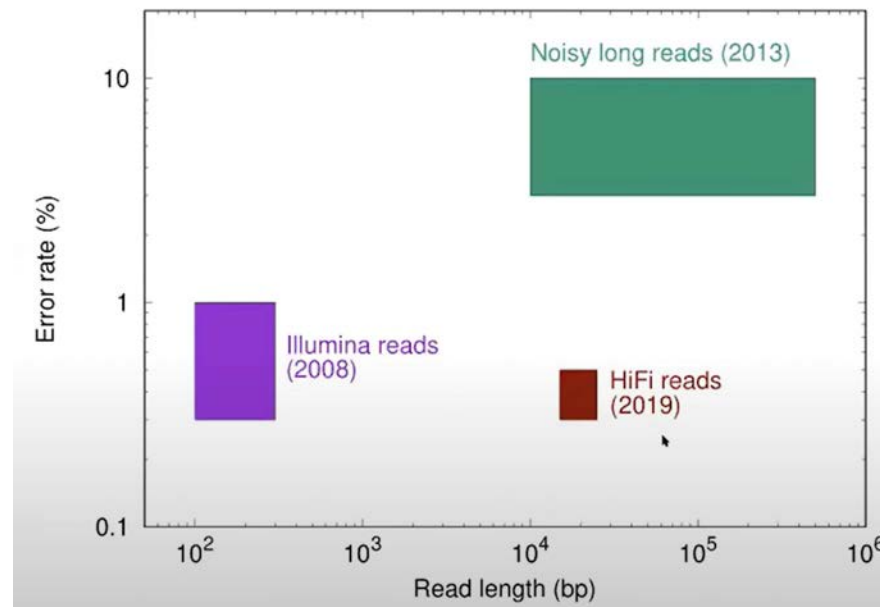
**Table 1**

Conceptual Relationships between the Fields of Comparative Population Genomics, Landscape Genomics, and Comparative Phylogeography

Concept/Parameter	Comparative Population Genomics	Landscape Genomics	Comparative Phylogeography
Comparative perspective	Growing	Nascent	Mature
Emphasis on space	No	Yes	Yes
Geographic scale	Random mating population	Region	Biome
Temporal scale	Arbitrary	Recent	Deep
Focus on:			
selection versus neutrality	Both	Both	Neutrality
recombination	Yes	Not yet considered	Not yet considered
geography versus environment	Nuisance parameters	Environment	Both
Future use of whole-genome sequencing	Yes	Likely	Unlikely
Growth out of museum collections community	No	No	Partial

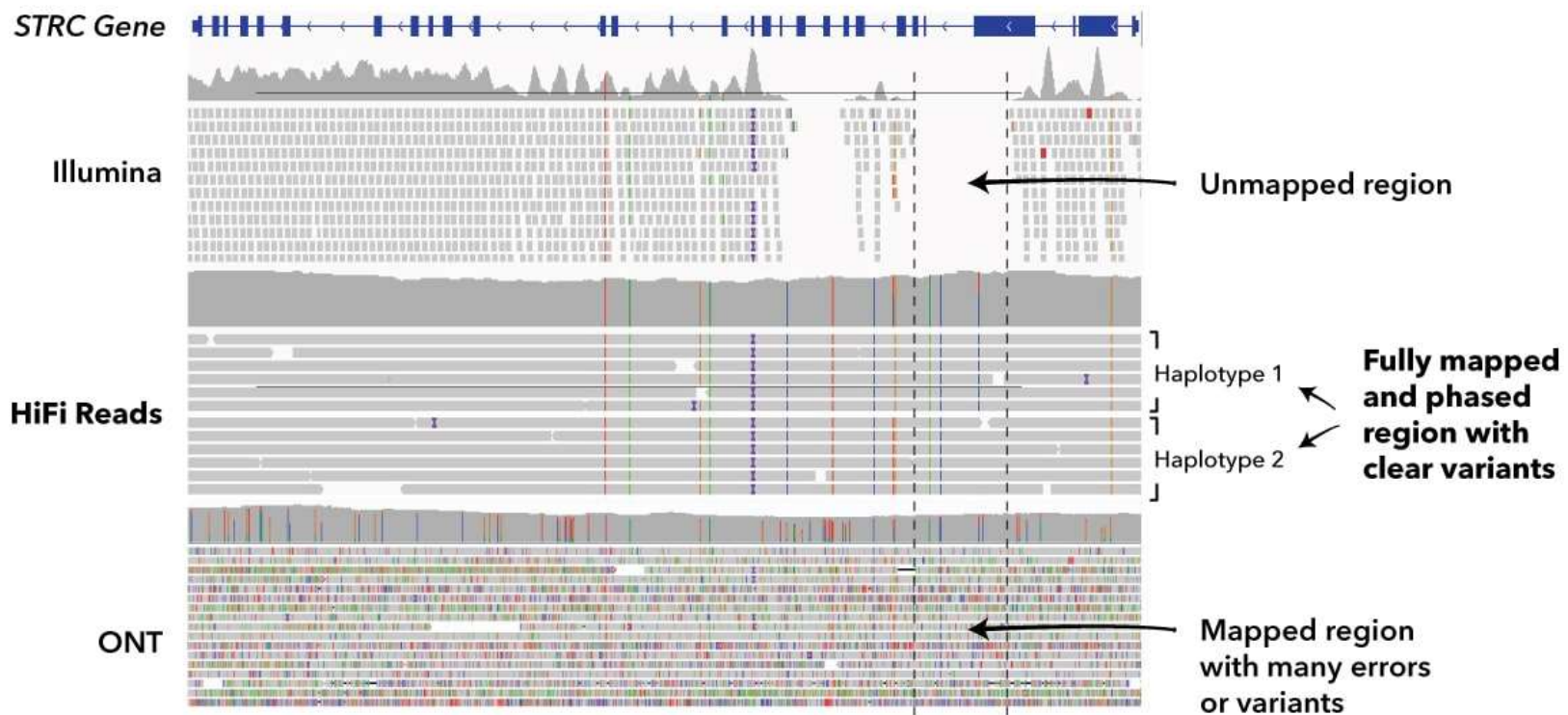
# PacBio HiFi reads are long and accurate

- ▶ HiFi reads: long & accurate
- ▶ A breakthrough every ~5 years
- ▶ Most existing assemblers cannot make full use of the accuracy



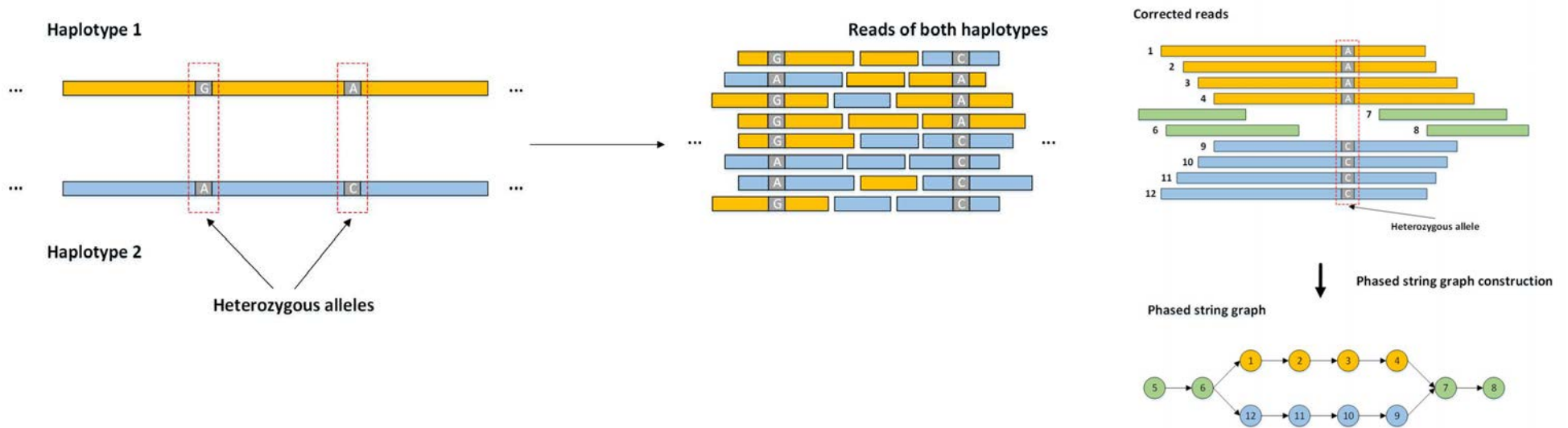
Courtesy Haoyu Cheng, Dana Farber Cancer Institute

# PacBio HiFi reads are long and accurate



HG002 GRCh38 chr15:43,599,422-43,619,001 (19 kb)

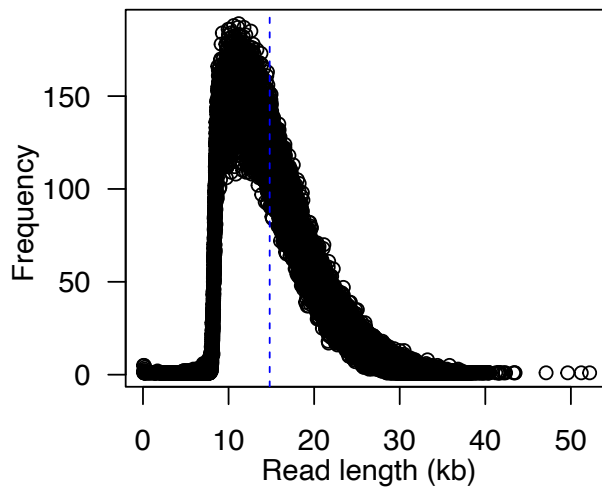
# Hifiasm – a HiFi accurate read assembler that resolves haplotypes



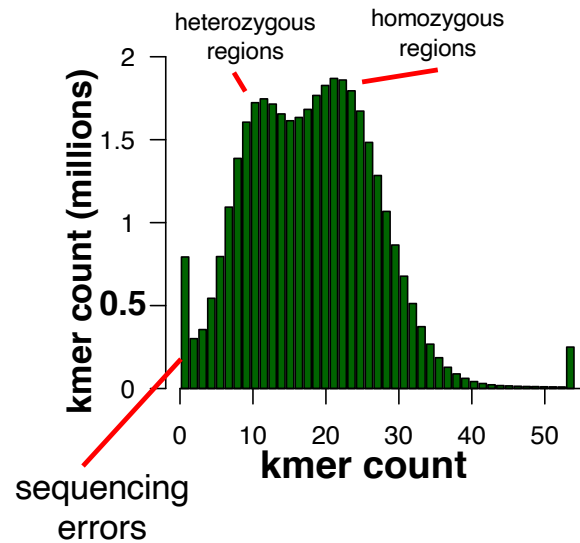
Courtesy Haoyu Cheng, Dana Farber Cancer Institute

# Scrub-jay PacBio HiFi data characteristics

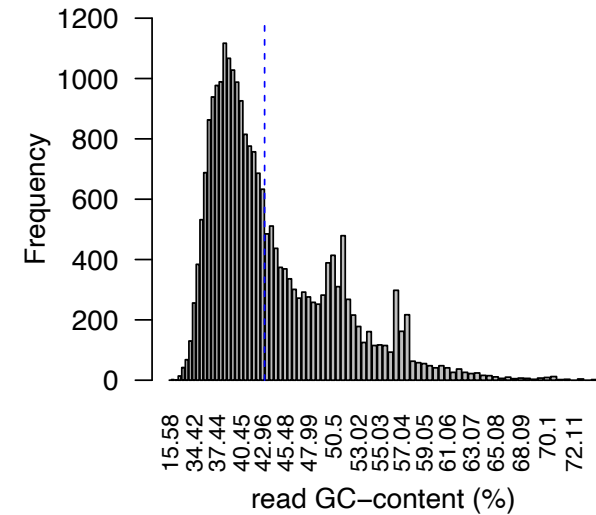
Average read length 14.8 kb



Average coverage ~40X



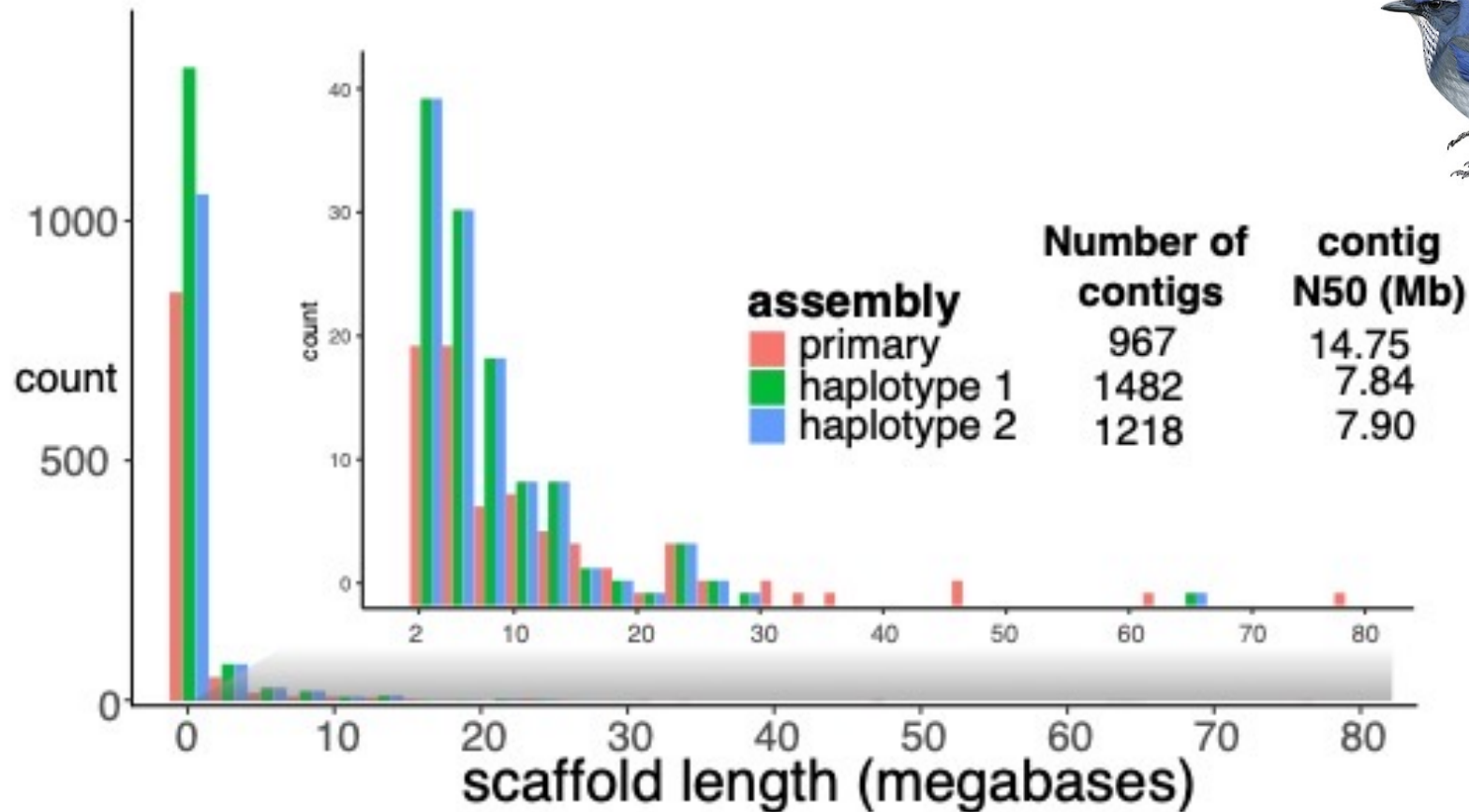
GC-content



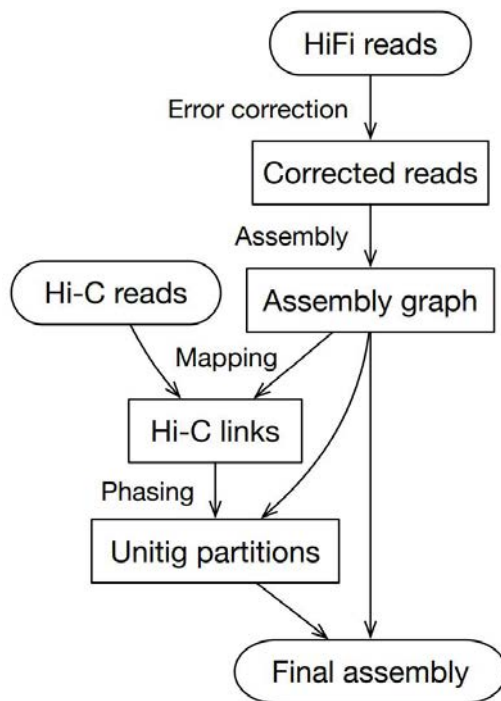
# Genome assembly with hifiasm yields ~1.3 Gb primary and haplotype assemblies



*A. woodhouseii*

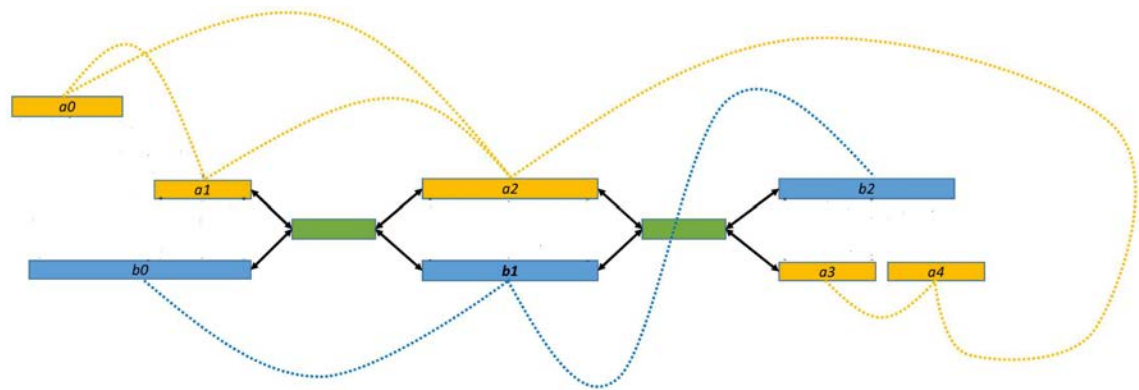


# Hifiasm – improved assemblies using HiC

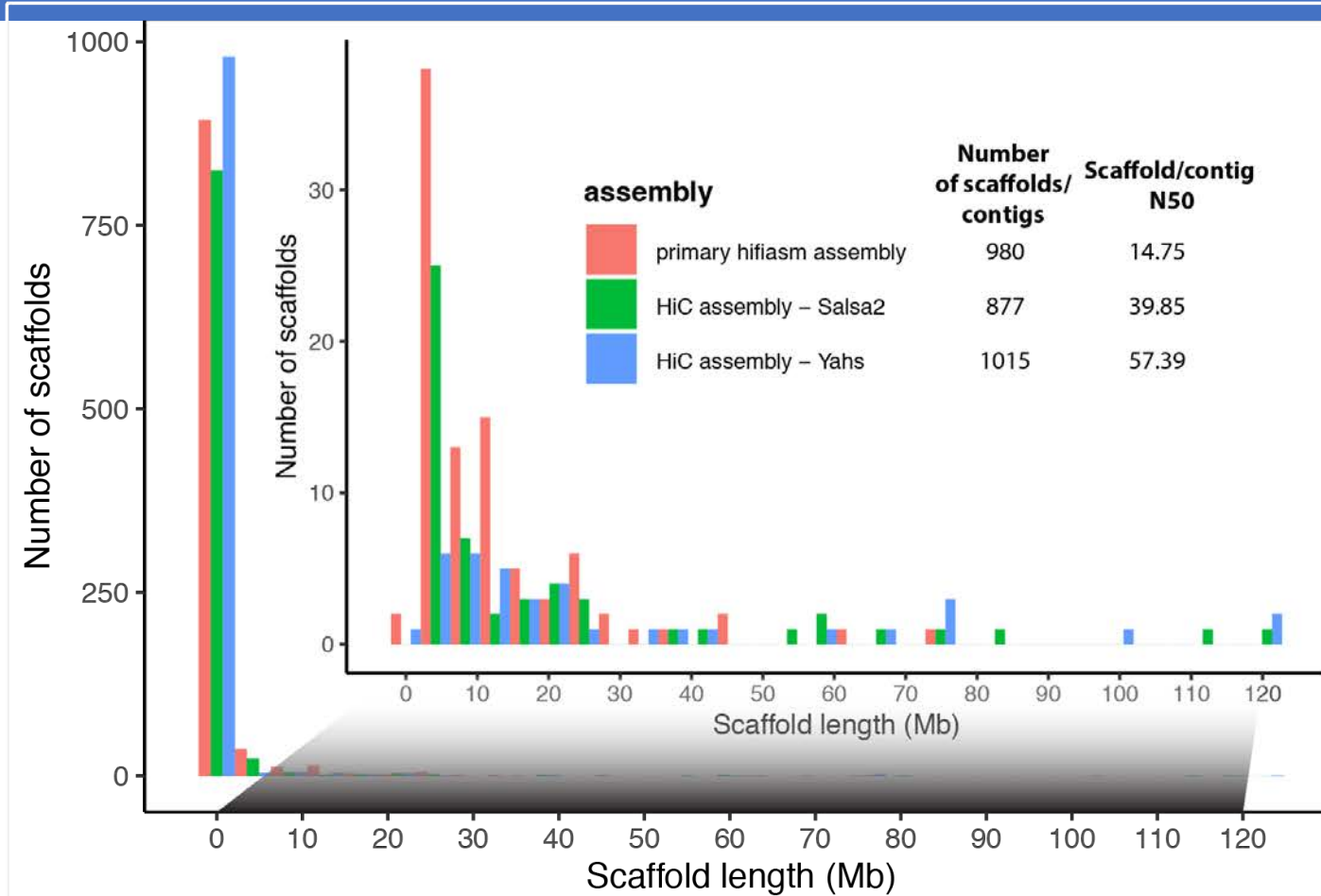


## ► Procedure:

- Identify heterozygous unitigs by coverage
- Build index by unique  $k$ -mers from heterozygous unitigs
- Align Hi-C reads using unique  $k$ -mers



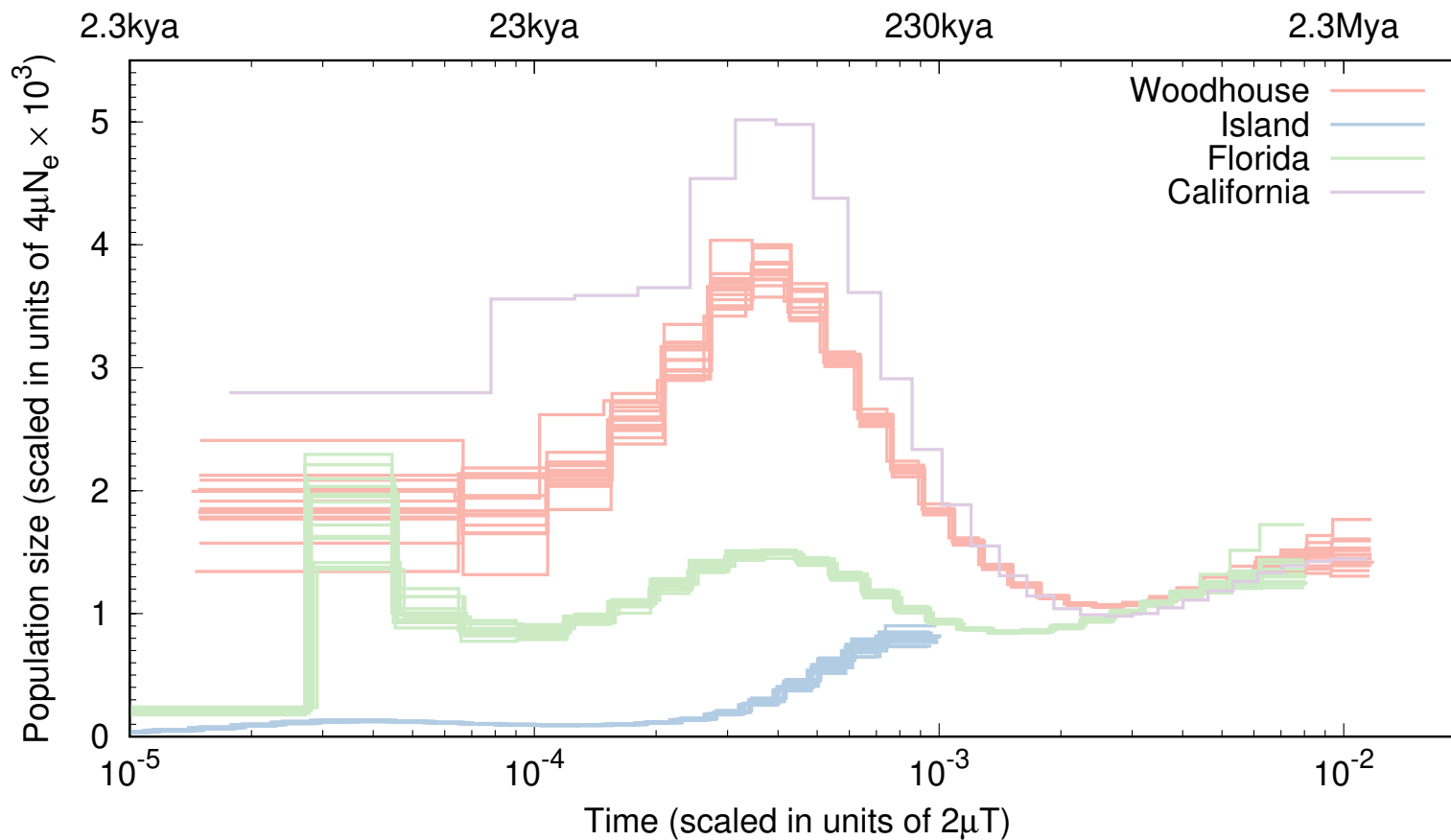
# HiC greatly improves contiguity of scrub jay assemblies



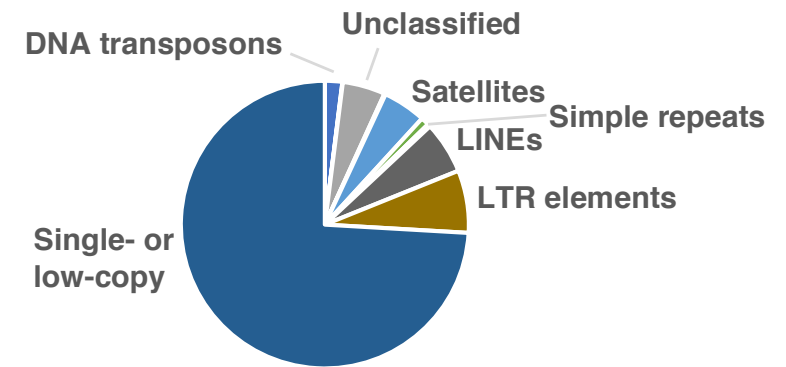
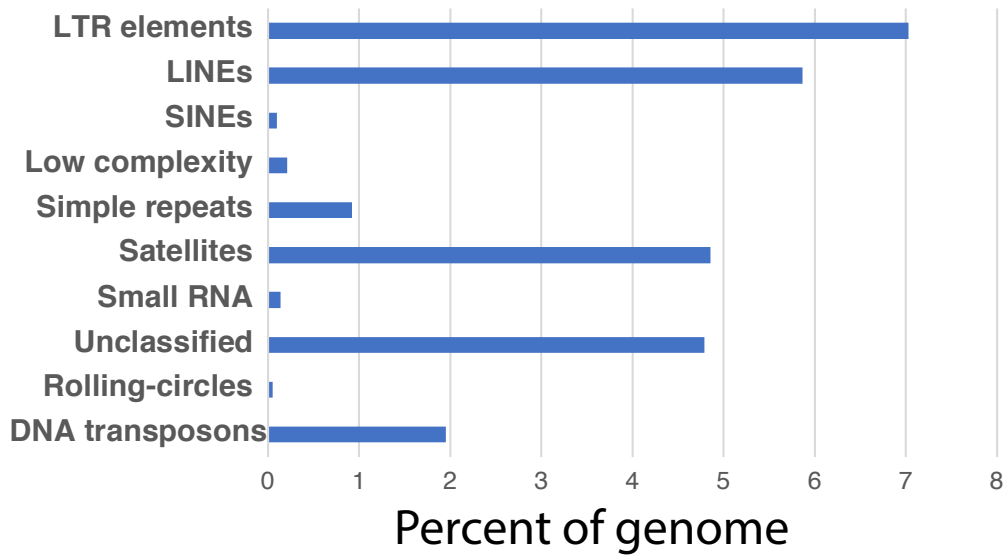


# PSMC analysis confirms variation in effective population size through time

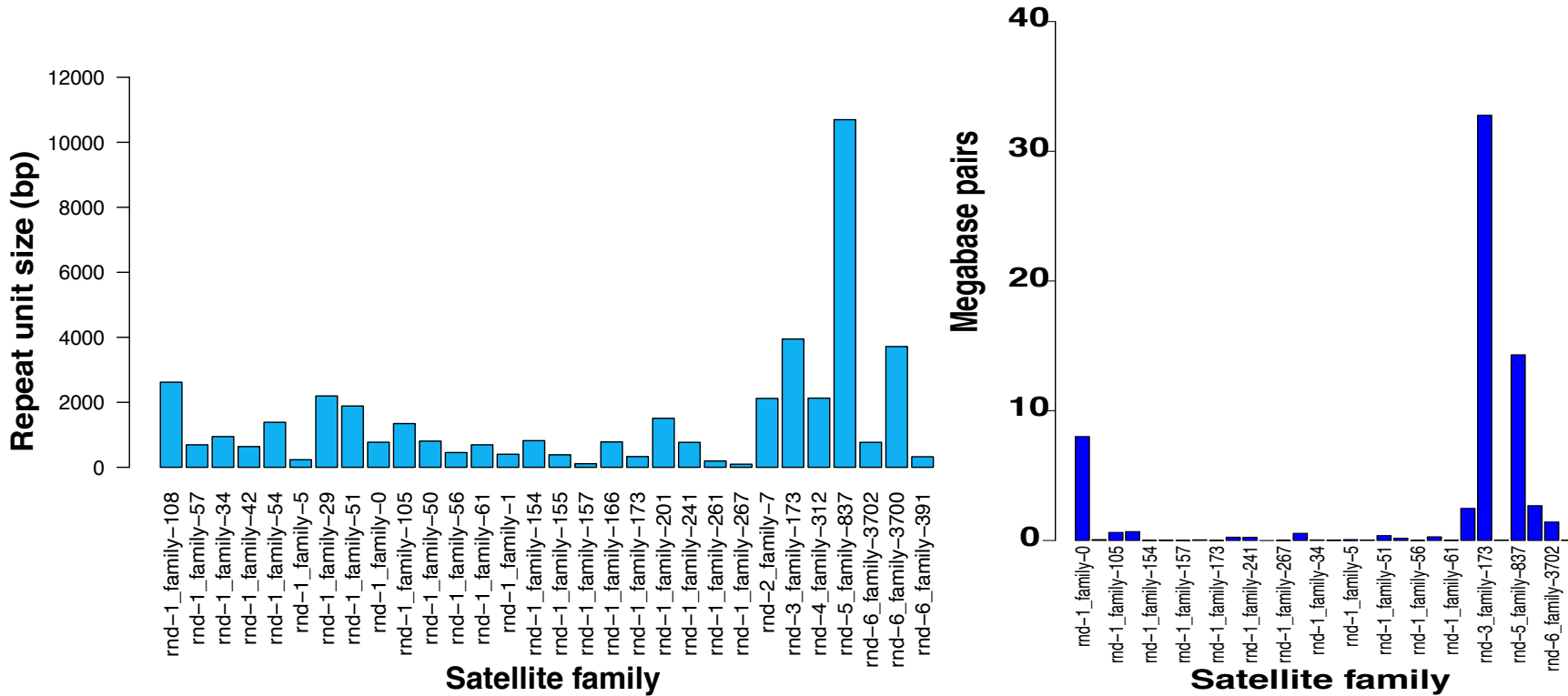
Pairwise Sequentially Markovian Coalescent (Li & Durbin 2011. Nature)



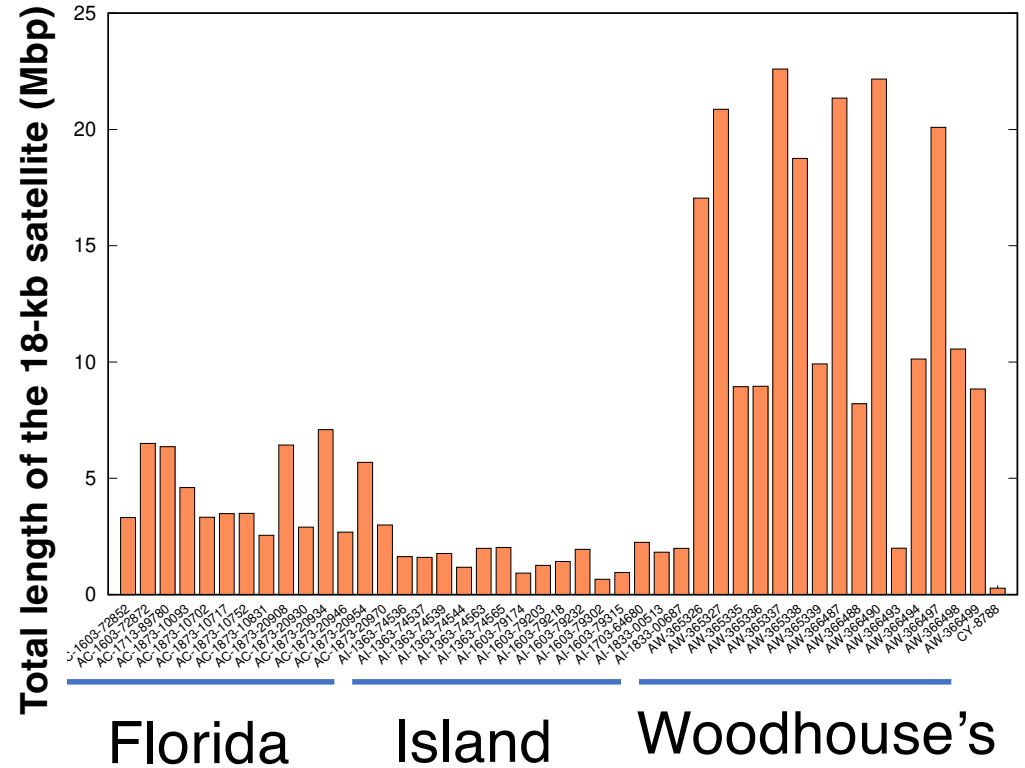
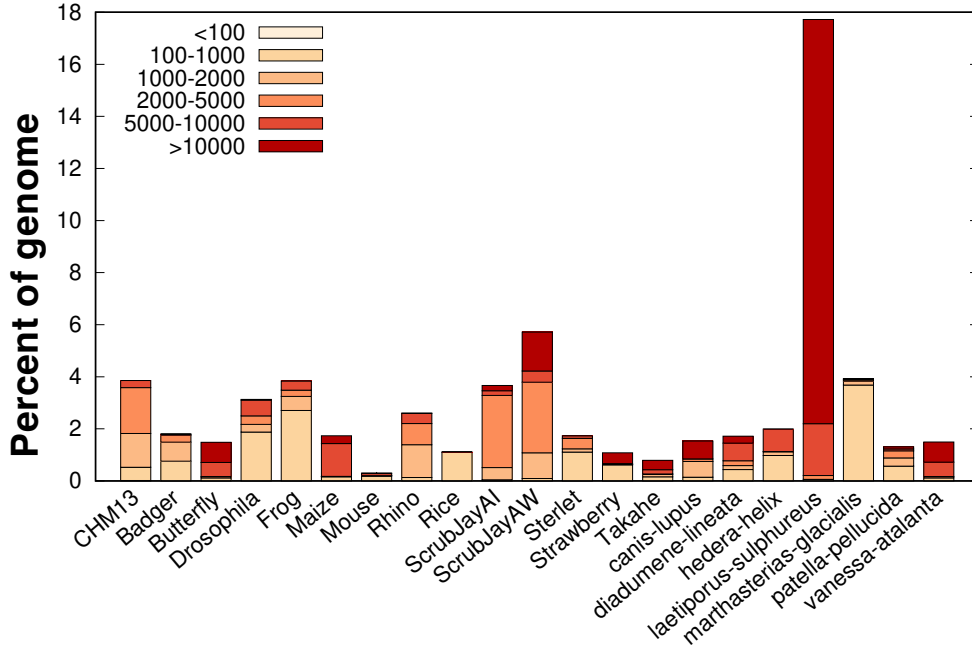
# RepeatMasker analysis suggests over 25% repeats and transposable elements



# Satellites are long and prevalent in scrub jay genomes



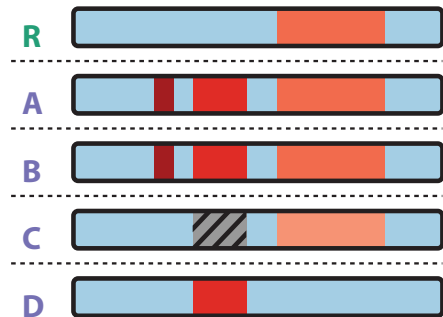
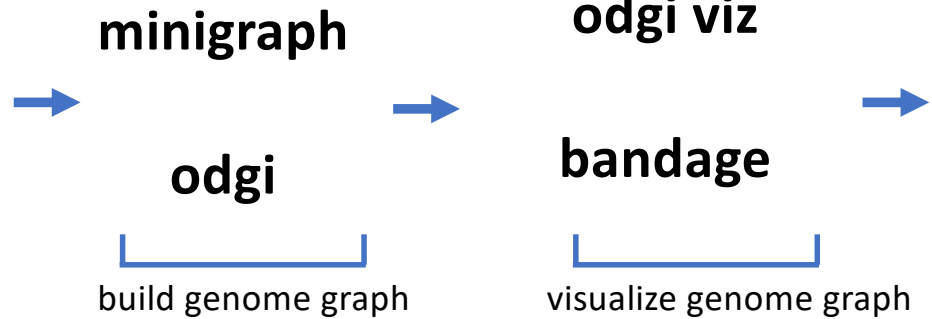
# Abundance of an 18-kb unit repeat satellite varies strongly among species



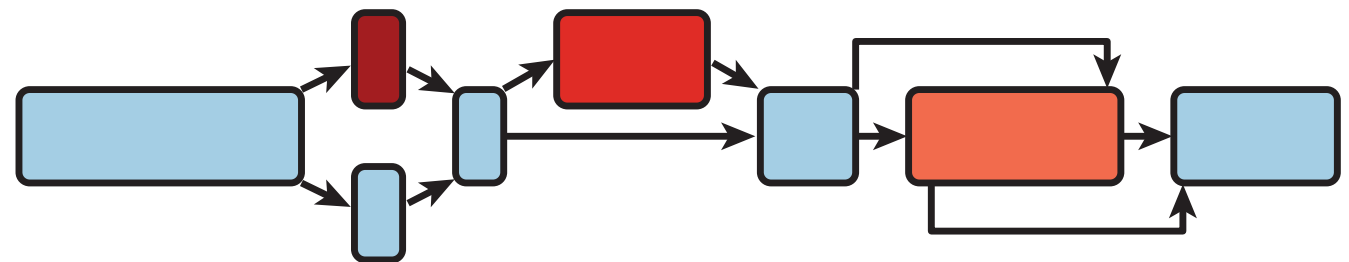
# Pangenome graphs capture structural variation within species

```
>h1tg000104l  
GGCGGGGCCCGGAGGGGCCGGGGCCGCTGAGGGGCCGCGGGTGCAGAGCC  
>h1tg000528l  
ATGGATACTTTCCAGTCAGAGCTTTATAATAATTTCCATAATTTAAATATTT  
>h1tg000795l  
ACTTTGGGGACACCTTTGGGGACACCTCGGGGACACTTTGGGCCACAAATCC
```

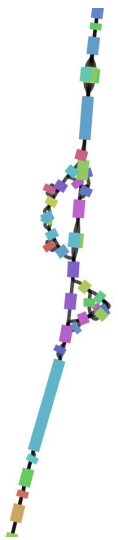
unaligned fasta files



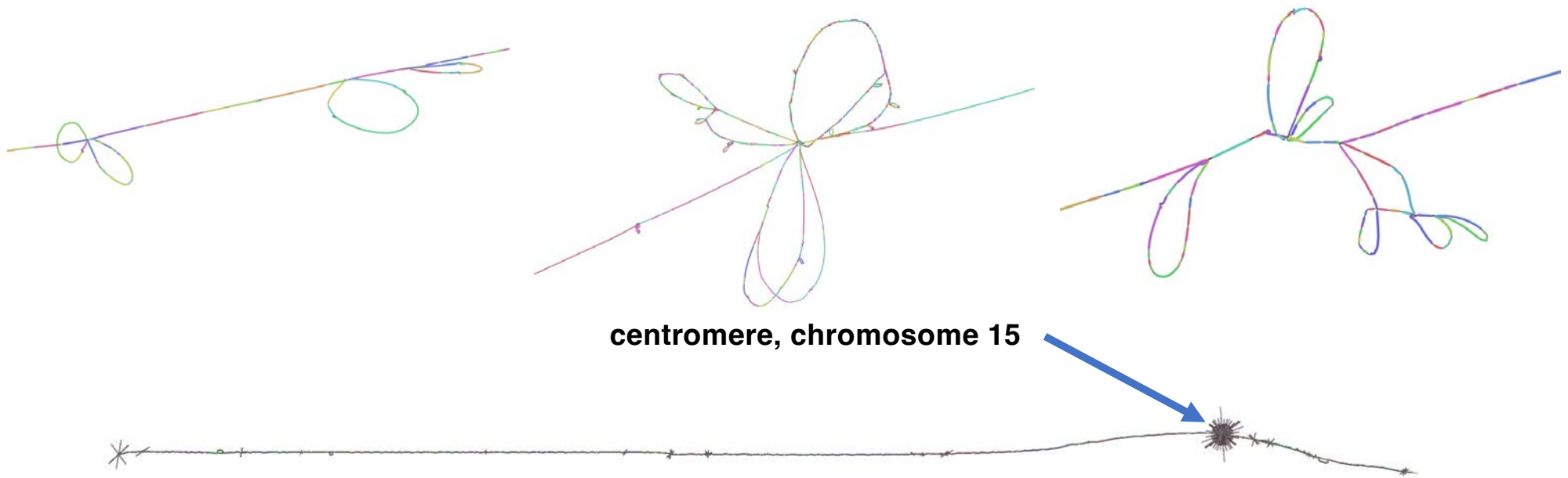
Multiple sequence alignment



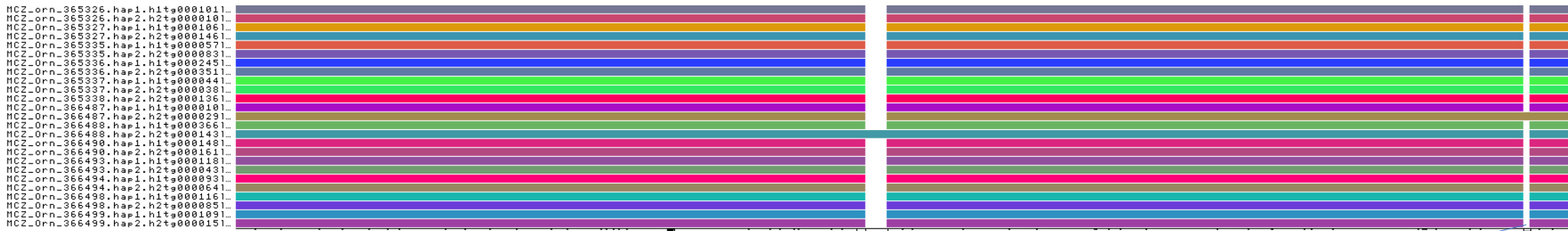
Bidirected genome graph



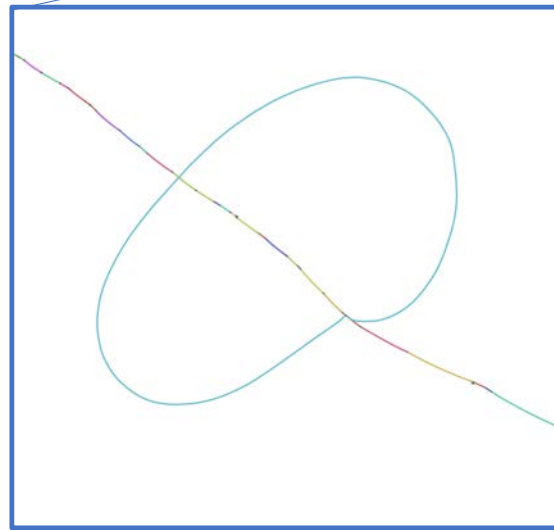
# Pangenome graphs of haplotype variation in Scrub Jays



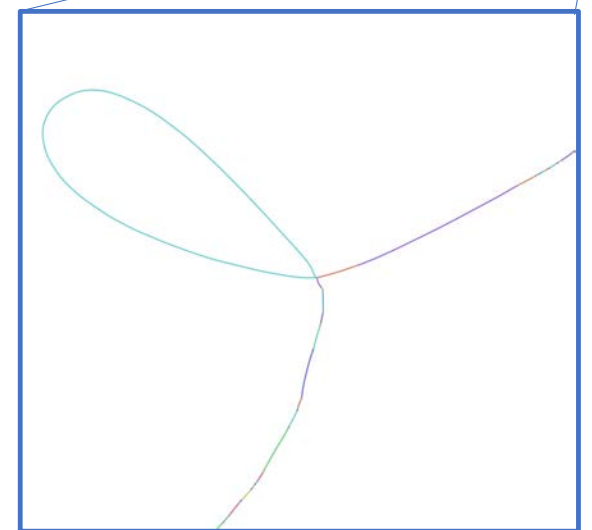
# Genomic stability of 400-kb *hox1a* region in Western Scrub Jays



Pangenome graphs  
generated with odgi  
and visualized with Bandage



7.5 kb polymorphic indel



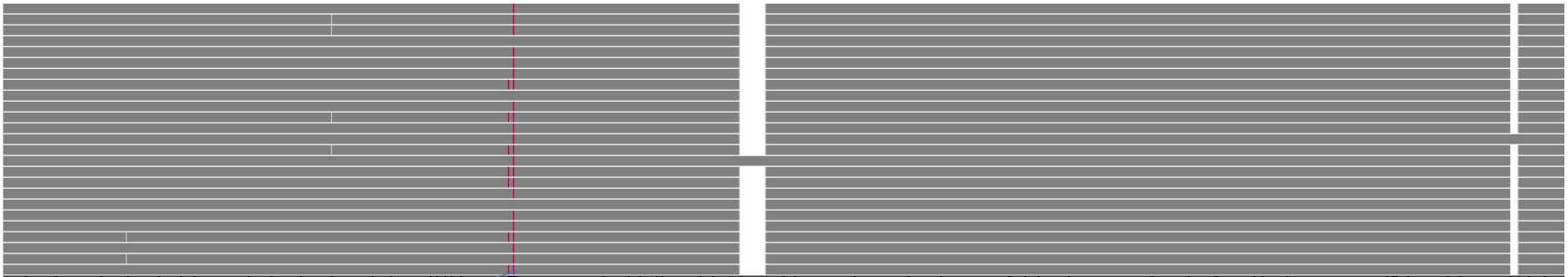
2.5 kb polymorphic indel

Guarracino et al. 2021.  
*Bioinformatics*, in press.  
Wick et al. 2015.  
*Bioinformatics* 31:3350.

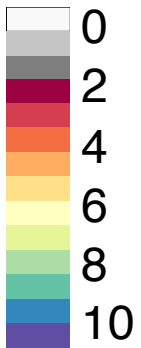
# Smaller regions of complexity in *hox1a* region

```

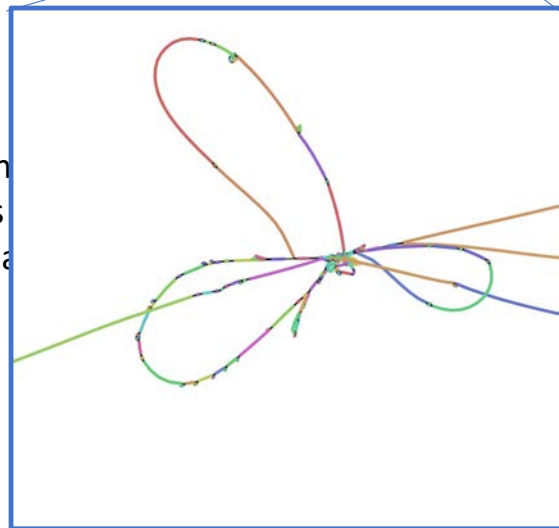
MCZ_orn_365326.hap1.hitg0001011...
MCZ_orn_365326.hap2.h2tg0001011...
MCZ_orn_365327.hap1.hitg0001061...
MCZ_orn_365327.hap2.h2tg0001461...
MCZ_orn_365335.hap1.hitg0000571...
MCZ_orn_365335.hap2.h2tg0000831...
MCZ_orn_365336.hap1.hitg0002451...
MCZ_orn_365336.hap2.h2tg0003511...
MCZ_orn_365337.hap1.hitg0000441...
MCZ_orn_365337.hap2.h2tg0000381...
MCZ_orn_365338.hap1.hitg0001361...
MCZ_orn_365338.hap2.h2tg0001911...
MCZ_orn_365487.hap1.hitg0000191...
MCZ_orn_365487.hap2.h2tg0000291...
MCZ_orn_365488.hap1.hitg0003661...
MCZ_orn_365488.hap2.h2tg0001431...
MCZ_orn_365490.hap1.hitg0001481...
MCZ_orn_365490.hap2.h2tg0001611...
MCZ_orn_365493.hap1.hitg0001181...
MCZ_orn_365493.hap2.h2tg0000431...
MCZ_orn_365494.hap1.hitg0000391...
MCZ_orn_365494.hap2.h2tg0000641...
MCZ_orn_365498.hap1.hitg0001161...
MCZ_orn_365498.hap2.h2tg0000851...
MCZ_orn_365499.hap1.hitg0001091...
MCZ_orn_365499.hap2.h2tg0000151...
    
```



depth of *hox1a* region graph (x)



White regions  
Gray regions  
Red regions and

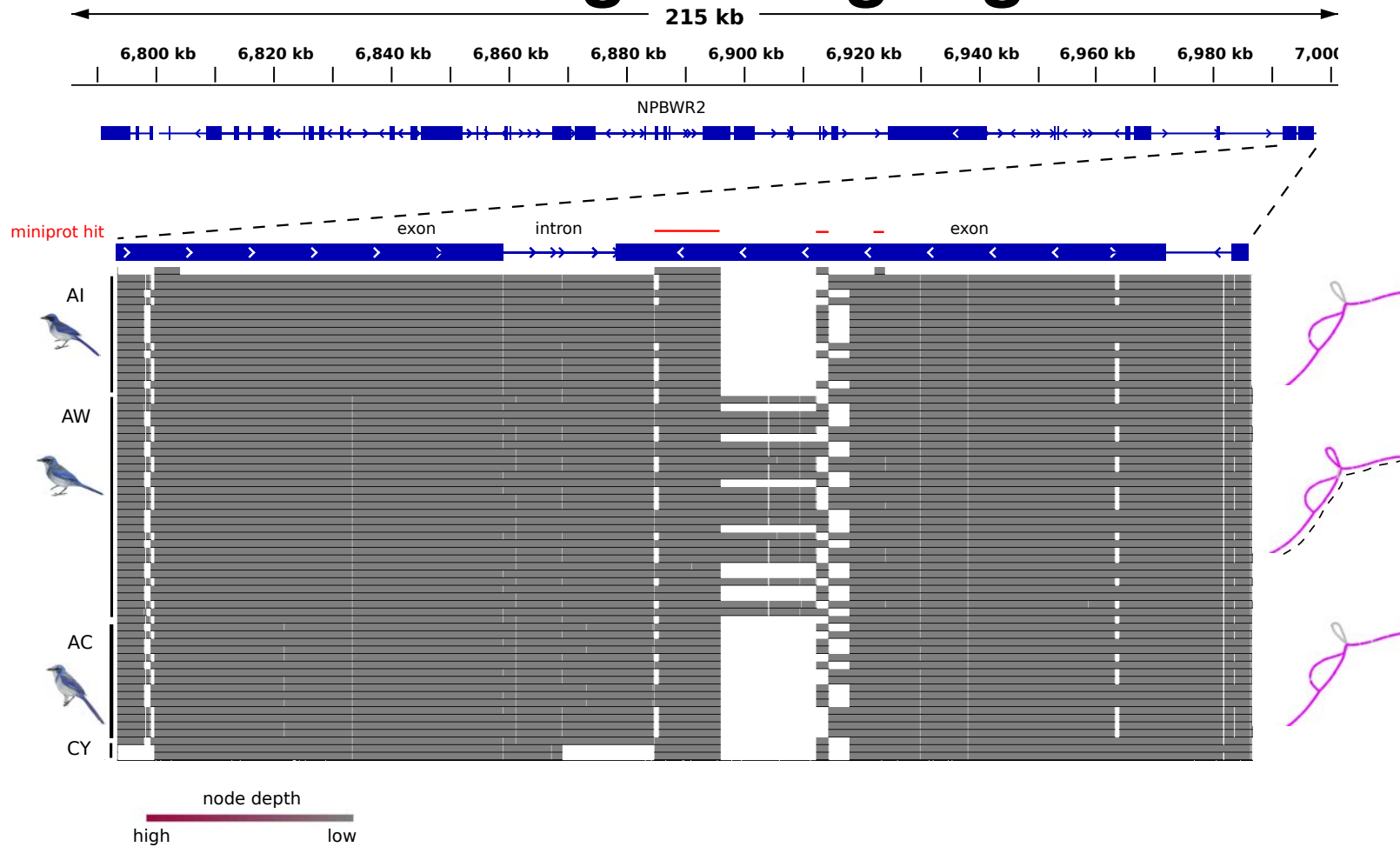


regions





# Structural variation in another conservatively evolving coding region



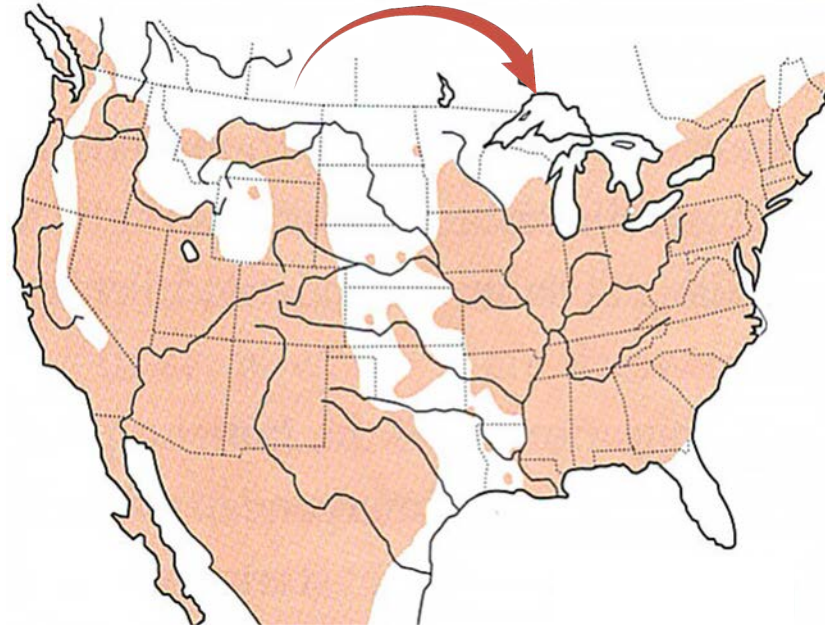
# Demographic history of House Finches

**WEST**

Historic range

**EAST**

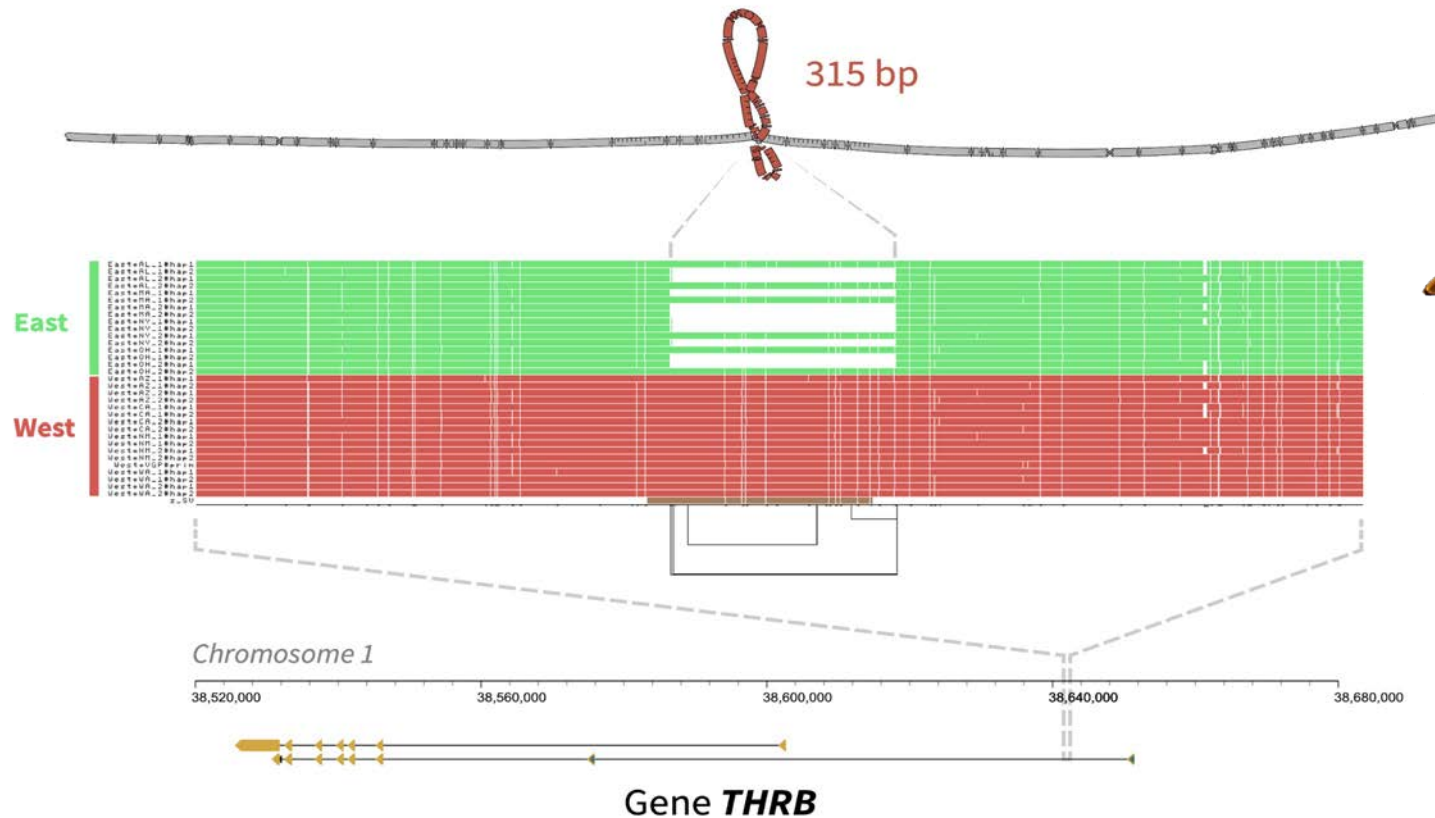
Introduced range



**Species distribution in ~1990**

*Maps: Birds of North America*

# Structural variants in the thyroid receptor- $\beta$ gene



# Major histocompatibility complex

```
graph TD; A([Major histocompatibility complex]) --> B([Molecular evolution]); A --> C([Sexual selection/parasites]); A --> D([Kin recognition]); A --> E([Conservation genetics]);
```

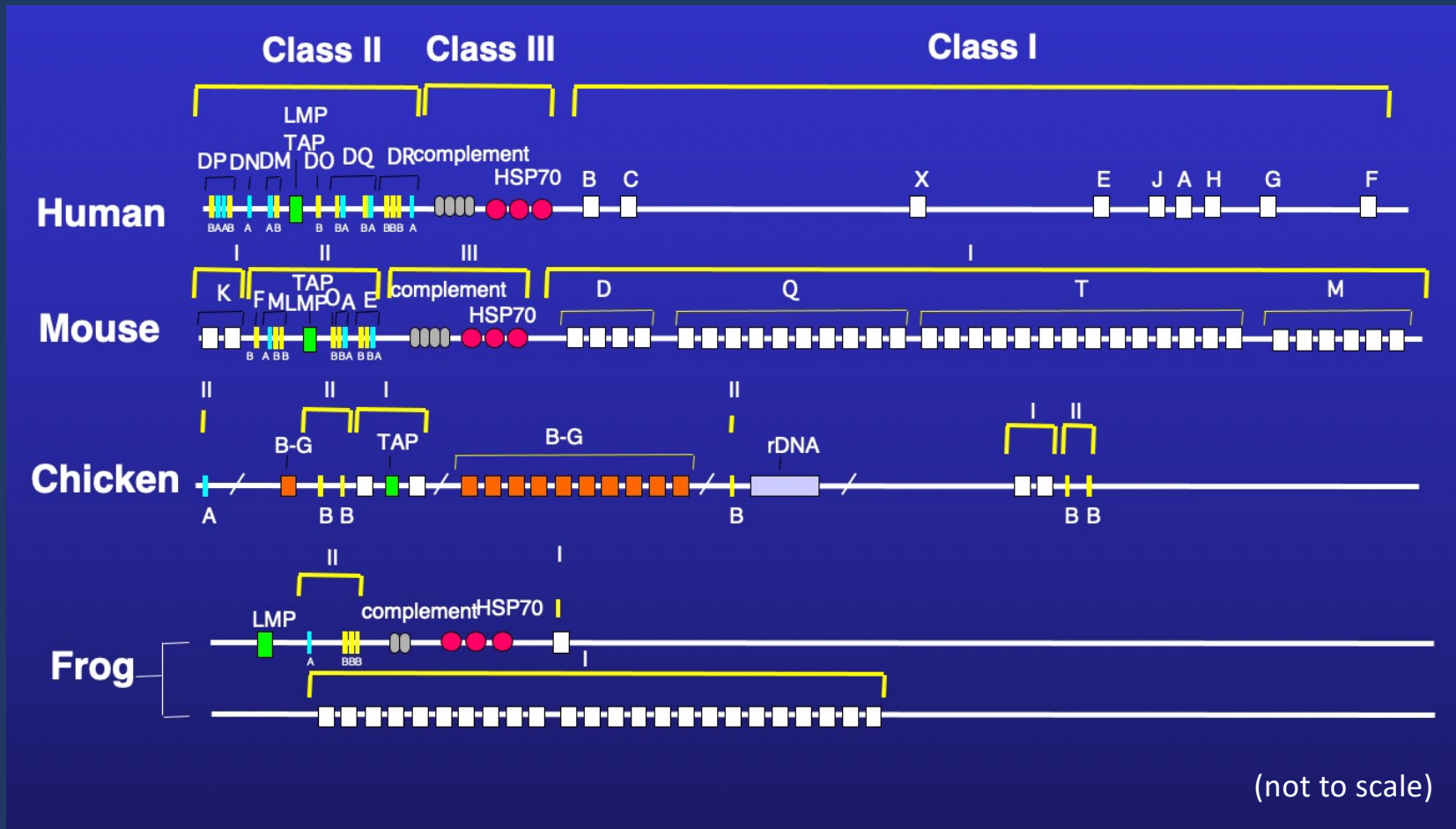
Molecular evolution

Sexual selection/  
parasites

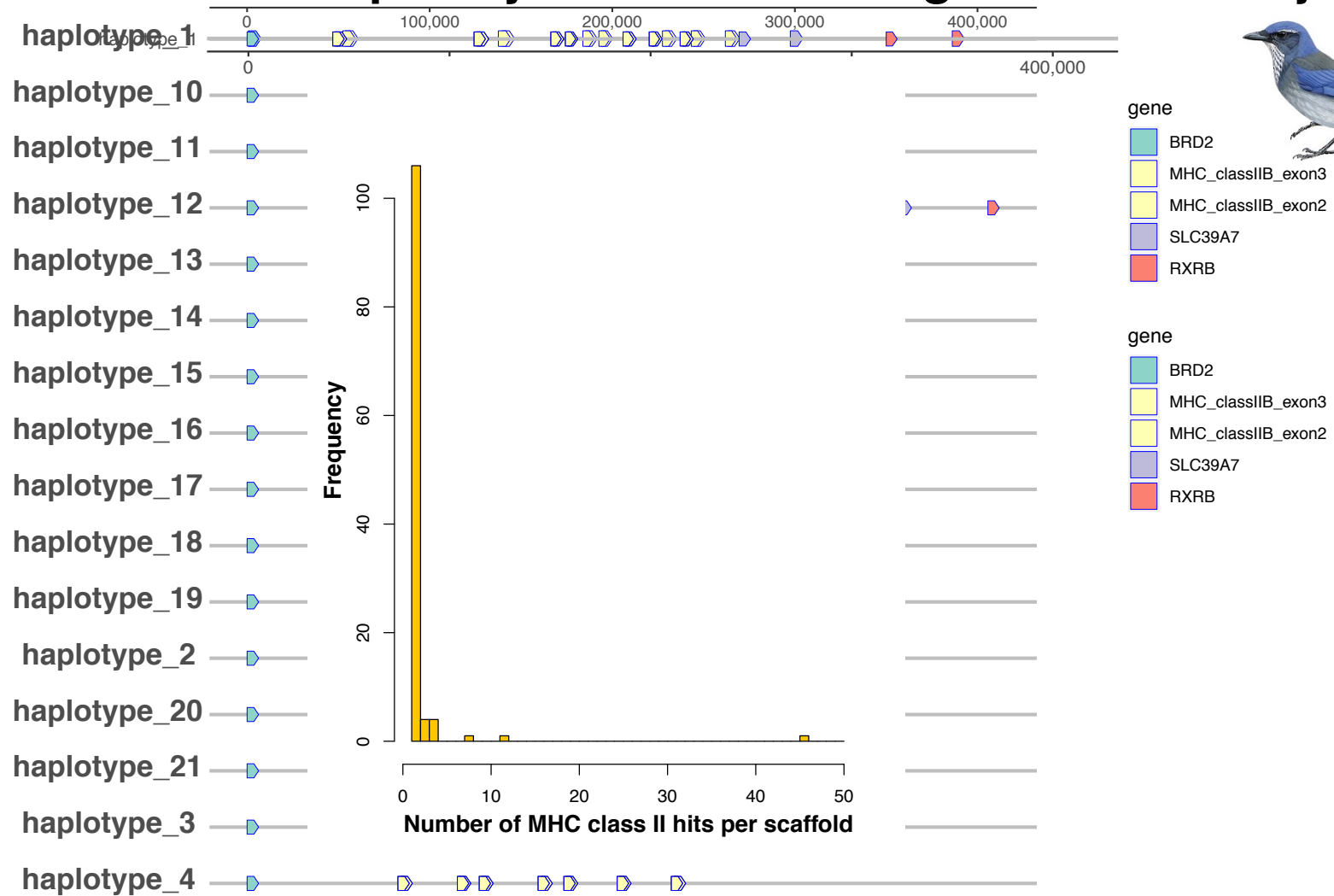
Kin recognition

Conservation genetics

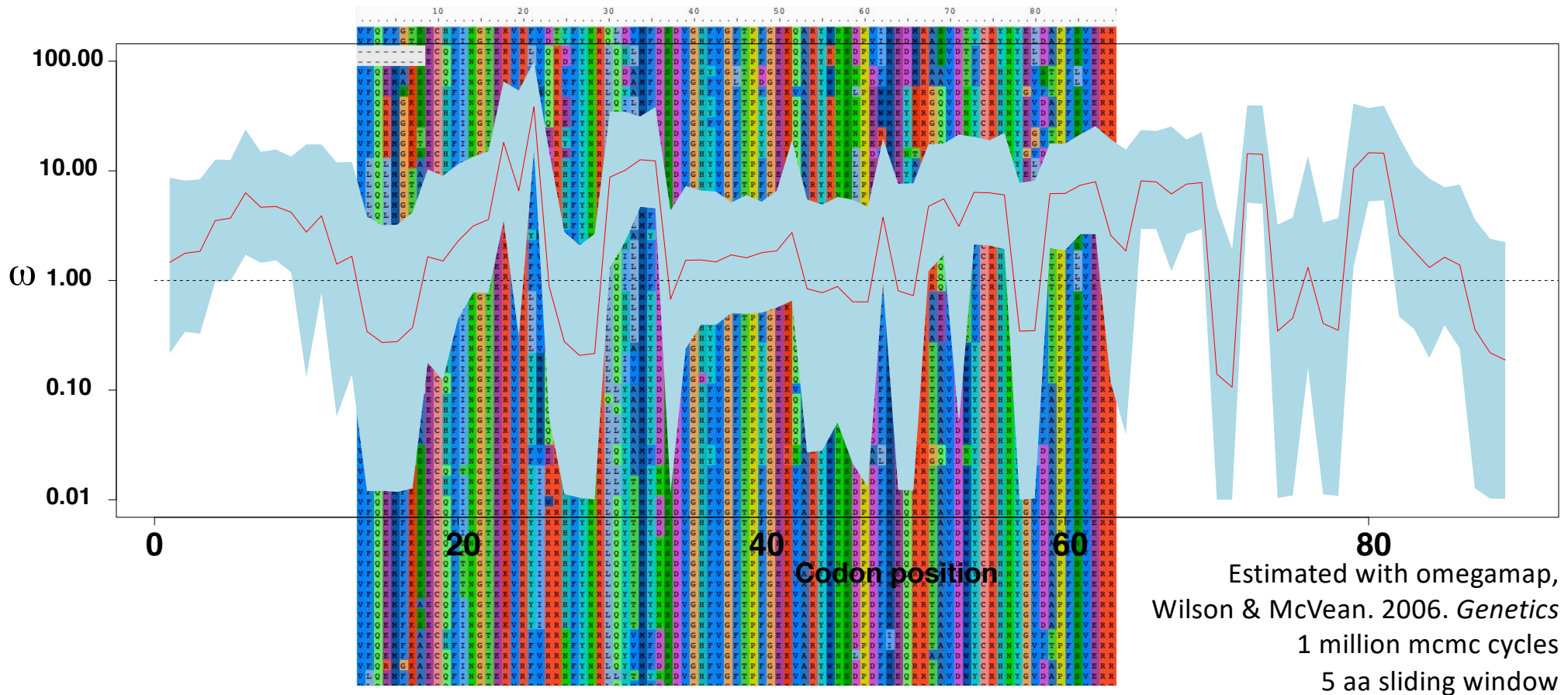
# The chicken MHC is small (~99 kb) and compact



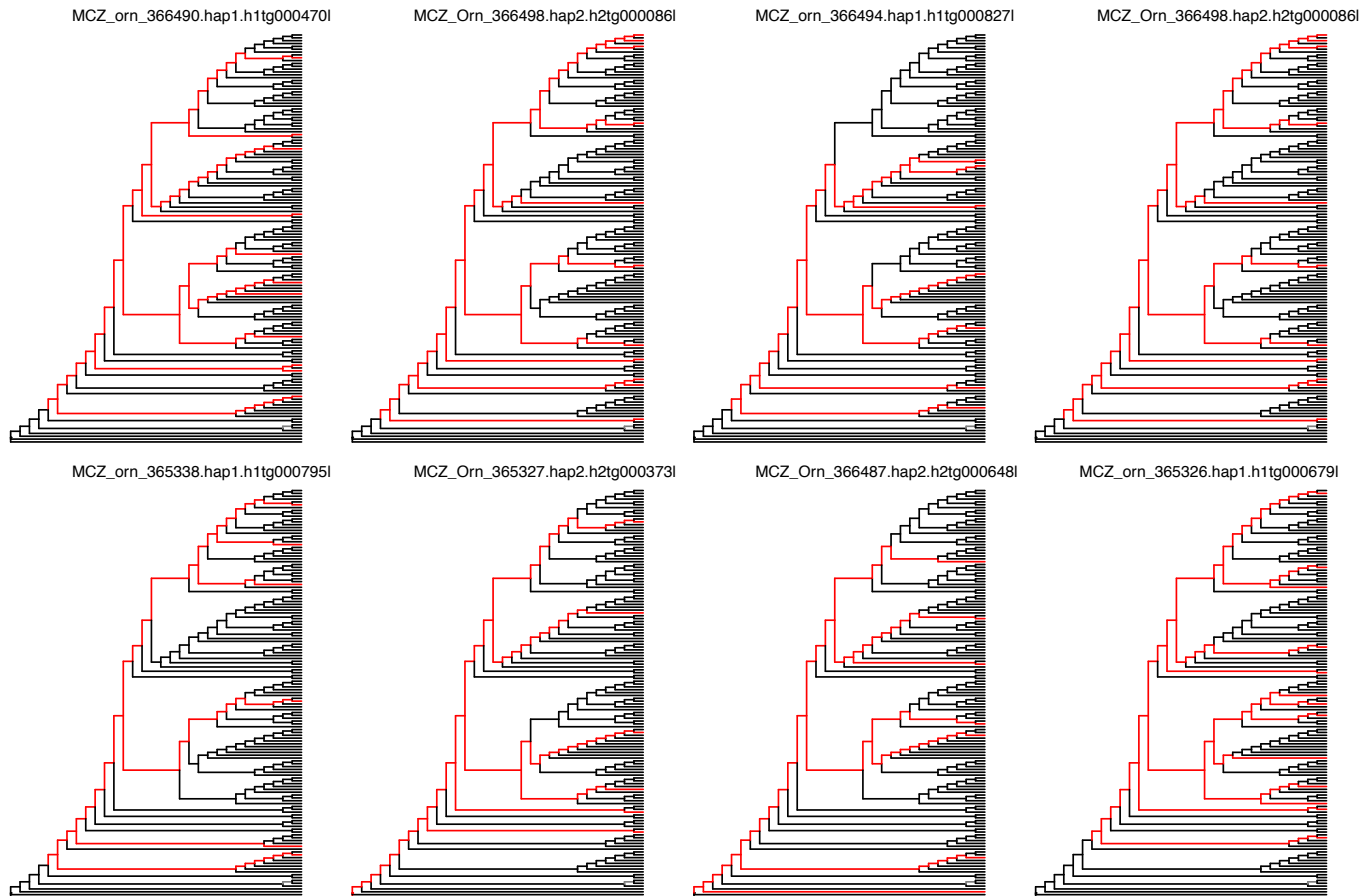
# Unprecedented complexity of MHC class II genes in scrub jays



# Mhc class II peptide-binding region shows solid evidence of balancing selection



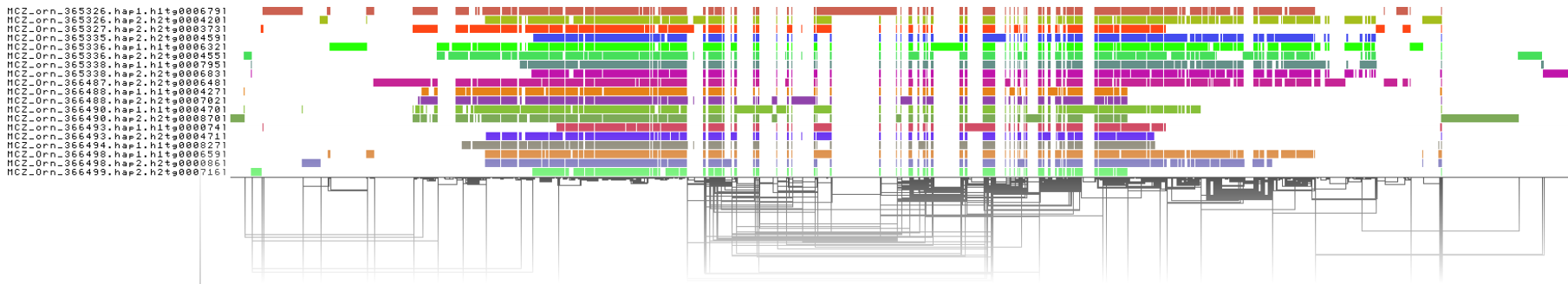
# Mhc class II peptide binding regions are phylogenetically diverse on individual haplotypes



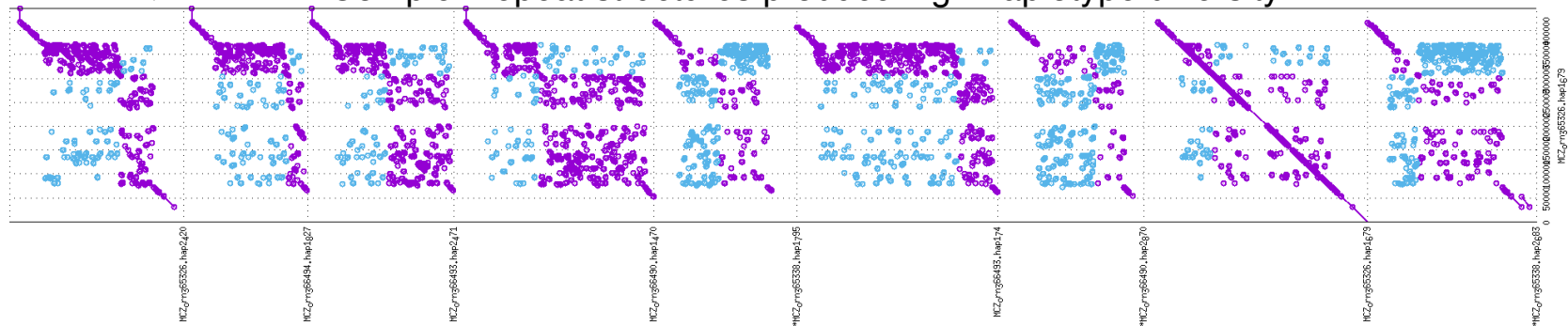
Phylogenetic paths of  
Mhc exon2 alleles  
on individual haplotypes



# Visualization of MHC class II region in 22 haplotypes of Woodhouse's scrub-jays with odgi



Complex repeat structures produce high haplotype diversity

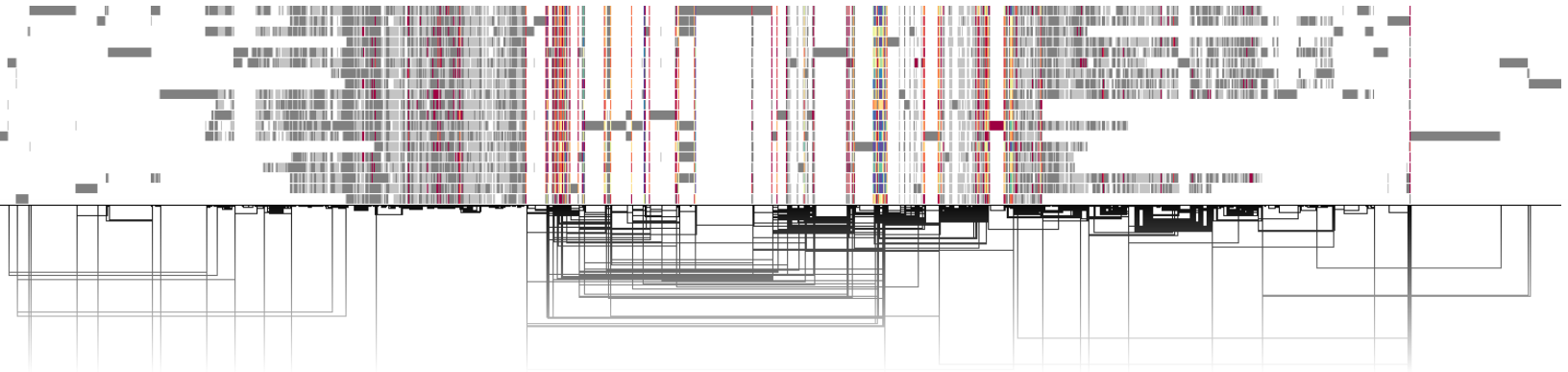


made with odgi and pangenome graph builder pipeline  
Guarracino et al. 2021. *Bioinformatics*, in press.

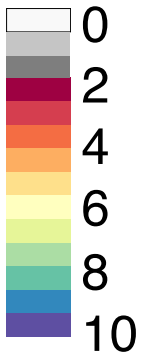


# Pangenome graph depth shows single-copy regions surrounded by complex VNTRs

MCZ\_orn\_365326.hap1.h1t@0005791  
MCZ\_orn\_365326.hap2.h2t@0004201  
MCZ\_orn\_365327.hap2.h2t@0003731  
MCZ\_orn\_365329.hap2.h2t@0004591  
MCZ\_orn\_365336.hap1.h1t@0006321  
MCZ\_orn\_365336.hap2.h2t@0004551  
MCZ\_orn\_365338.hap1.h1t@0007951  
MCZ\_orn\_365338.hap2.h2t@0006831  
MCZ\_orn\_365487.hap2.h2t@0006481  
MCZ\_orn\_365488.hap1.h1t@0004271  
MCZ\_orn\_365488.hap2.h2t@0007021  
MCZ\_orn\_365490.hap1.h1t@0004701  
MCZ\_orn\_365490.hap2.h2t@0005701  
MCZ\_orn\_365493.hap1.h1t@000741  
MCZ\_orn\_365493.hap2.h2t@0004711  
MCZ\_orn\_365494.hap1.h1t@0008271  
MCZ\_orn\_365498.hap1.h1t@000591  
MCZ\_orn\_365498.hap2.h2t@0008061  
MCZ\_orn\_365499.hap2.h2t@0007161



depth of  
MHC  
graph (x)



odgi visualization of pangenome graph depth based on  
26 MHC-containing scrub jay haplotypes, up to ~480 kb

White regions are indels between haplotypes

Gray regions are SNP variation

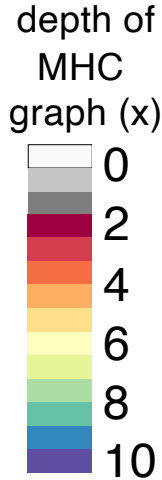
Red regions are low-complexity and repetitive regions



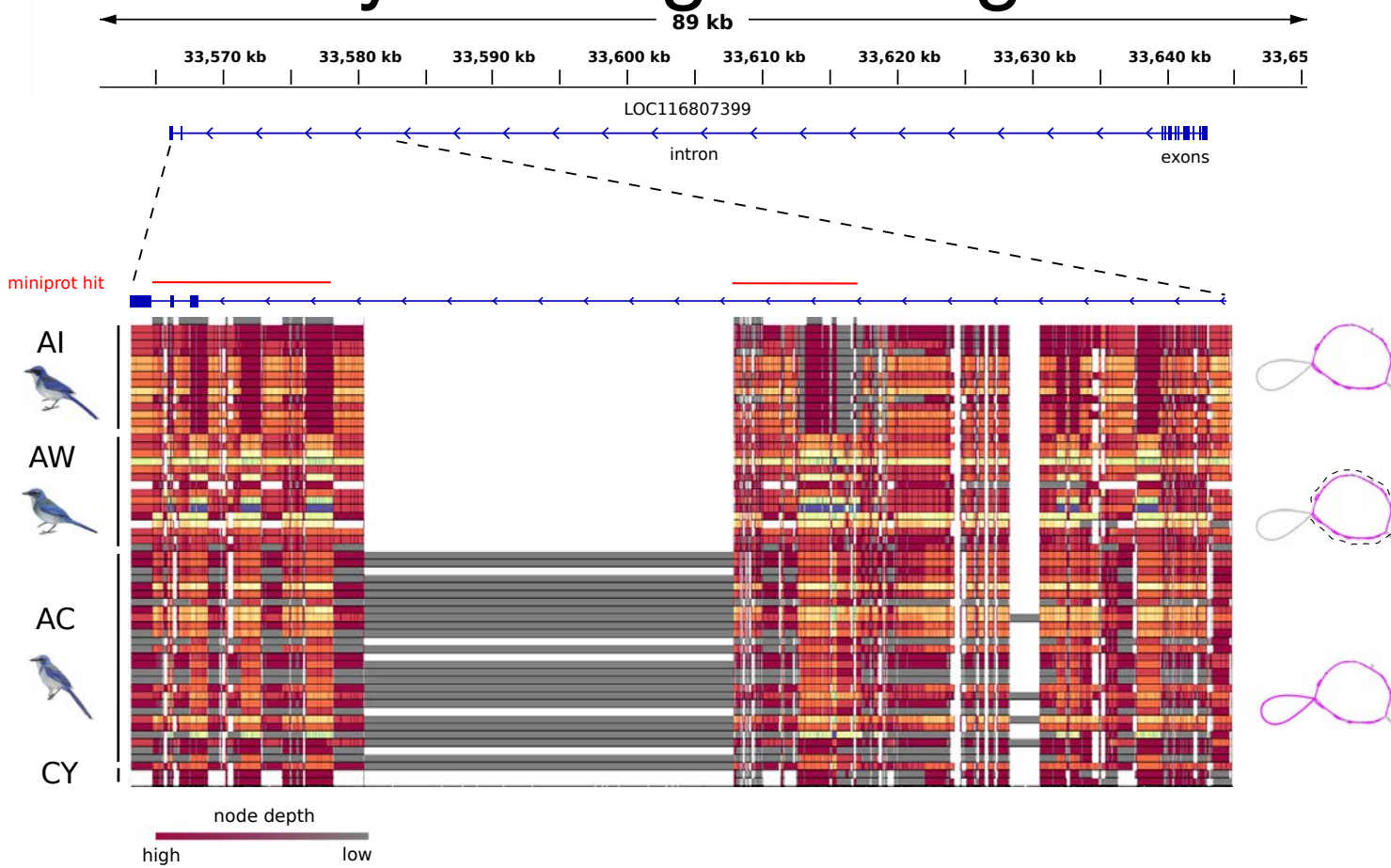
made with odgi and pangenome graph builder pipeline  
Guarracino et al. 2021. *Bioinformatics*, in press.

# Graph depth shows single-copy regions surrounded by complex VNTRs

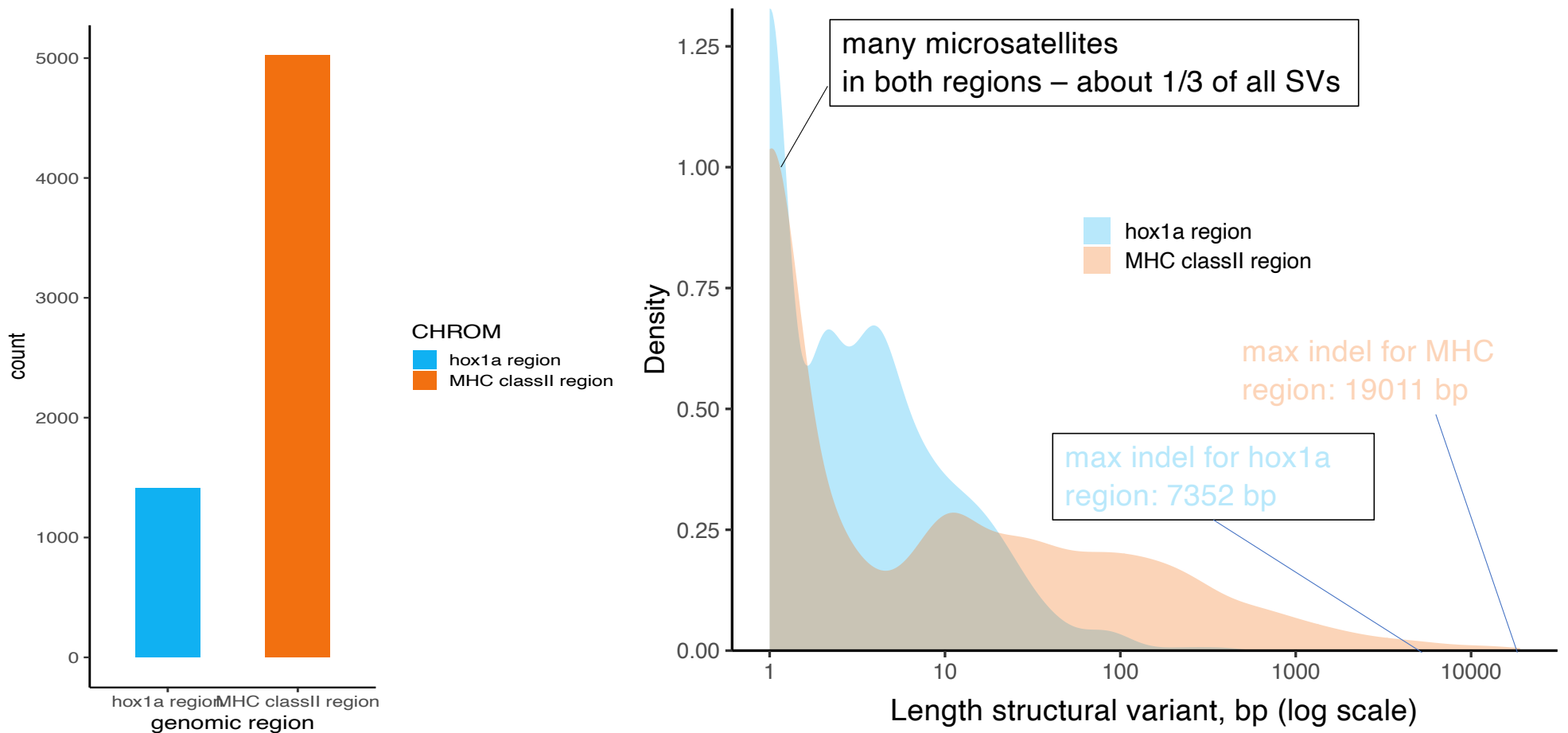
MCZ\_orn\_366494.hap1\_827  
MCZ\_orn\_366493.hap2\_471  
MCZ\_orn\_366490.hap2\_370  
MCZ\_orn\_366490.hap1\_470  
MCZ\_orn\_365338.hap2\_683  
MCZ\_orn\_365338.hap1\_795  
MCZ\_orn\_365326.hap2\_420  
MCZ\_orn\_365326.hap1\_679  
MCZ\_orn\_366493.hap1\_74



# Structural variation in another dynamic genic region

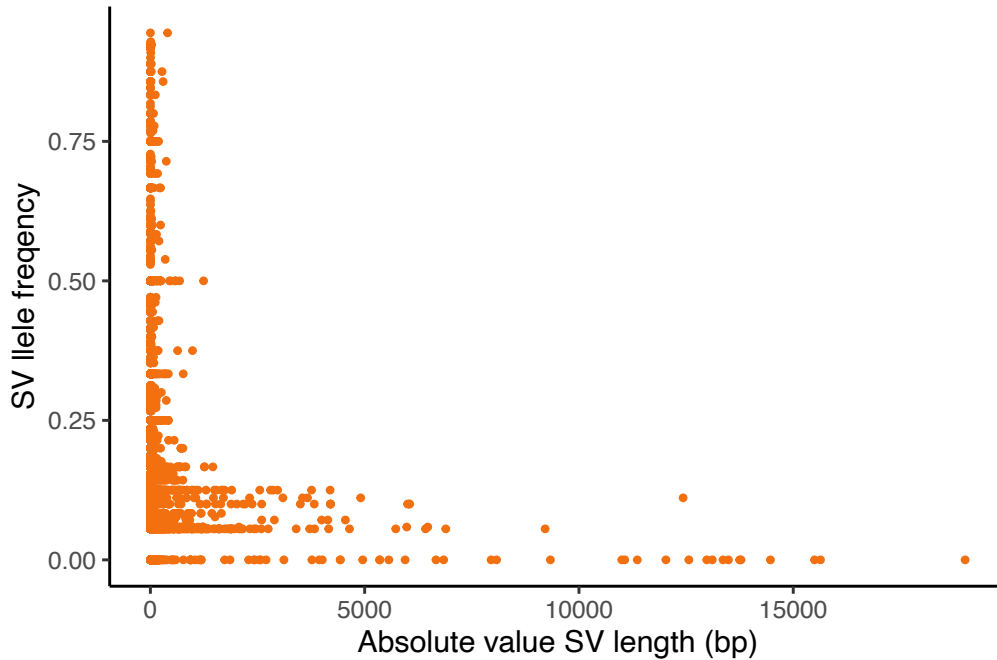


# MHC region has more numerous and longer structural variants than hox1a region

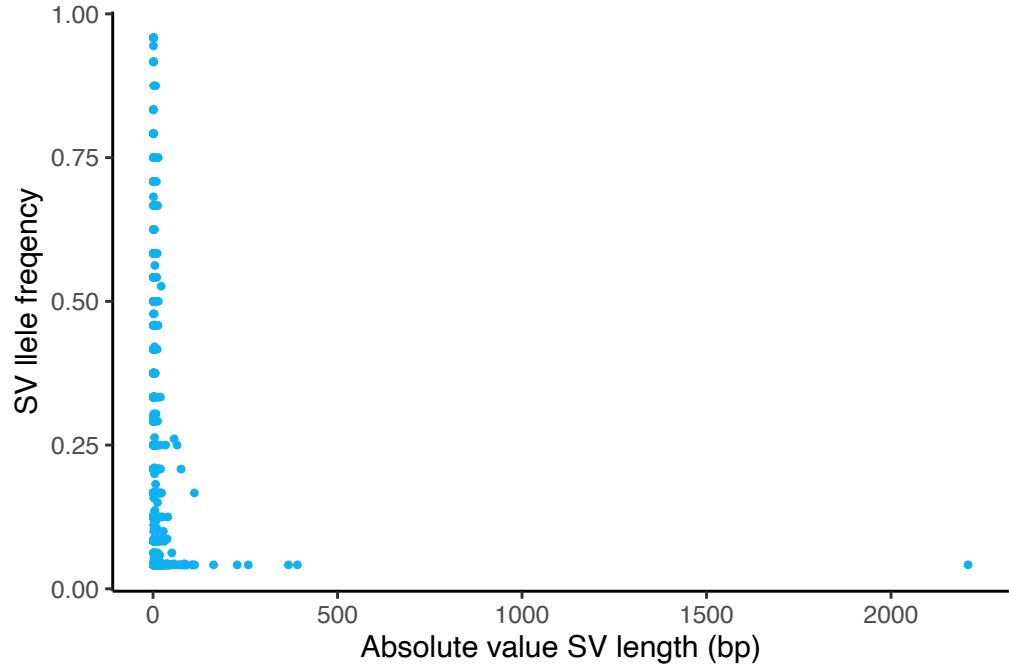


# Low frequency of large SVs in both MHC and hox1a regions (~400-kb)

Structural variants in MHC class II region



Structural variants in hox1a region

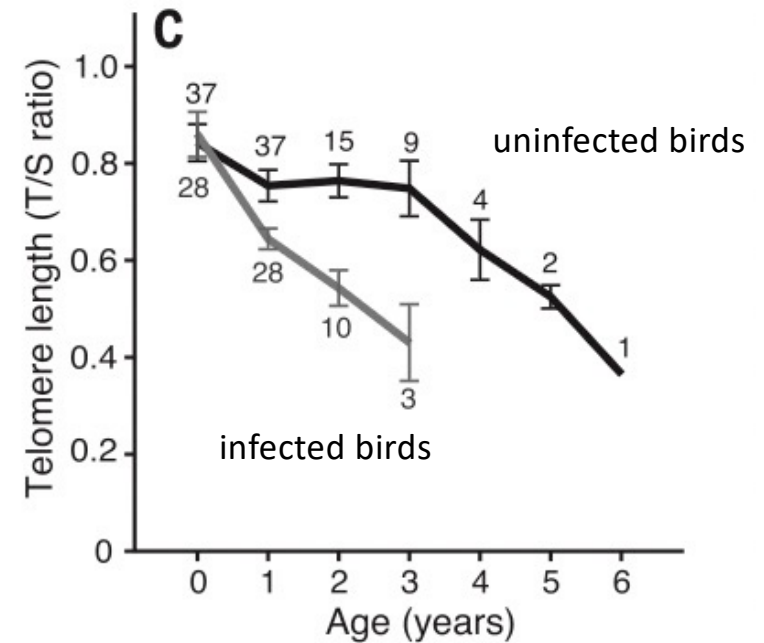
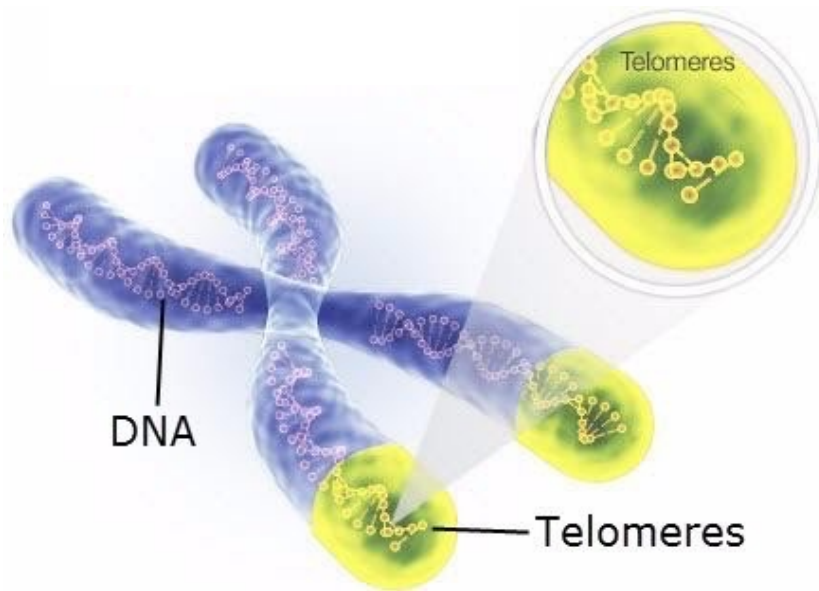


# Telomeres – barometers of age and stress in birds

RESEARCH | REPORTS

## CHRONIC INFECTION

### Hidden costs of infection: Chronic malaria accelerates telomere degradation and senescence in wild birds



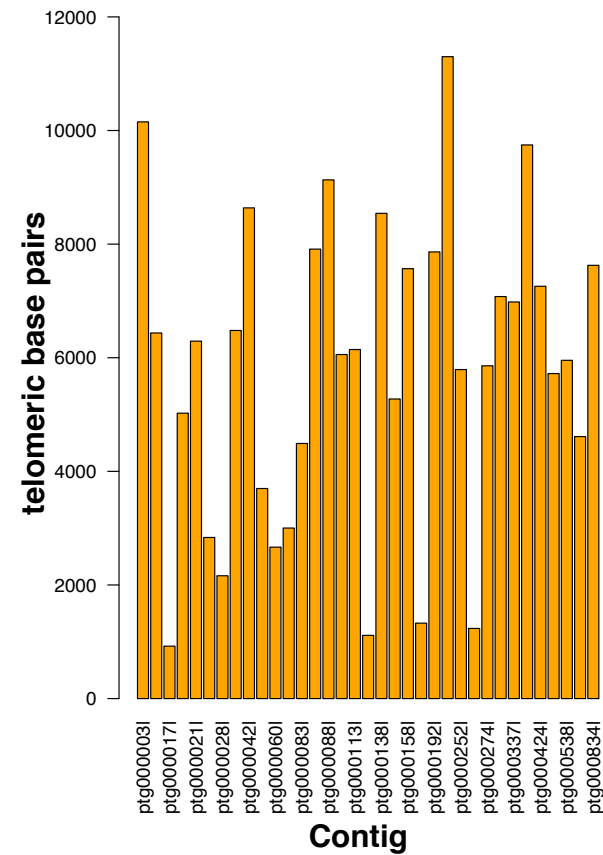
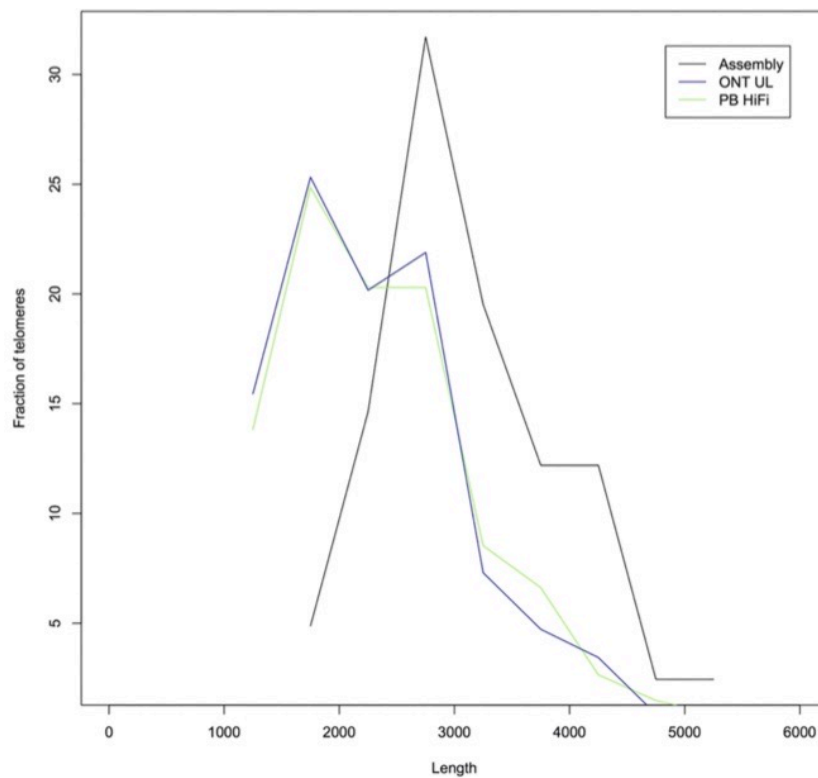
Ashgar et al. 2015. Science 347:436-438

<https://medibalans.com/telomere/>



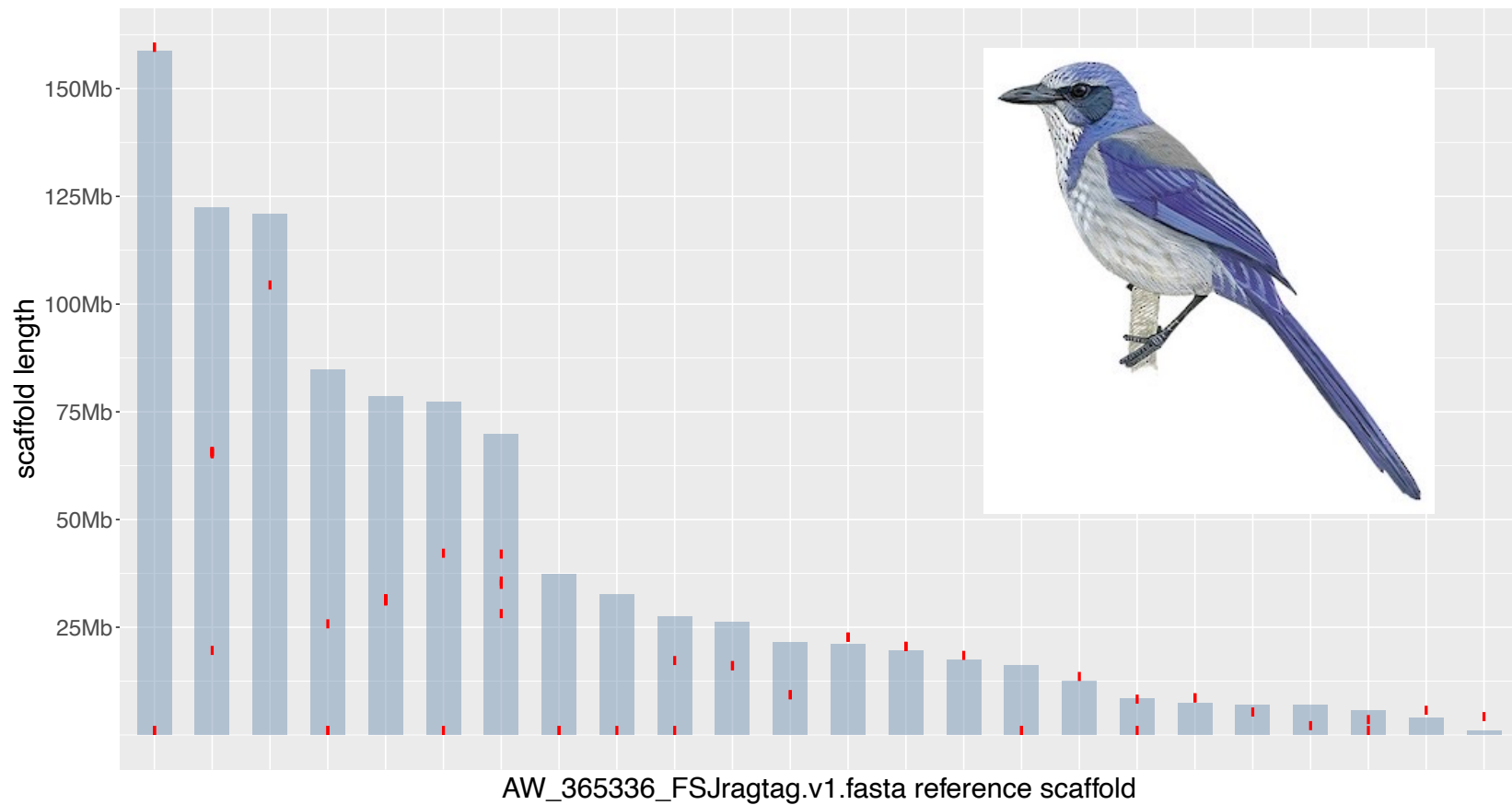


# Scrub jay telomeres are usually ~3-10 kb long

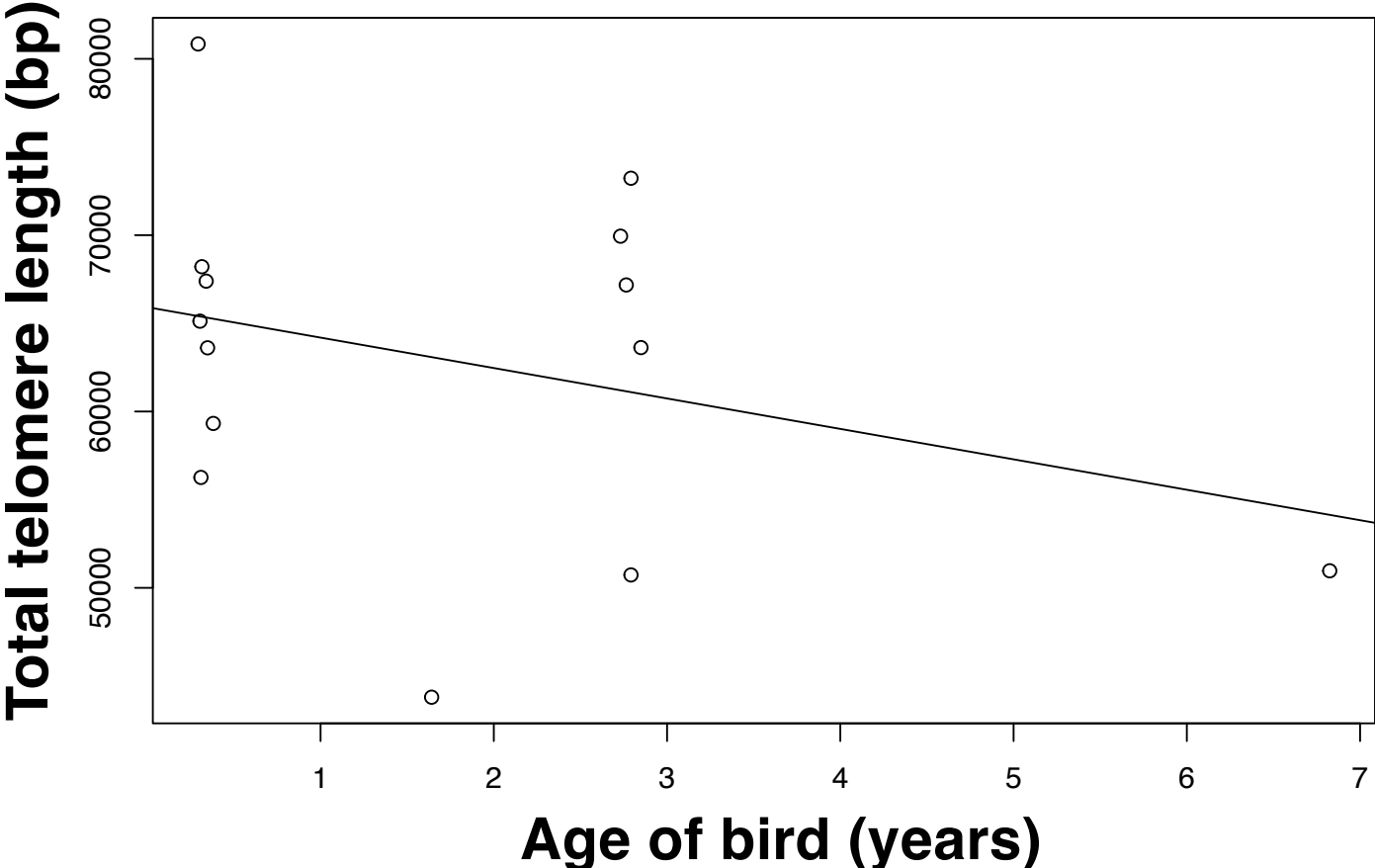


Miga et al. 2020. Nature 485:79-87.

# Telomere sequences are generally found at chromosome ends



# Telomere abundance declines with age in Florida birds



# Conclusions



- Scrub-jay genomes are repeat-rich
- The MHC class II region is much more complex than chicken and likely dispersed on multiple contigs and chromosomes
- Pangenome graph analysis illustrates dynamic and conserved regions of the scrub-jay genome
- Large structural variants appear in lower frequency than small ones
- Pangenome analysis will likely become the common standard

# Acknowledgements

## Colorado team - Island Scrub Jay

Chris Funk  
Rebecca Cheek  
Paul Hohenlohe  
Cameron Ghalambor

## Florida team - Florida Scrub Jay

Nancy Chen  
Reed Bowman  
John Fitzpatrick

## Harvard team - Woodhouse's Scrub Jay and Informatics

Tim Sackton  
Danielle Khost  
Heng Li  
Bohao Fang  
George Kolyfetis

## Pangenome informatics

Erik Garrison  
Andrea Guarracino



## Fieldwork

Greg and Donna