

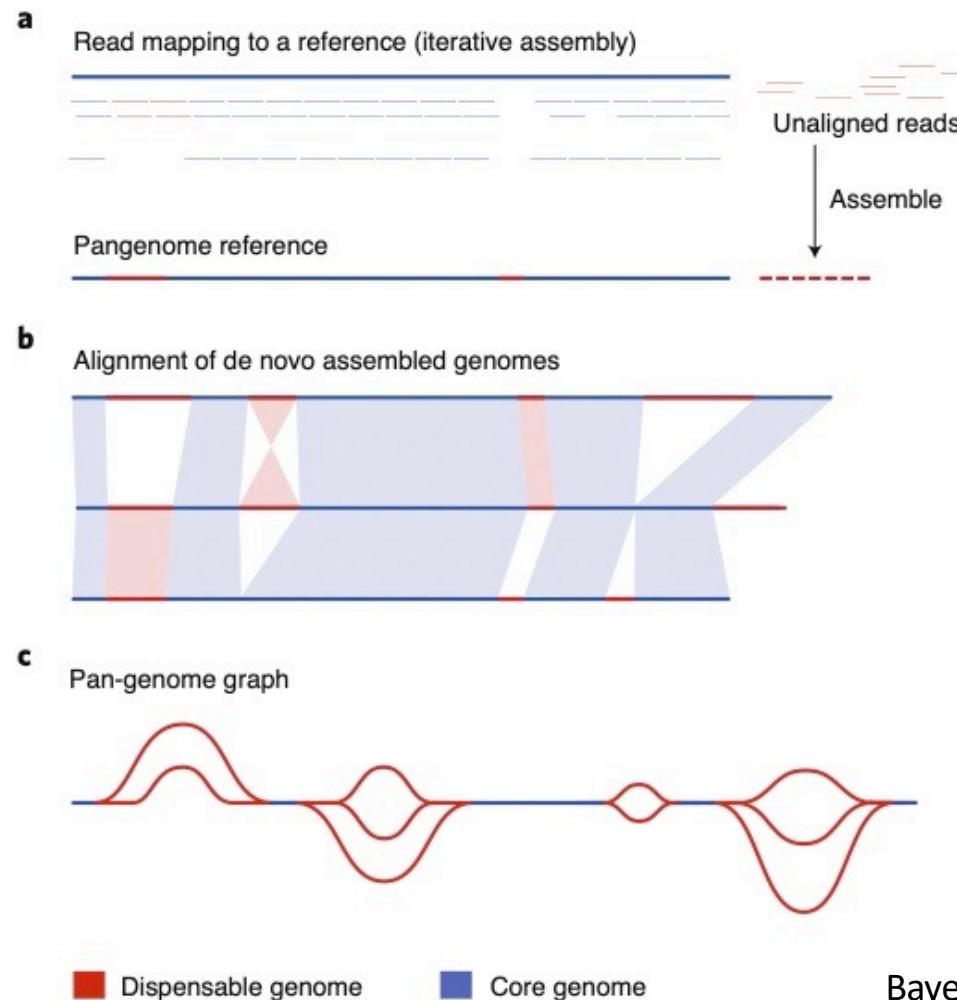
# Pangenomes as a new tool for studying ecology and evolution of natural populations

Scott V. Edwards

Museum of Comparative Zoology, Harvard University, Cambridge, USA



# Pangenomes: moving beyond reference-based genomics

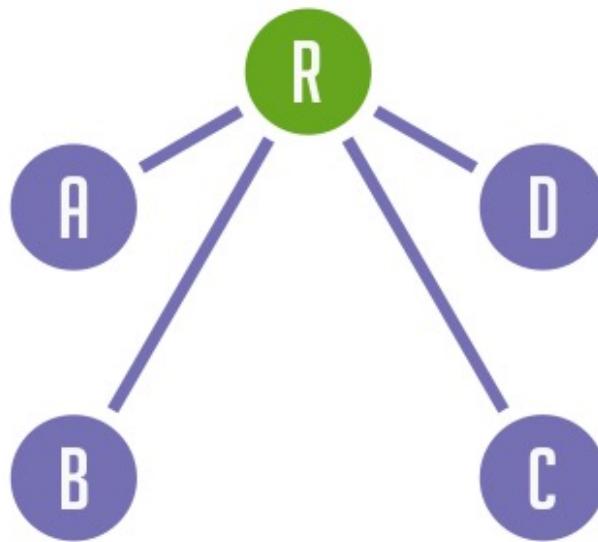


Bayer et al. 2020. *Nature Plants* 6: 914-920.

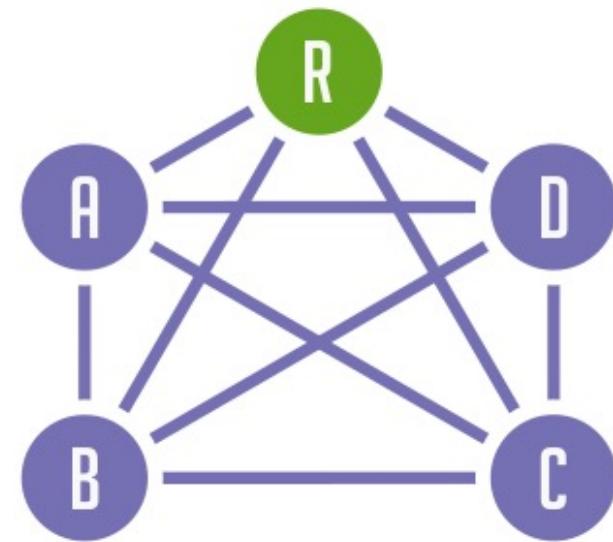
# Reference-free genomics

Reference model

Genomic



Pangenomic



Eizenga et al. 2021. *Ann. Rev. Genomics and Human Genetics*

A photograph of a sunset over a body of water, likely a lake or coastal area. The sky is filled with warm orange, yellow, and pink hues. In the foreground, dark, textured rocks are silhouetted against the water. A single bird, possibly a gull, stands on one of the rocks, facing the horizon. The water reflects the colors of the sky.

# Pangenomes

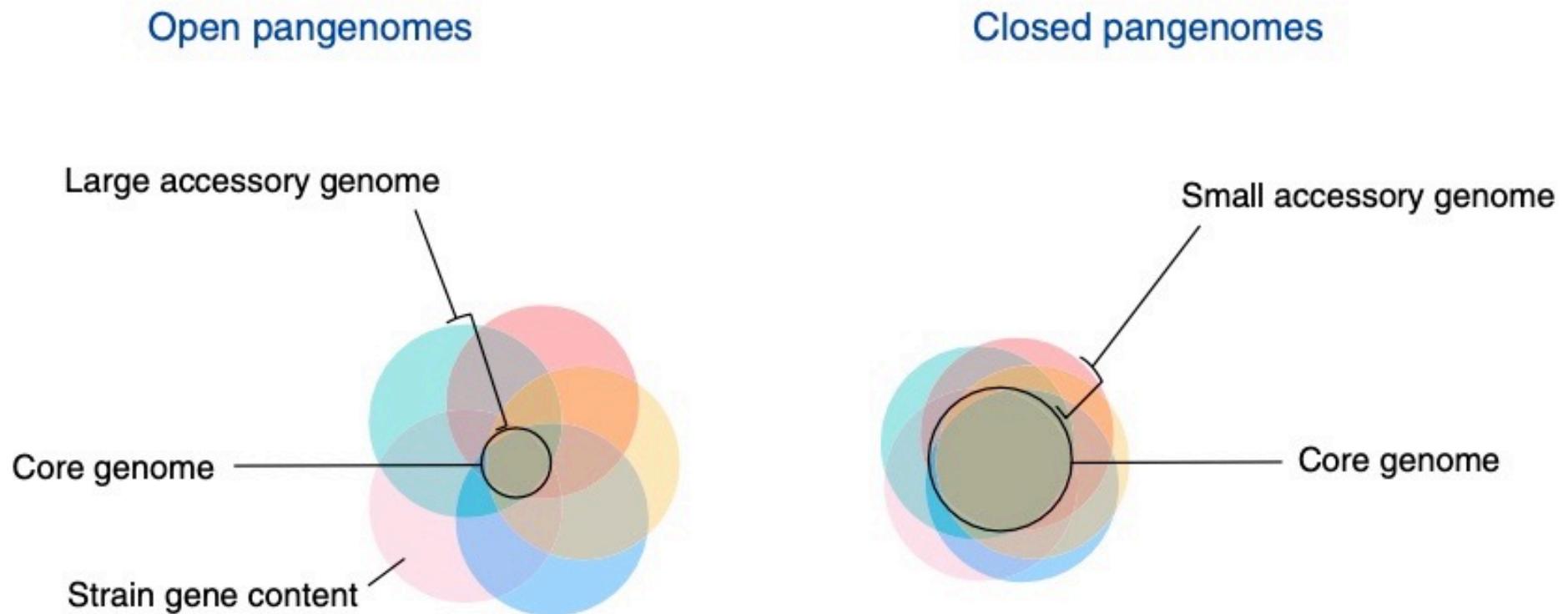
## Evolution and Computation

PGEC 2021

November 17, 2021 University of Gothenburg

<https://pgec2021.schlieplab.org/>

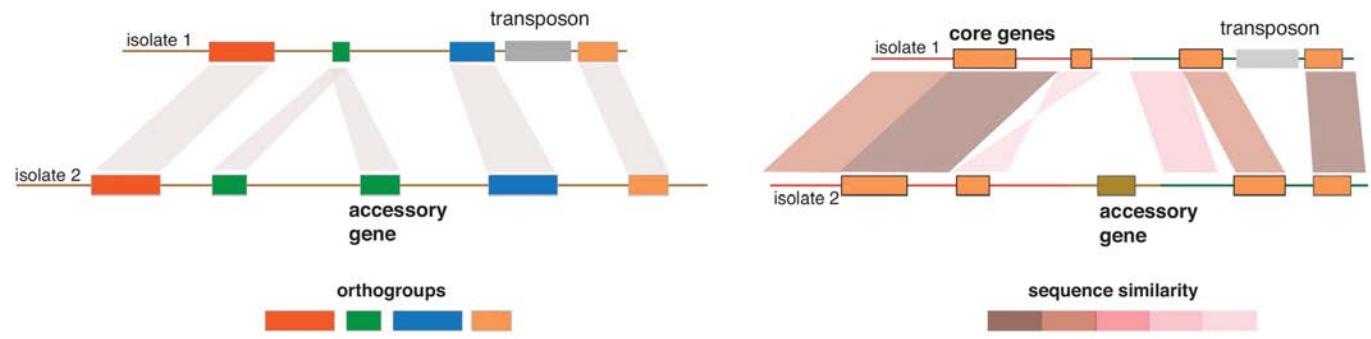
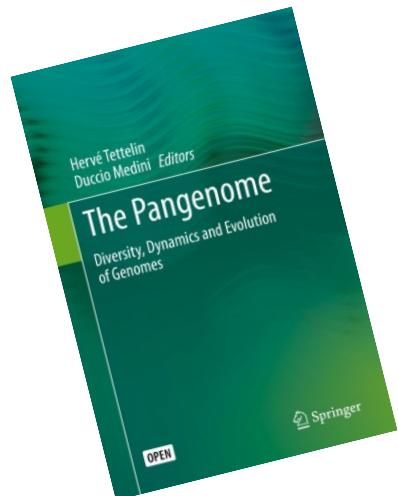
# Open and closed pangenomes



Brokhurst et al. 2019. *Curr. Biol.*

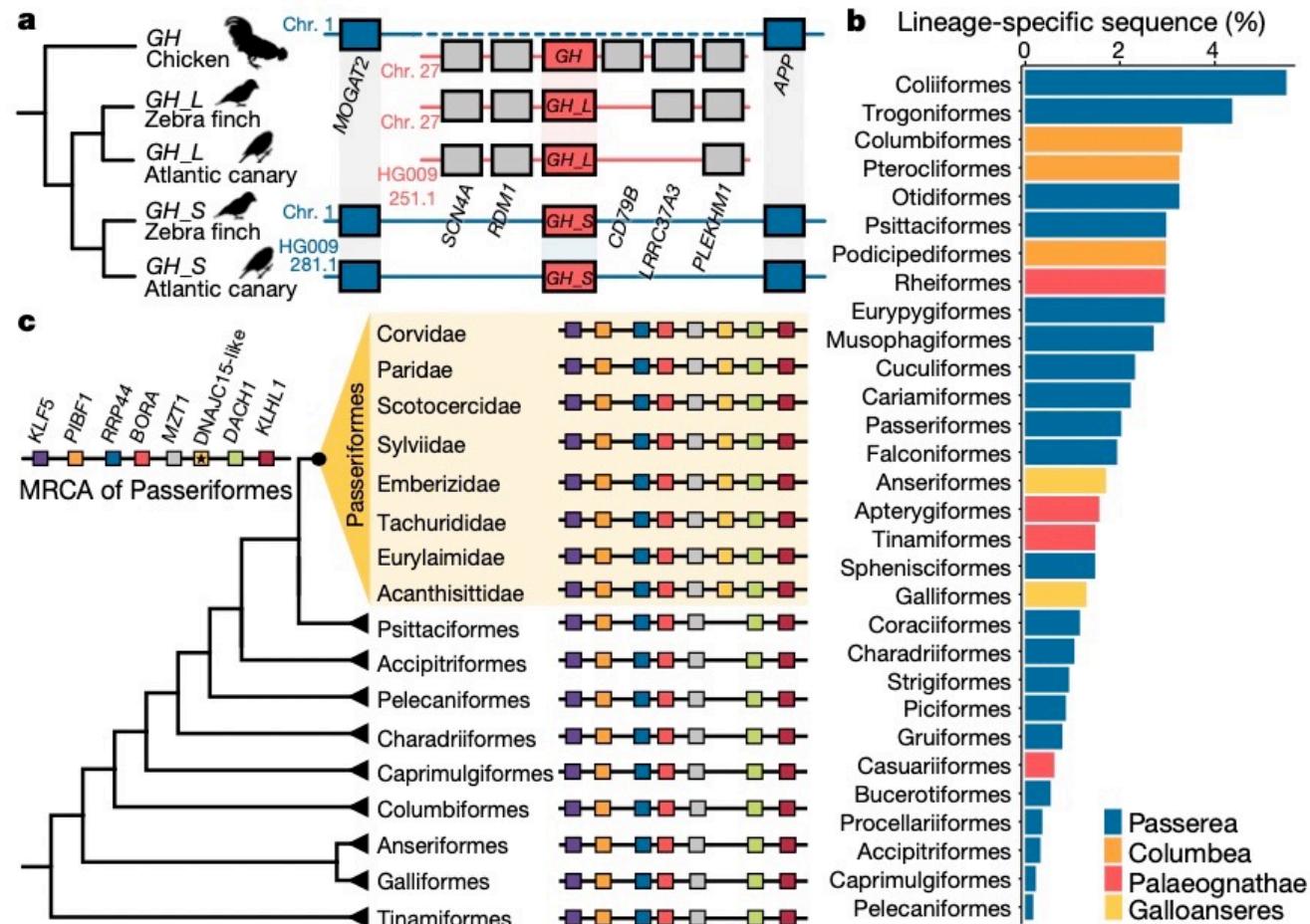
# The eukaryotic pangenome

- “The existence of pangenesomes in eukaryotes is debated...Pangome 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 bird evolution



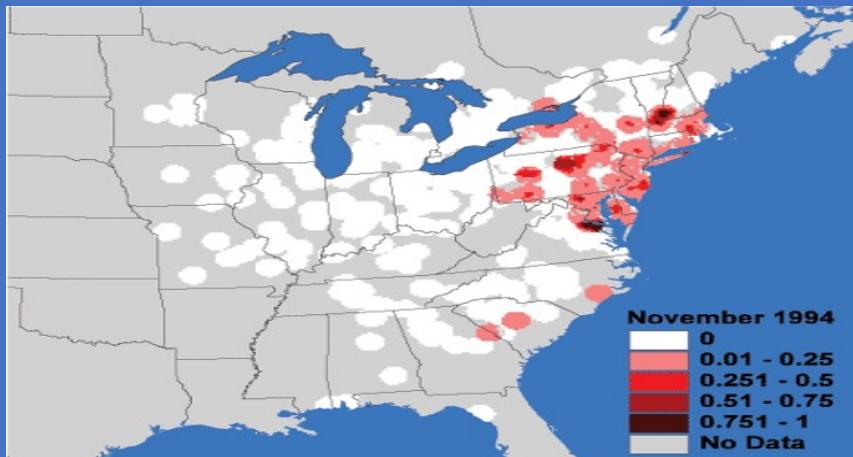
Feng et al. 2020. *Nature* 587:252-257.



# Recent history of House Finch populations



# 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

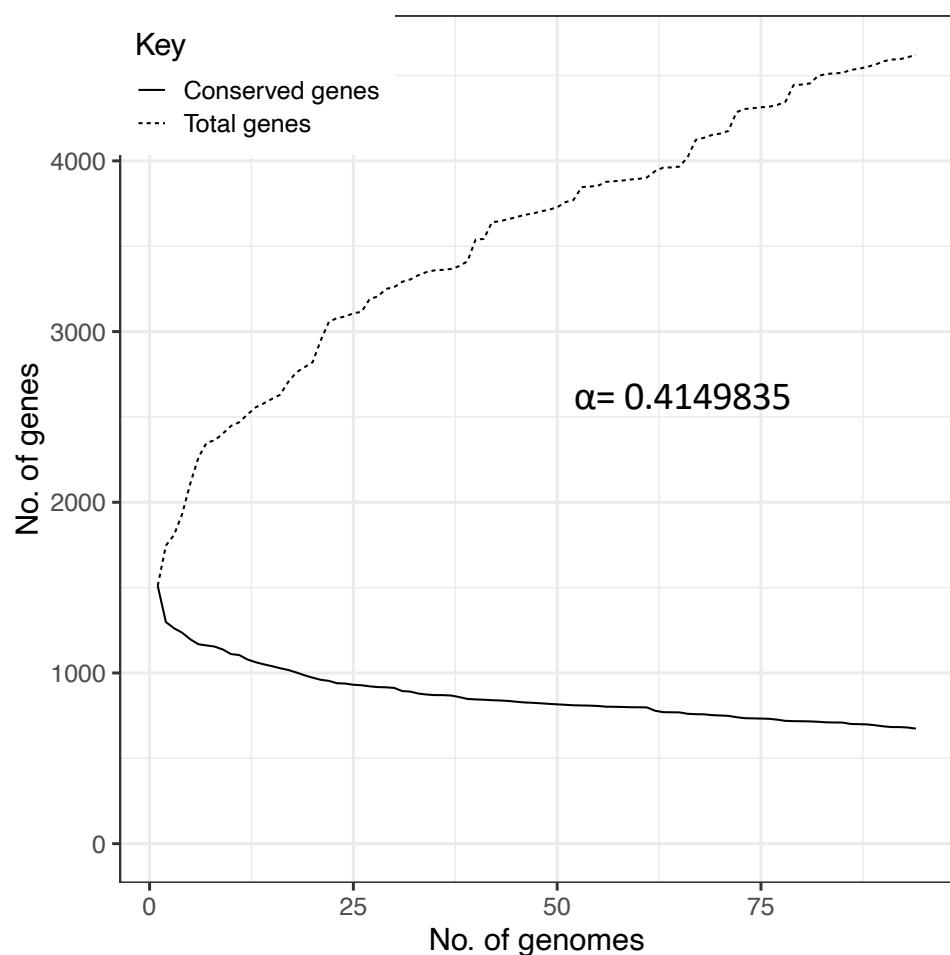


Analyzed 81 Mycoplasma strains from chicken, turkey and house finch, available on NCBI

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

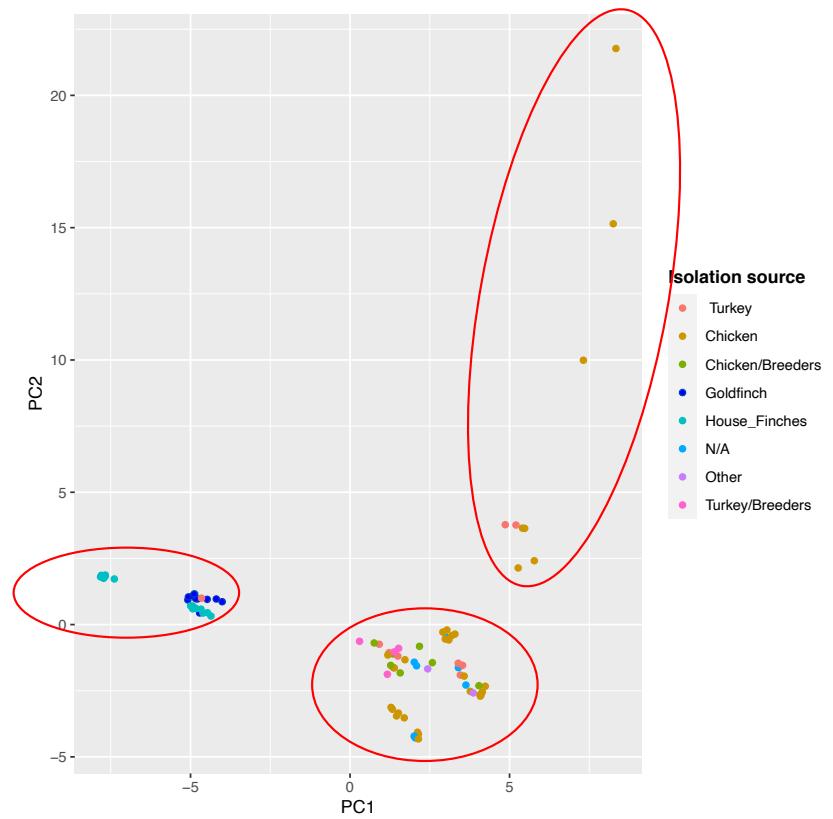
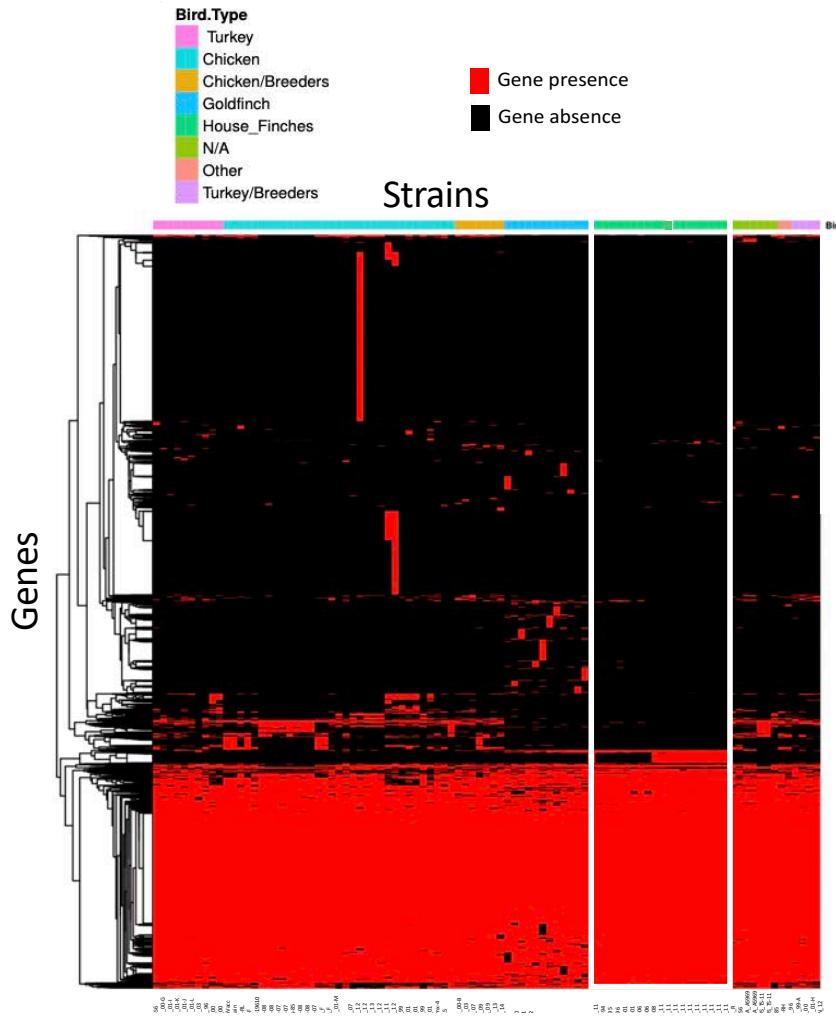
Feature	Info	Number of genes	Percentage
Core genes	(99% <= strains <= 100%)	674	14.586
Soft core genes	(95% <= strains < 99%)	464	10.041
Shell genes	(15% <= strains < 95%)	412	8.916
Cloud genes	(0% <= strains < 15%)	3071	66.457
SGF	one copy in all strains	141	3.051
SGF	without recombination signals	117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

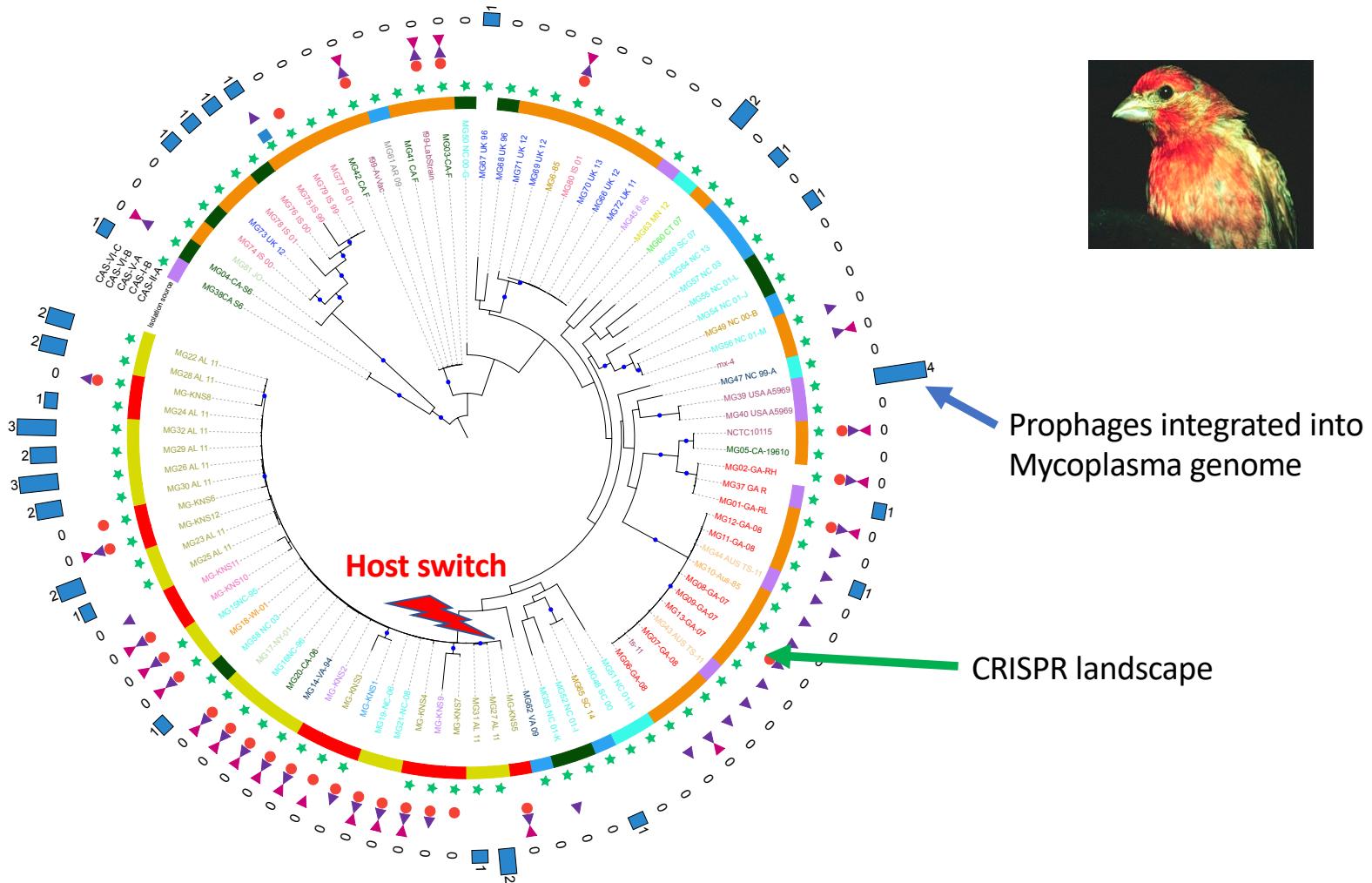
New data: Determine the alpha value using MicroPan

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

# Mycoplasma pangenome gene repertoire is highly strain-specific

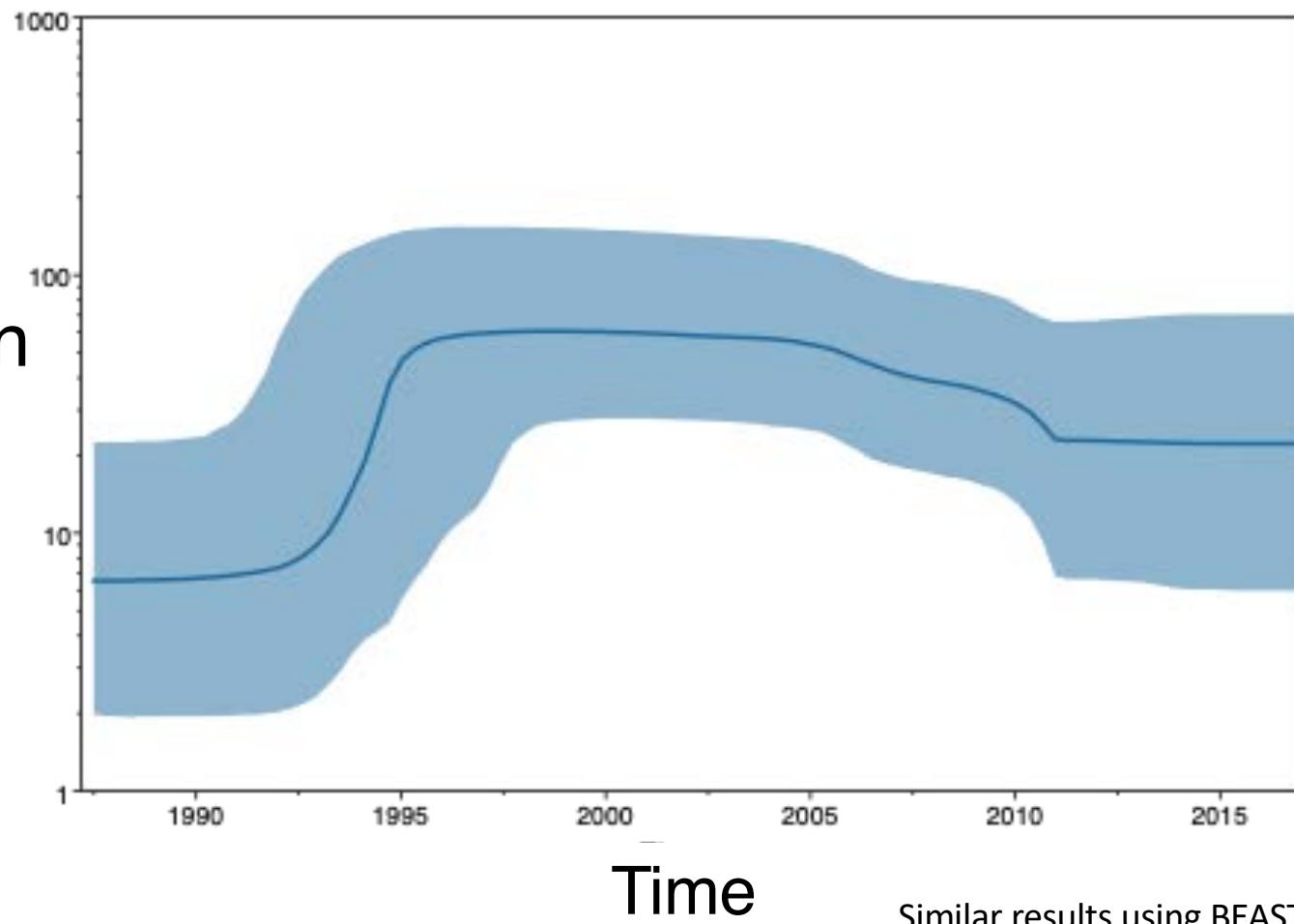


## House Finch *Mycoplasma* strains have distinct CRISPR and prophage landscapes



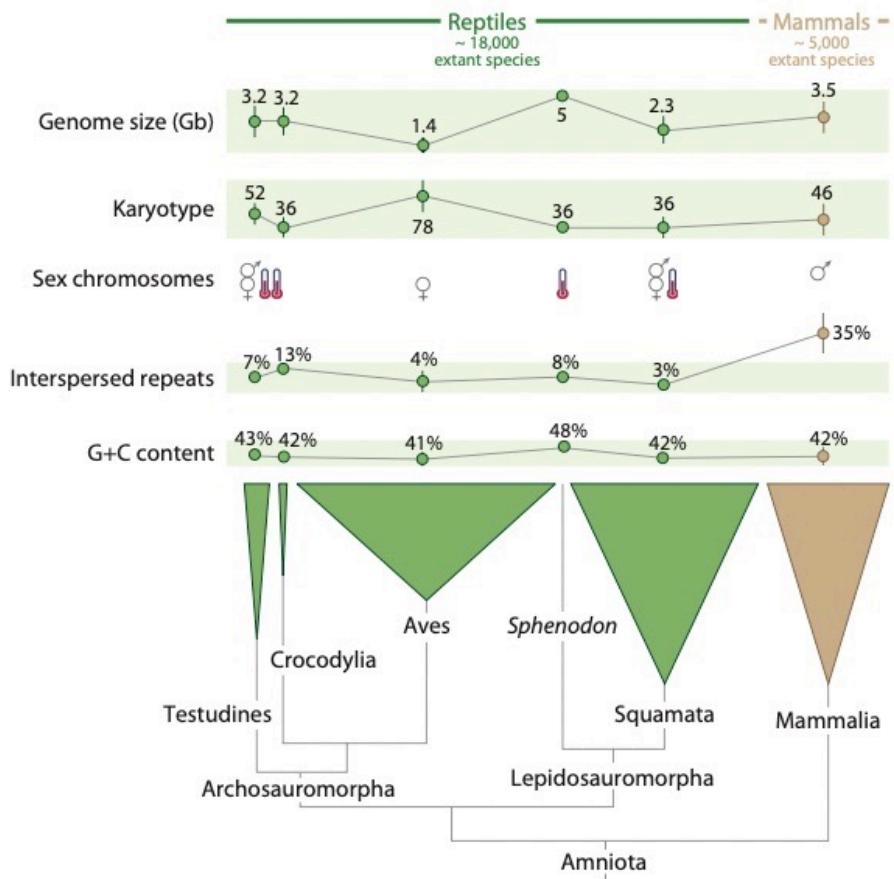
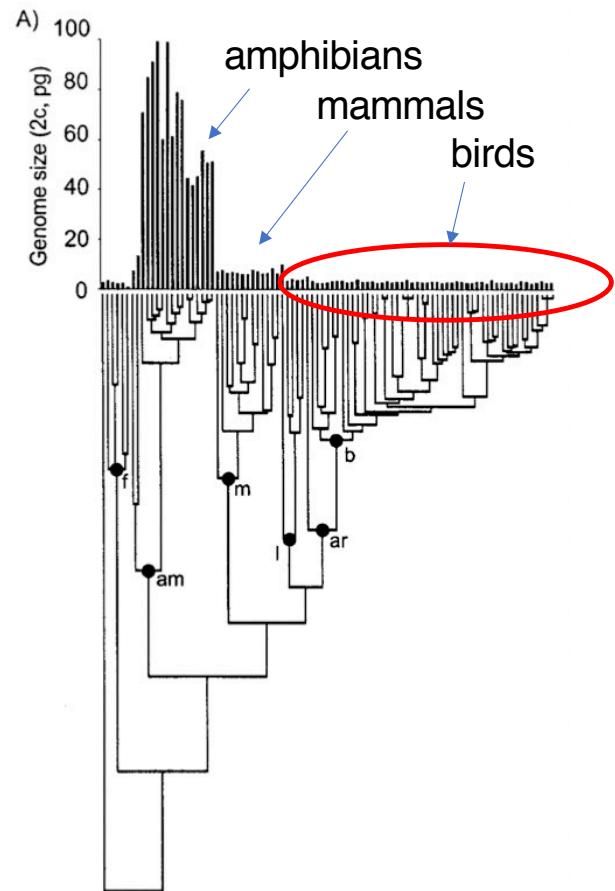
*Mycoplasma* epizootic likely began ~2 years before first detection

Effective population size



Similar results using BEAST and Stairway plot

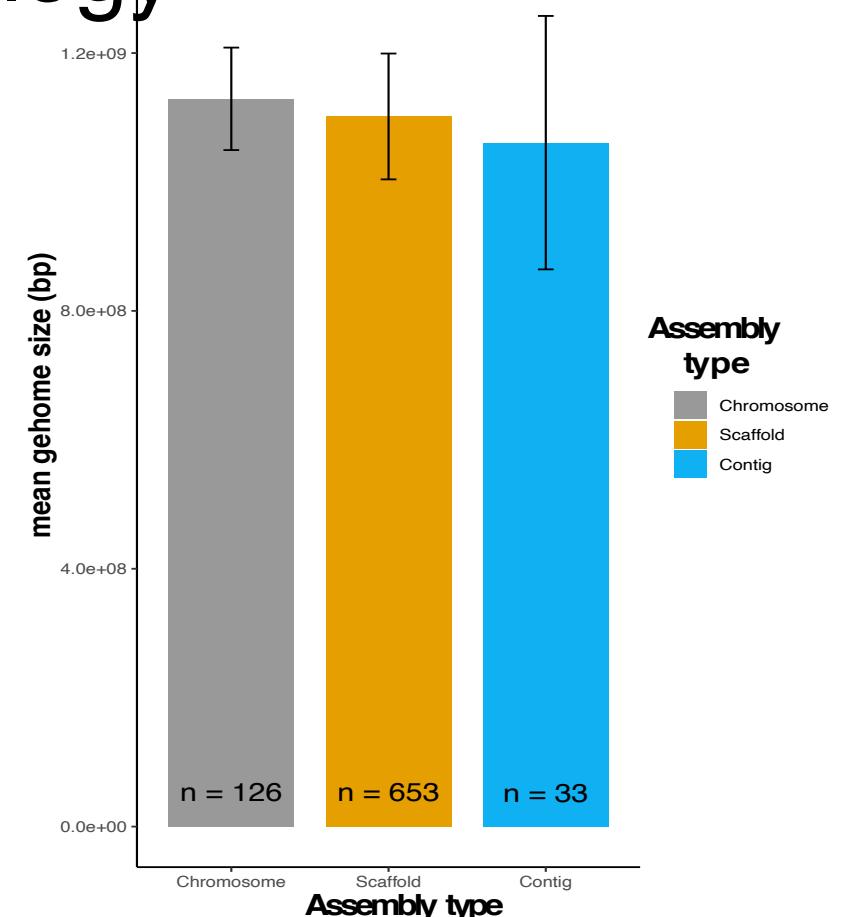
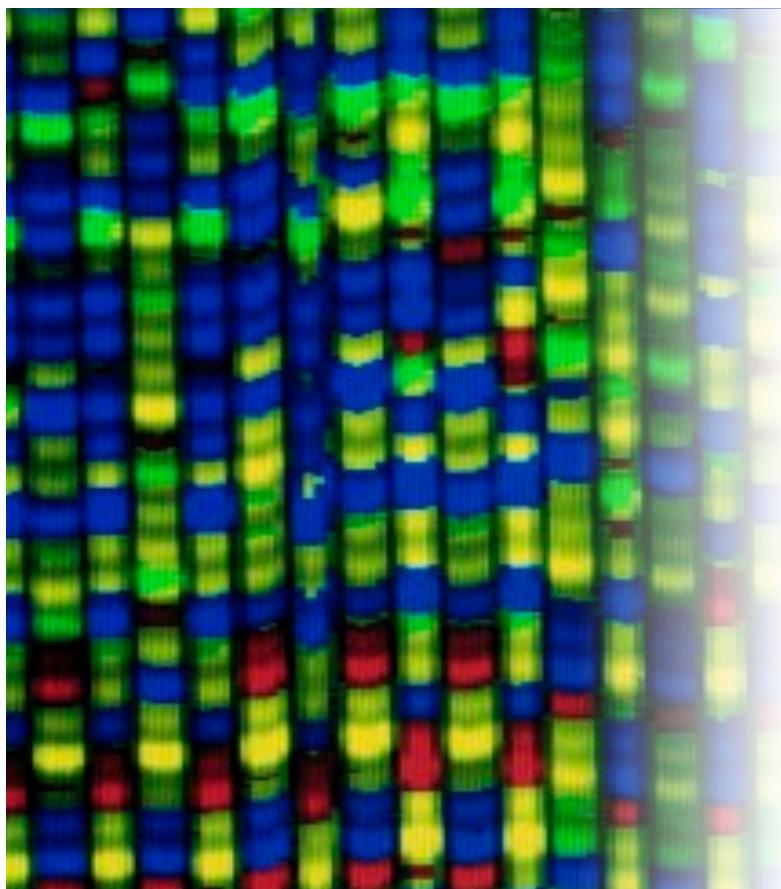
# Birds have small, streamlined genomes



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

decimal latitude  
45 -  
40 -  
35 -  
30 -  
25 -



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

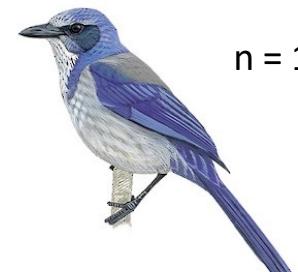
n = 15

Datapoints from gbif.org



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

n = 15



# 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**

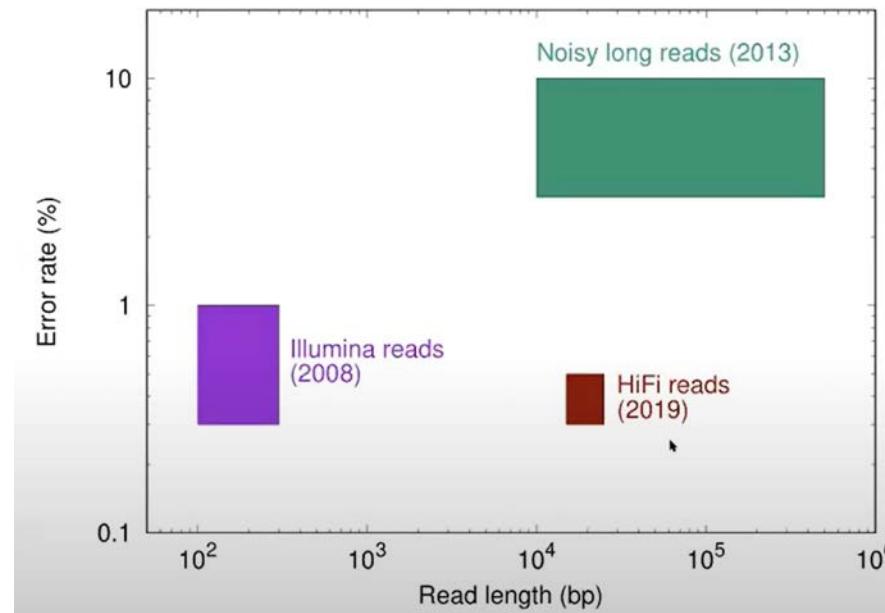
Conceptual Relationships between the Fields of Comparative Population Genomics, Landscape Genomics, and Comparative Phylogeography

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

Edwards et al. 2021. *Genome Biology and Evolution* 14: 10.1101/gbe/evab176

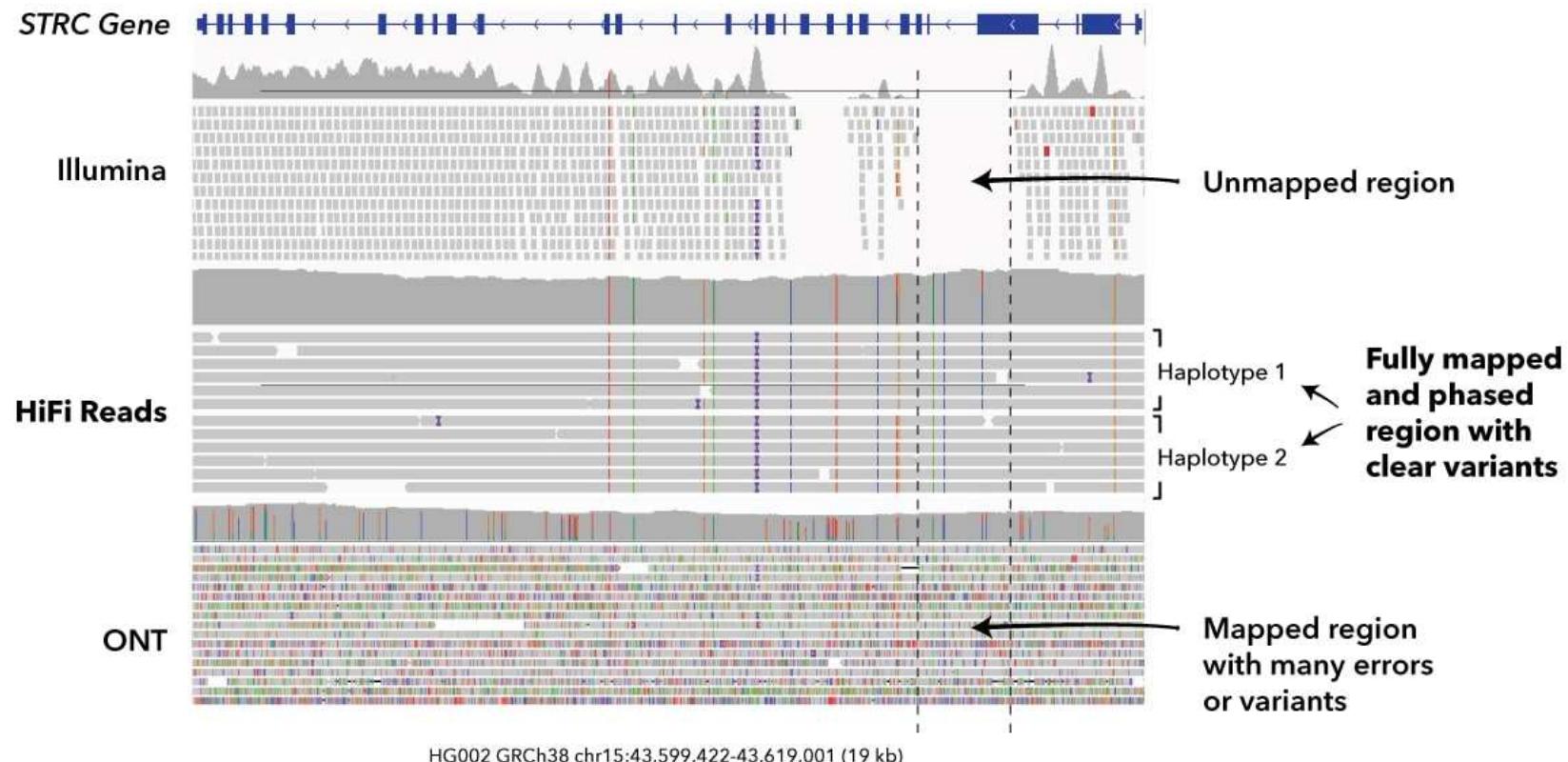
# PacBio HiFi reads are long and accurate

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

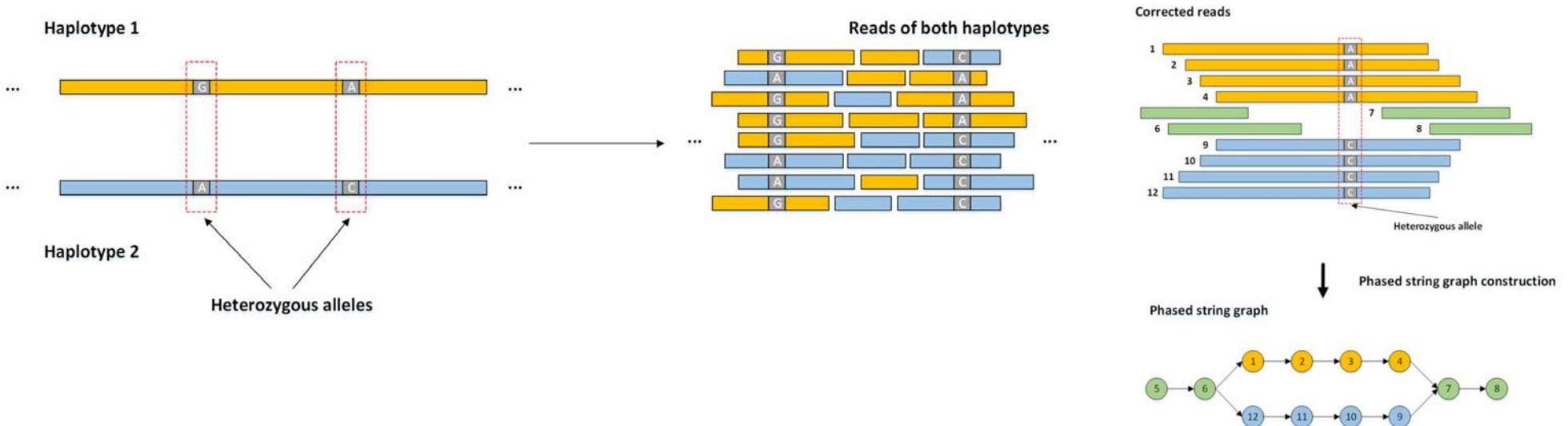


Courtesy Haoyu Cheng, Dana Farber Cancer Institute

# PacBio HiFi reads are long and accurate

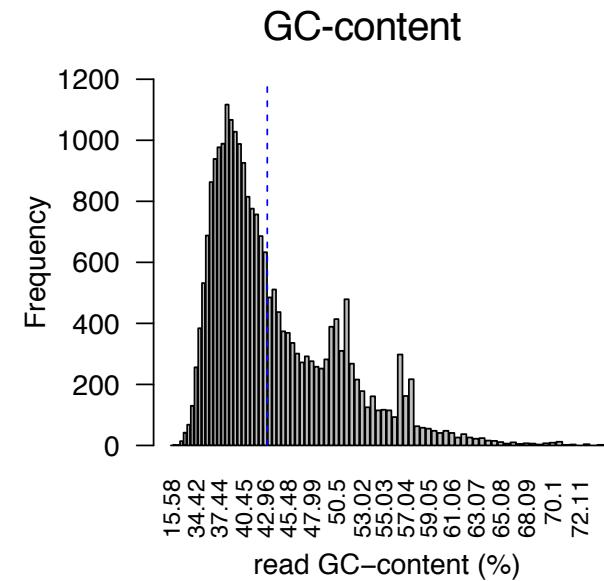
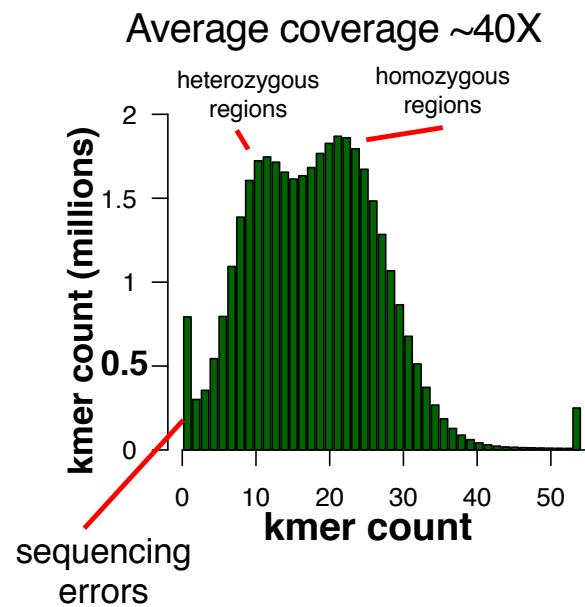
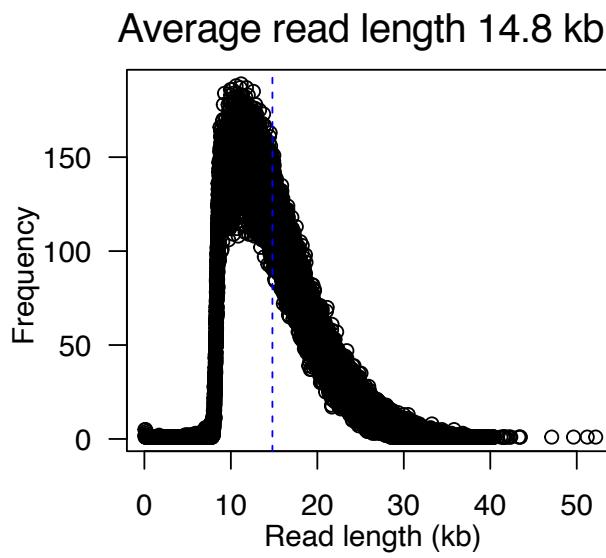


# Hifiasm – a HiFi accurate read assembler that resolves haplotypes

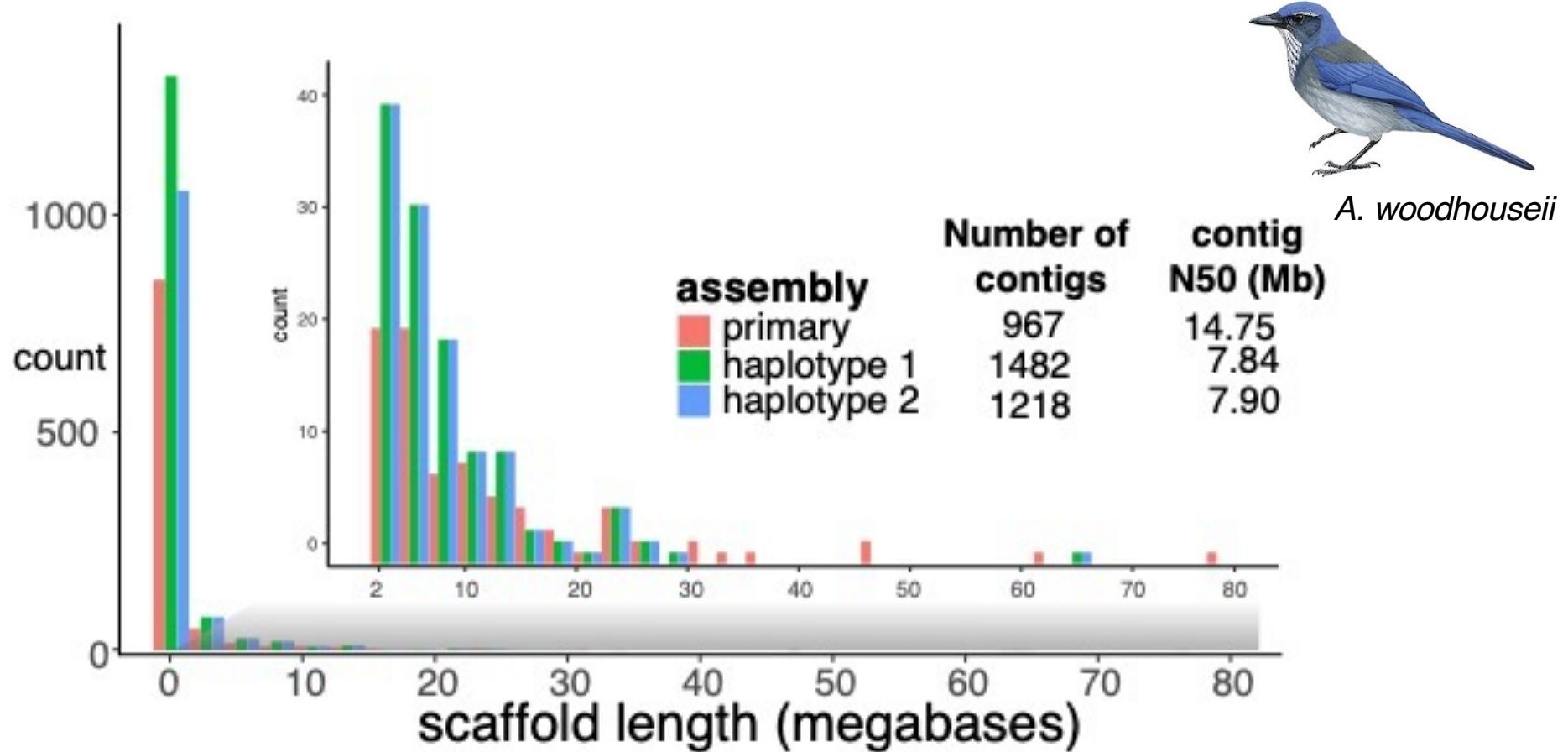


Courtesy Haoyu Cheng, Dana Farber Cancer Institute

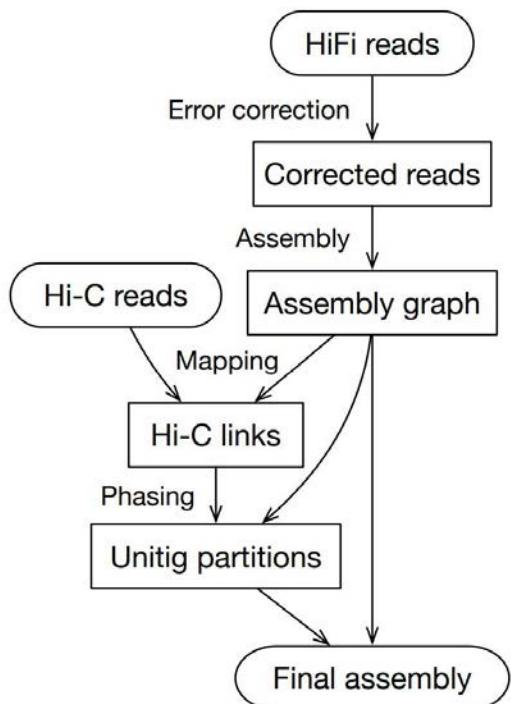
# Scrub-jay PacBio HiFi data characteristics



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

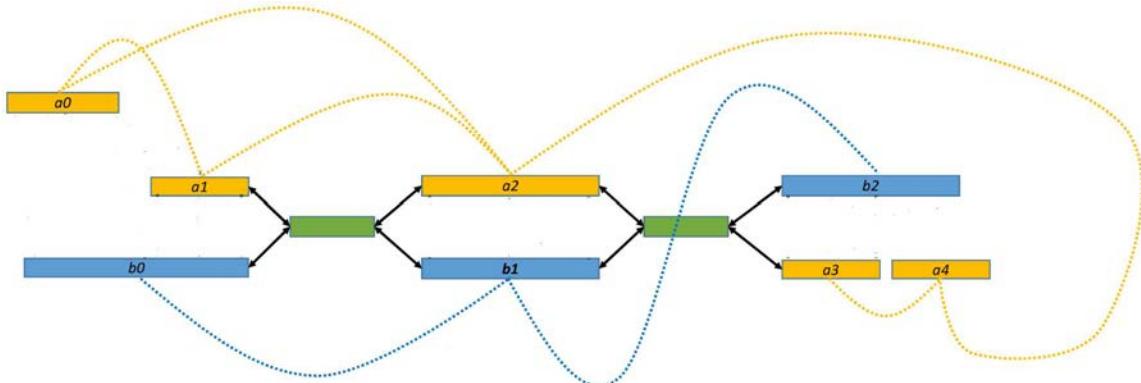


# Hifiasm – improved assemblies using HiC

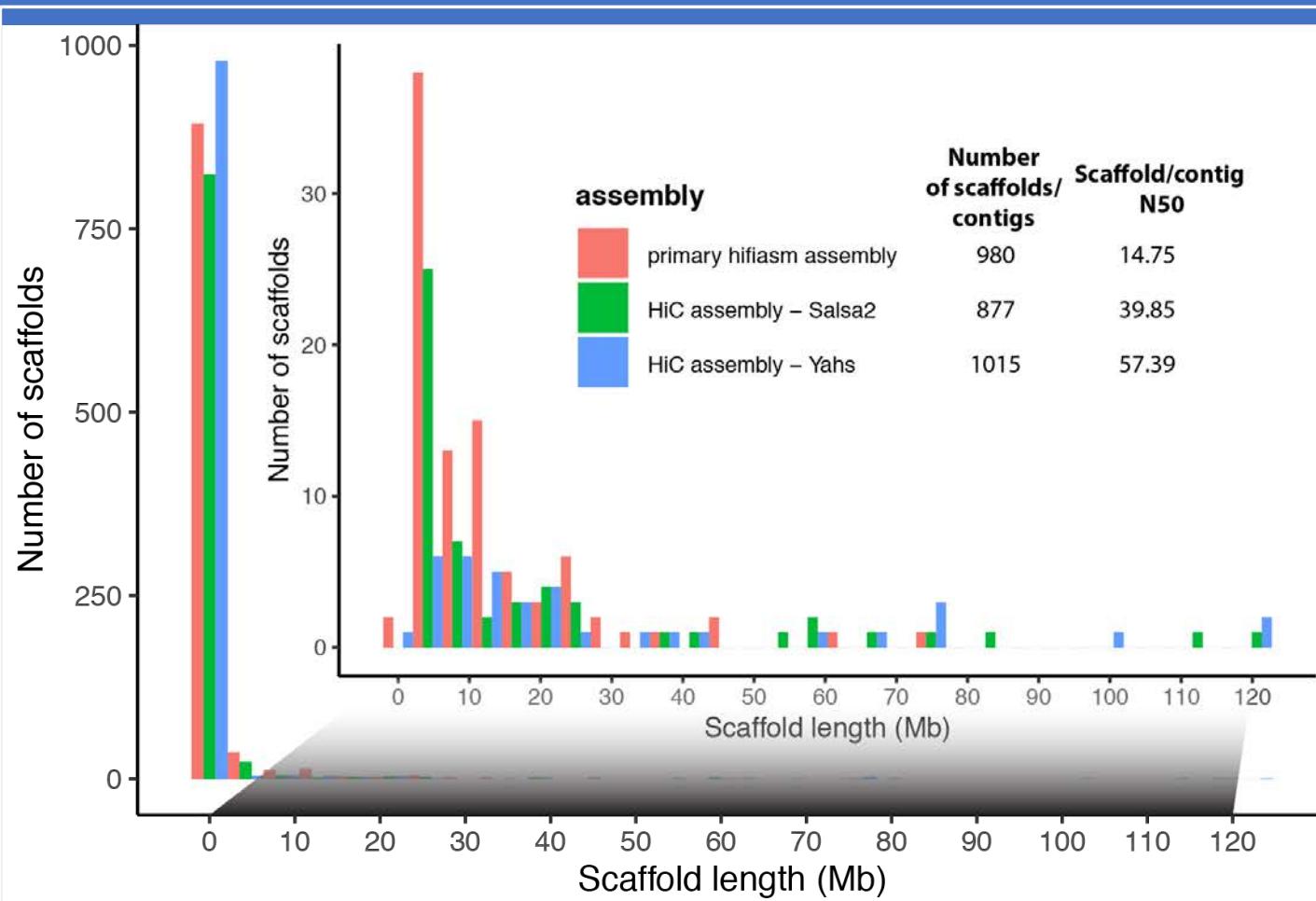


## Procedure:

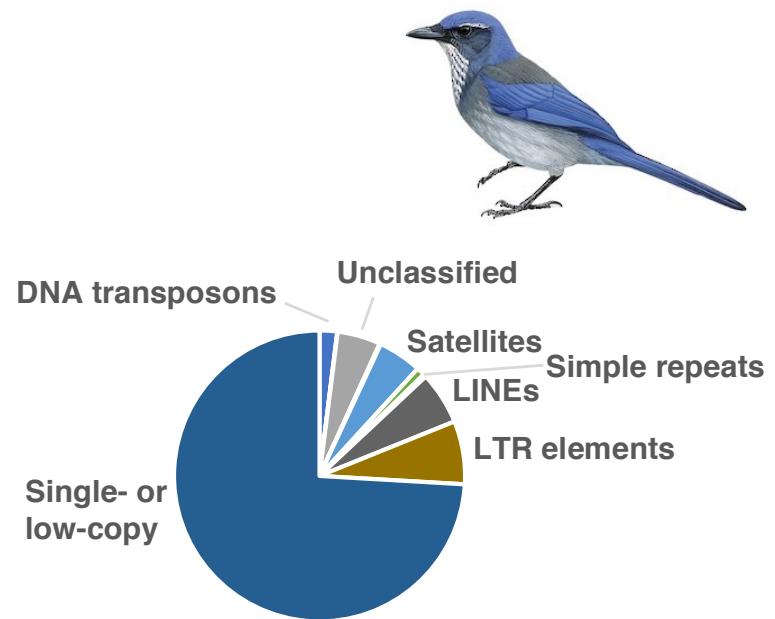
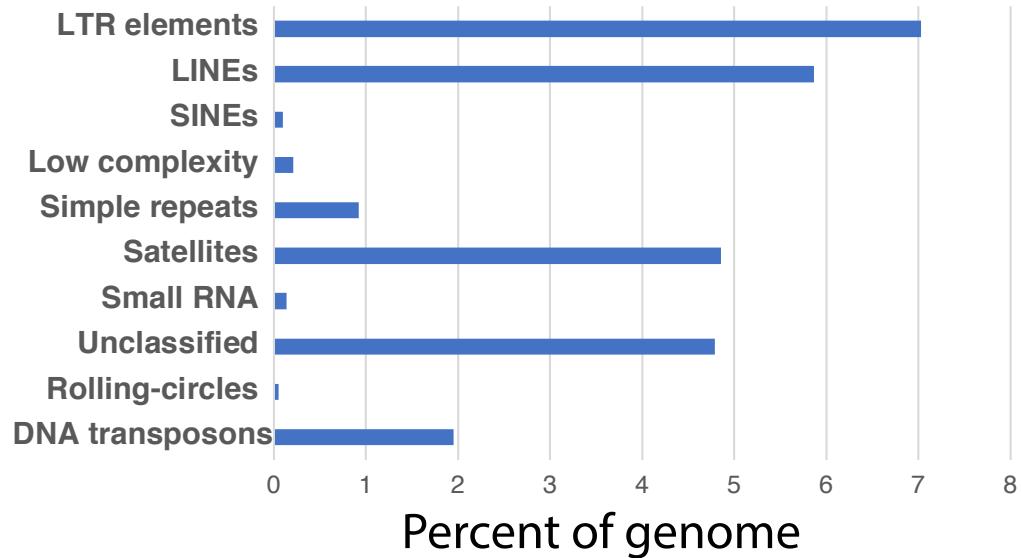
- ▶ Identify heterozygous unitigs by coverage
- ▶ Build index by unique  $k$ -mers from heterozygous unitigs
- ▶ Align Hi-C reads using unique  $k$ -mers



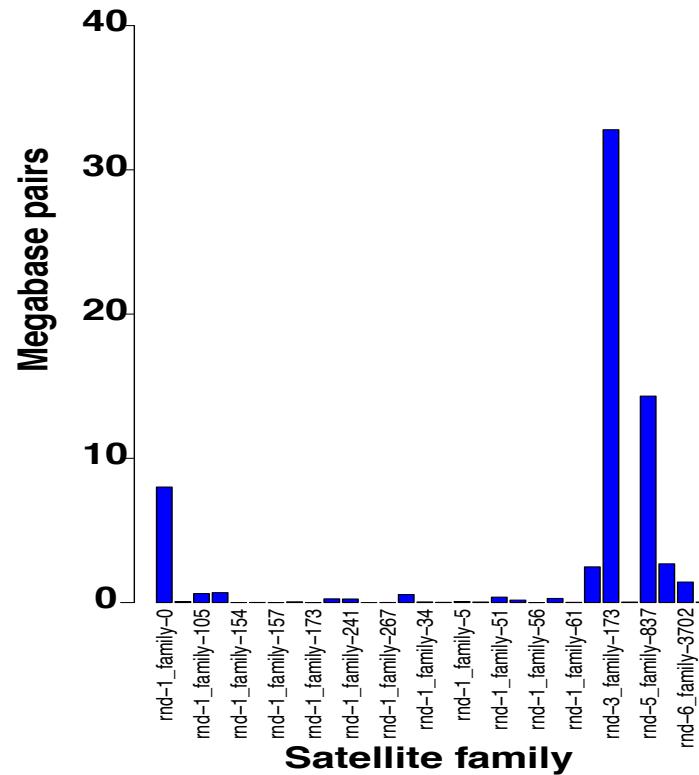
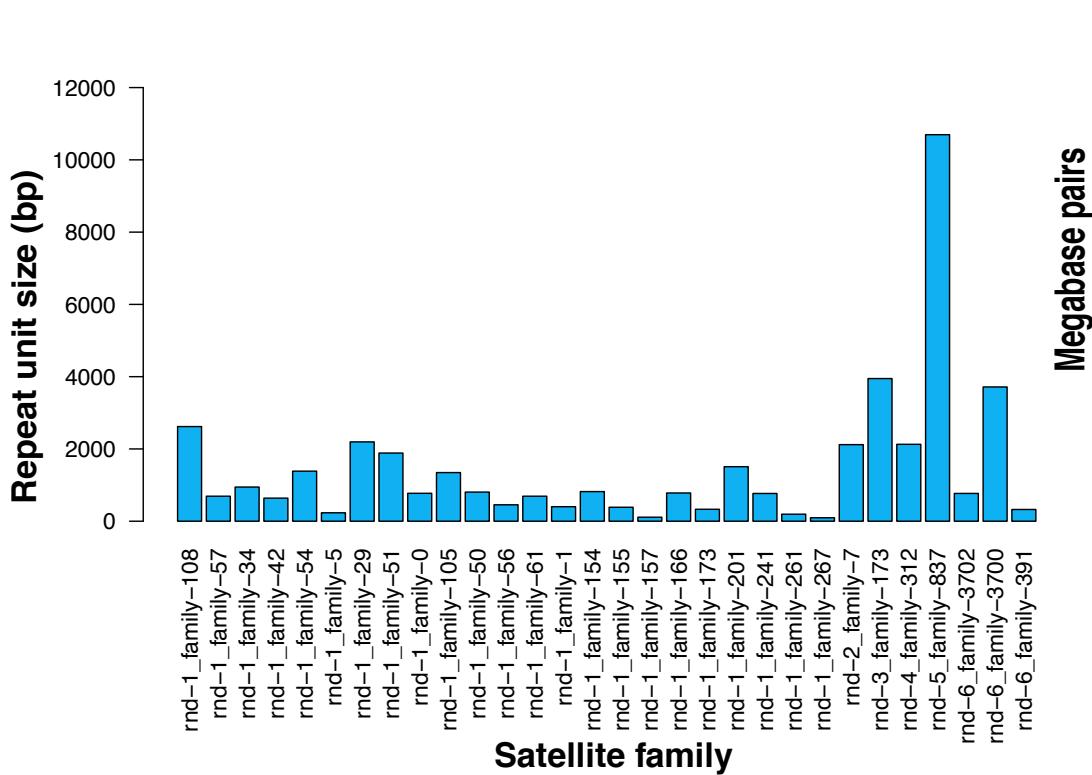
# HiC greatly improves contiguity of scrub jay assemblies



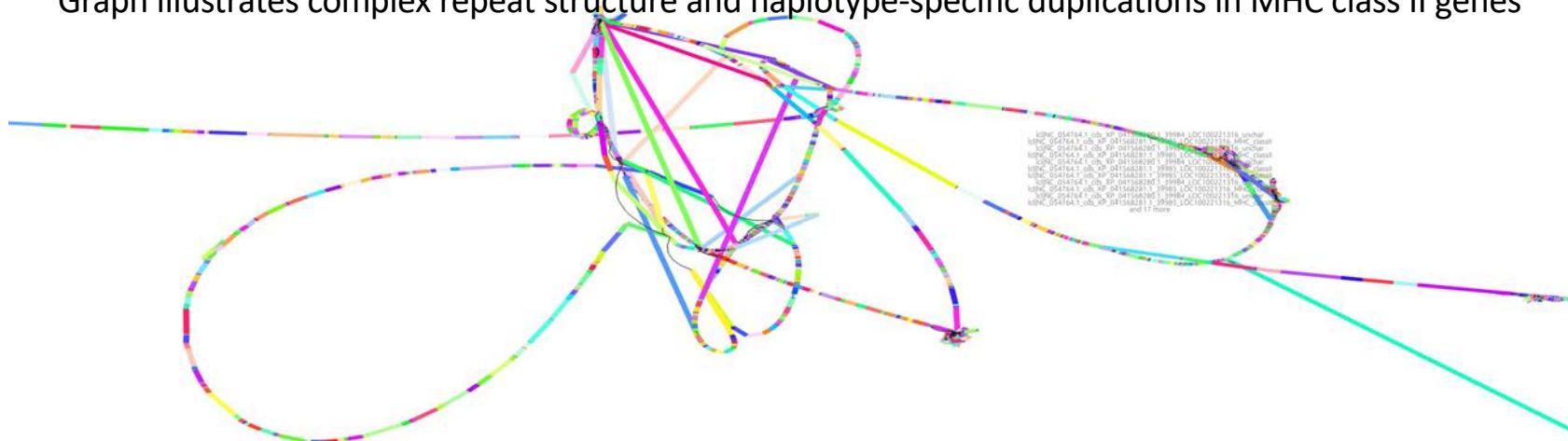
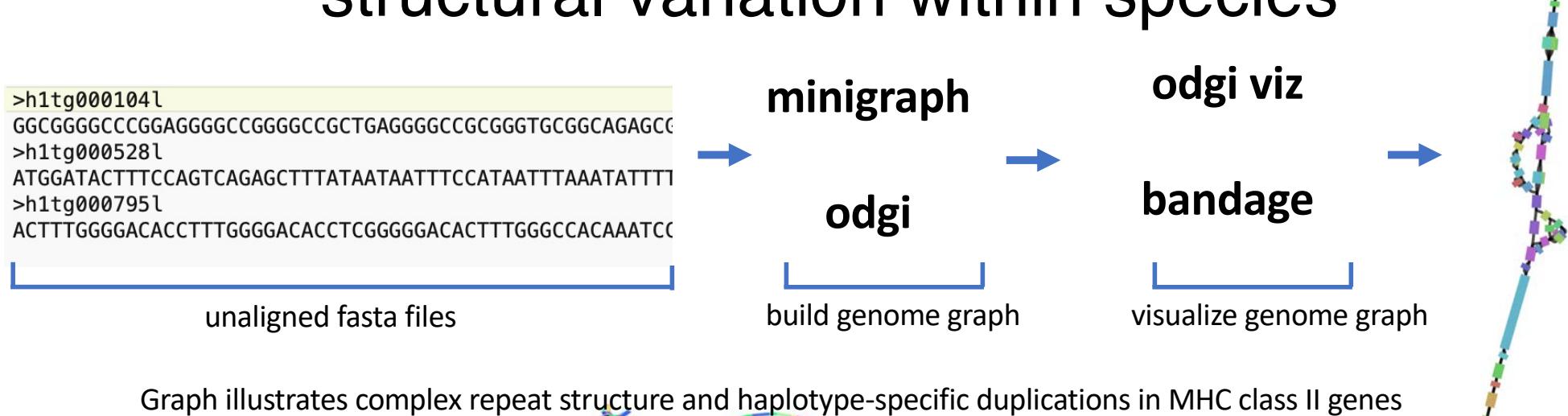
RepeatMasker analysis suggests over 25% repeats and transposable elements



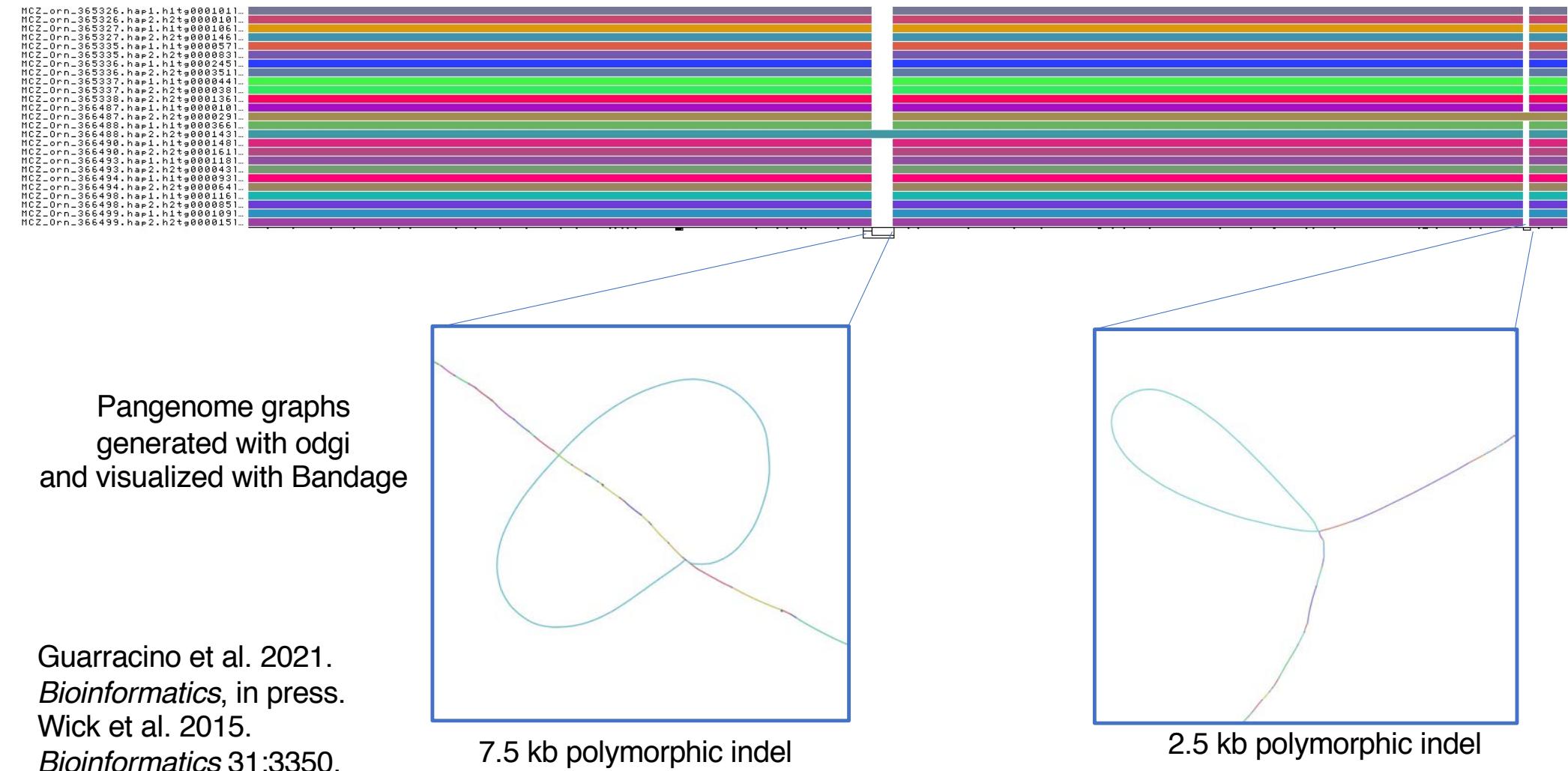
# Satellites are long and prevalent in scrub jay genomes



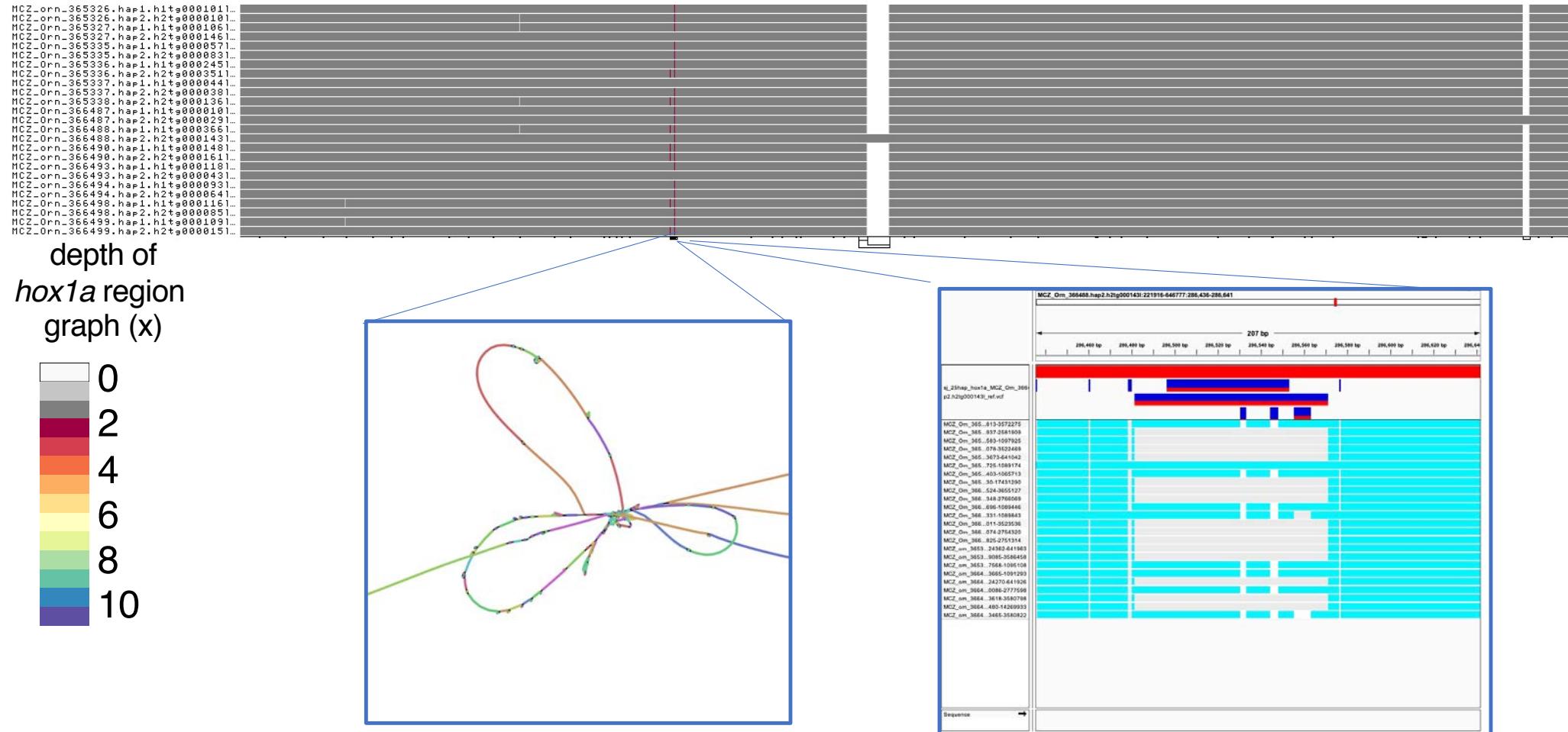
# Pangenome graphs capture structural variation within species



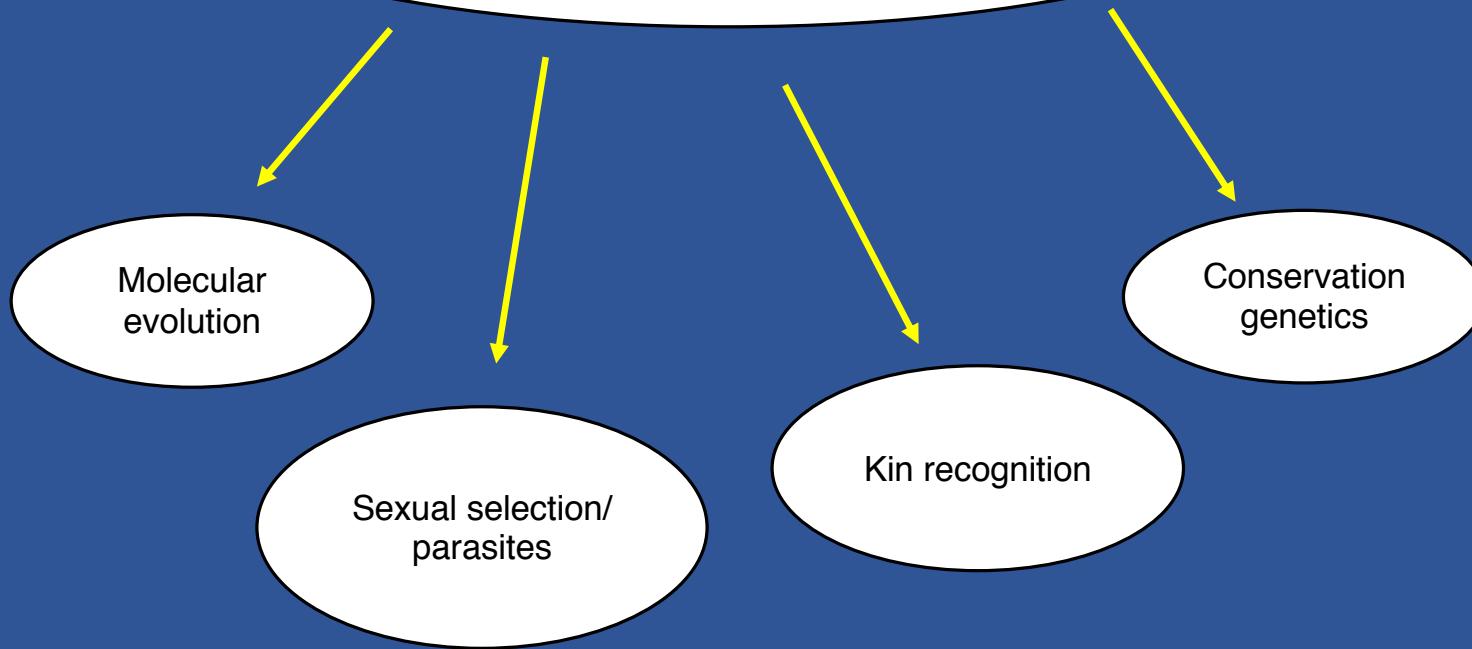
# Genomic stability of 400-kb hox1a region in Western Scrub Jays



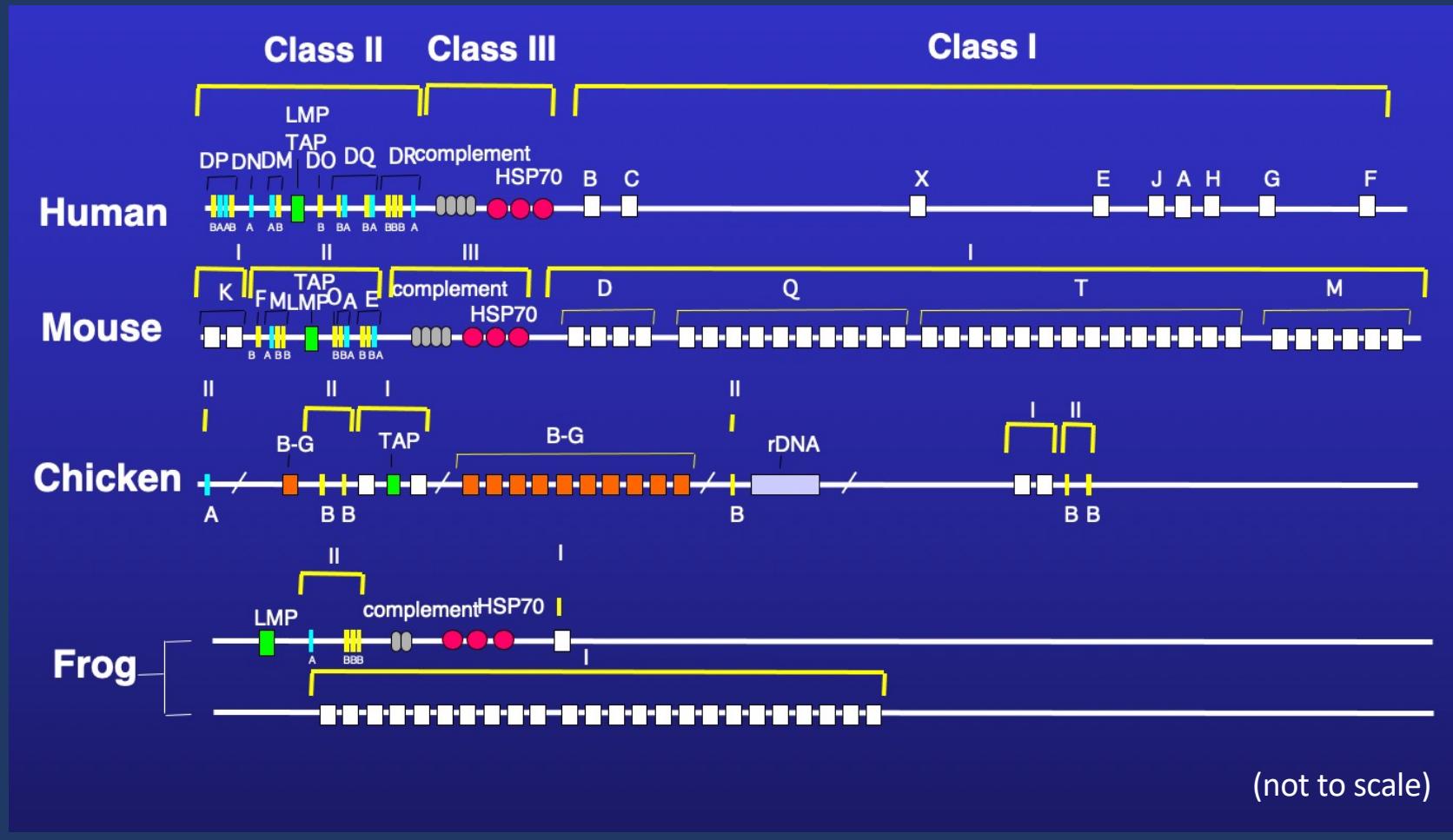
# Smaller regions of complexity in hox1a region



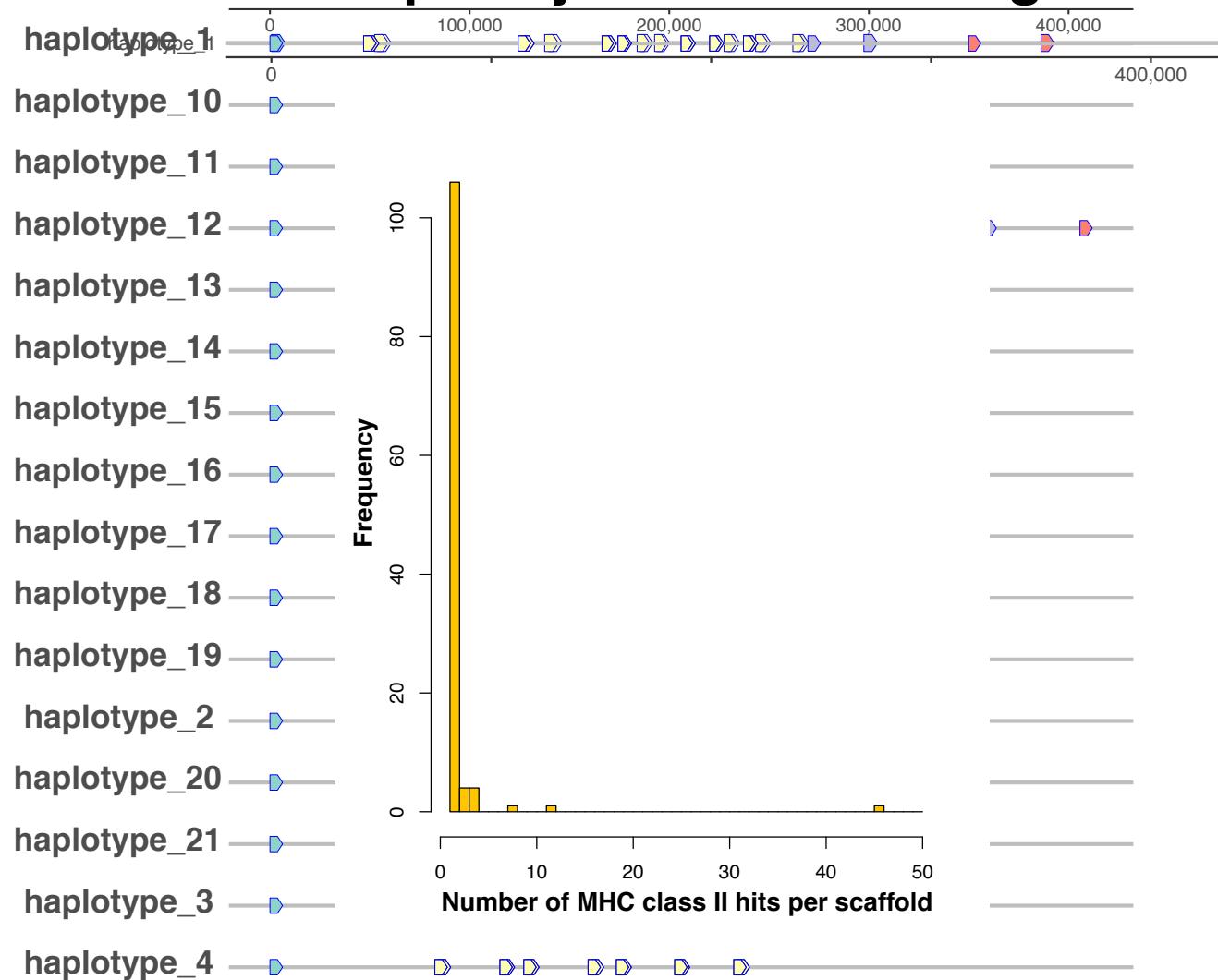
# Major histocompatibility complex



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



# Unprecedented complexity of MHC class II genes in scrub jays



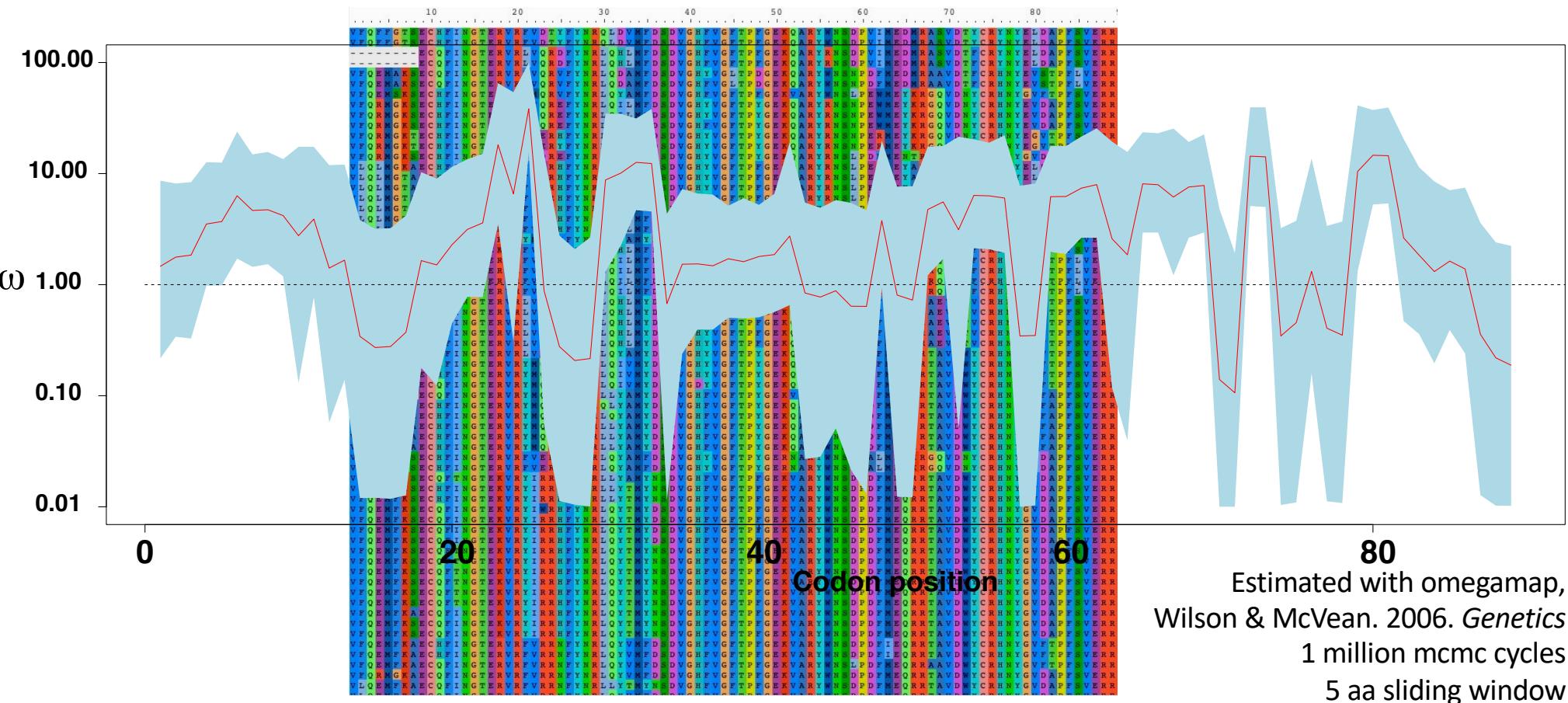
gene

- BRD2
- MHC\_classIIB\_exon3
- MHC\_classIIB\_exon2
- SLC39A7
- RXRB

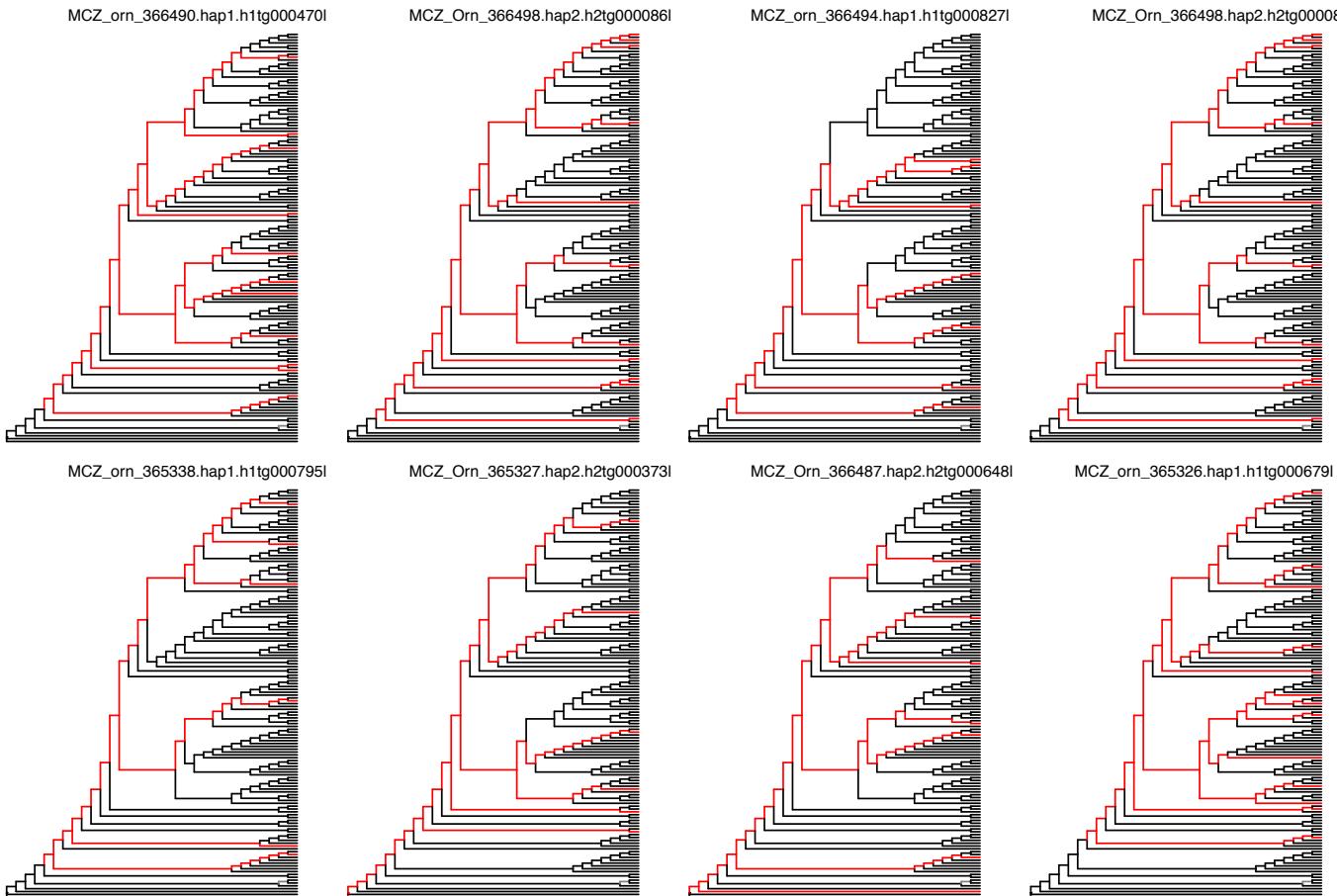
gene

- BRD2
- MHC\_classIIB\_exon3
- MHC\_classIIB\_exon2
- SLC39A7
- RXRB

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



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

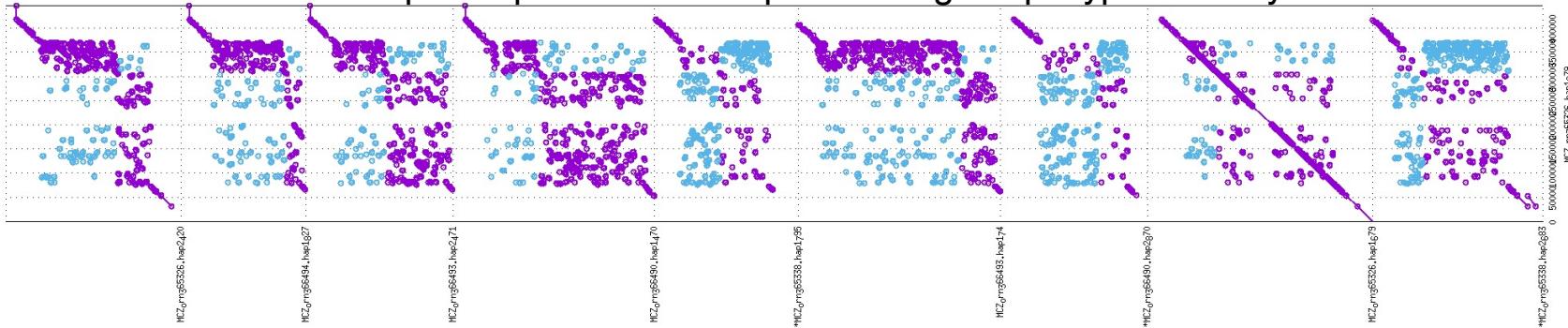


Phylogenetic paths of  
Mhc exon2 alleles  
on individual haplotypes

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



Complex repeat structures produce high haplotype diversity

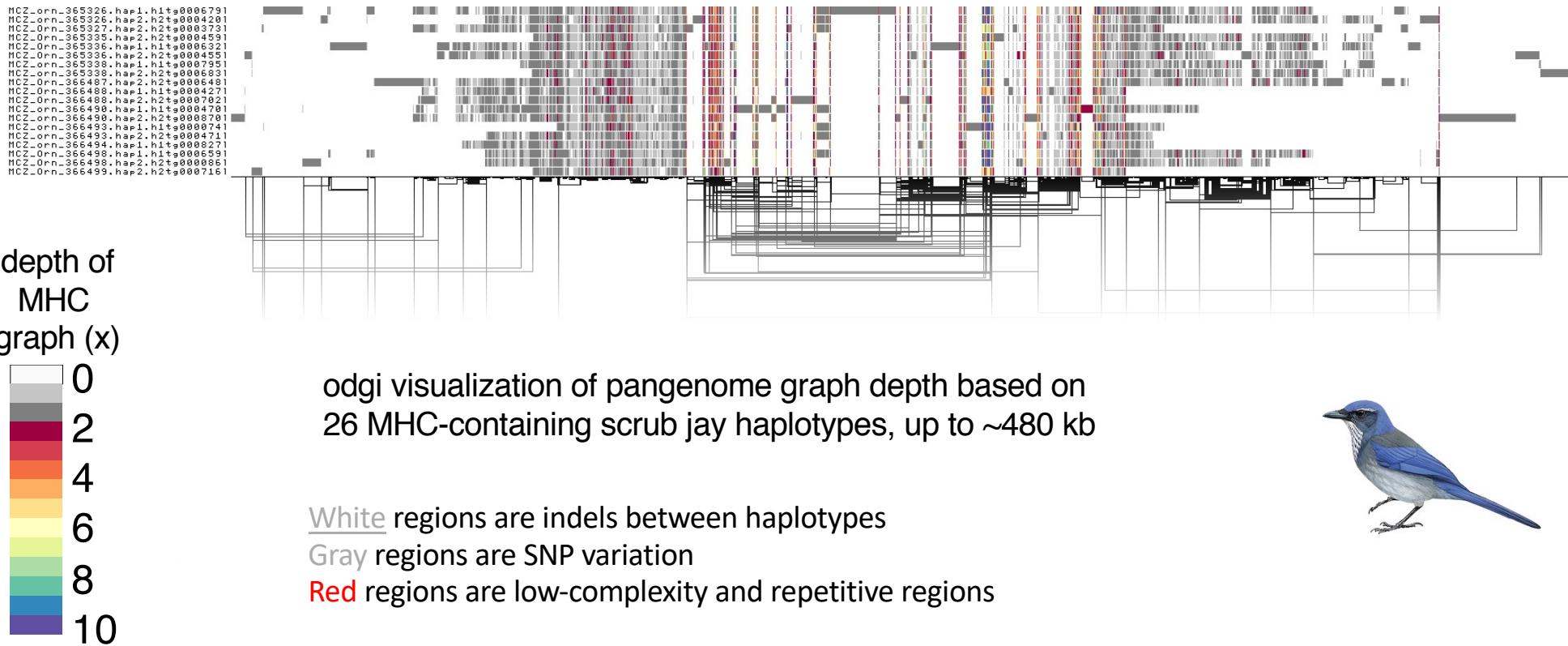


made with odgi and pangenome graph builder pipeline  
Guarracino et al. 2021. *Bioinformatics*, in press.

## Example satellites in MHC class II region of Woodhouse's scrub jay

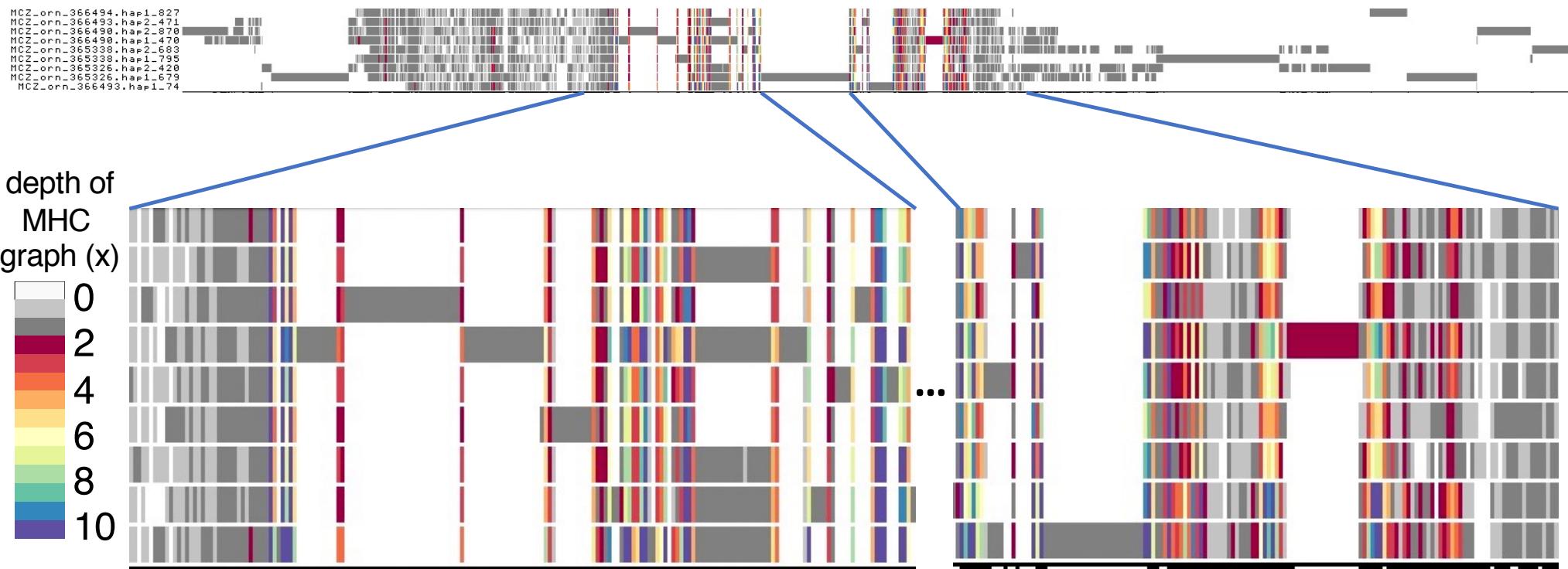


# Pangenome graph depth shows single-copy regions surrounded by complex VNTRs

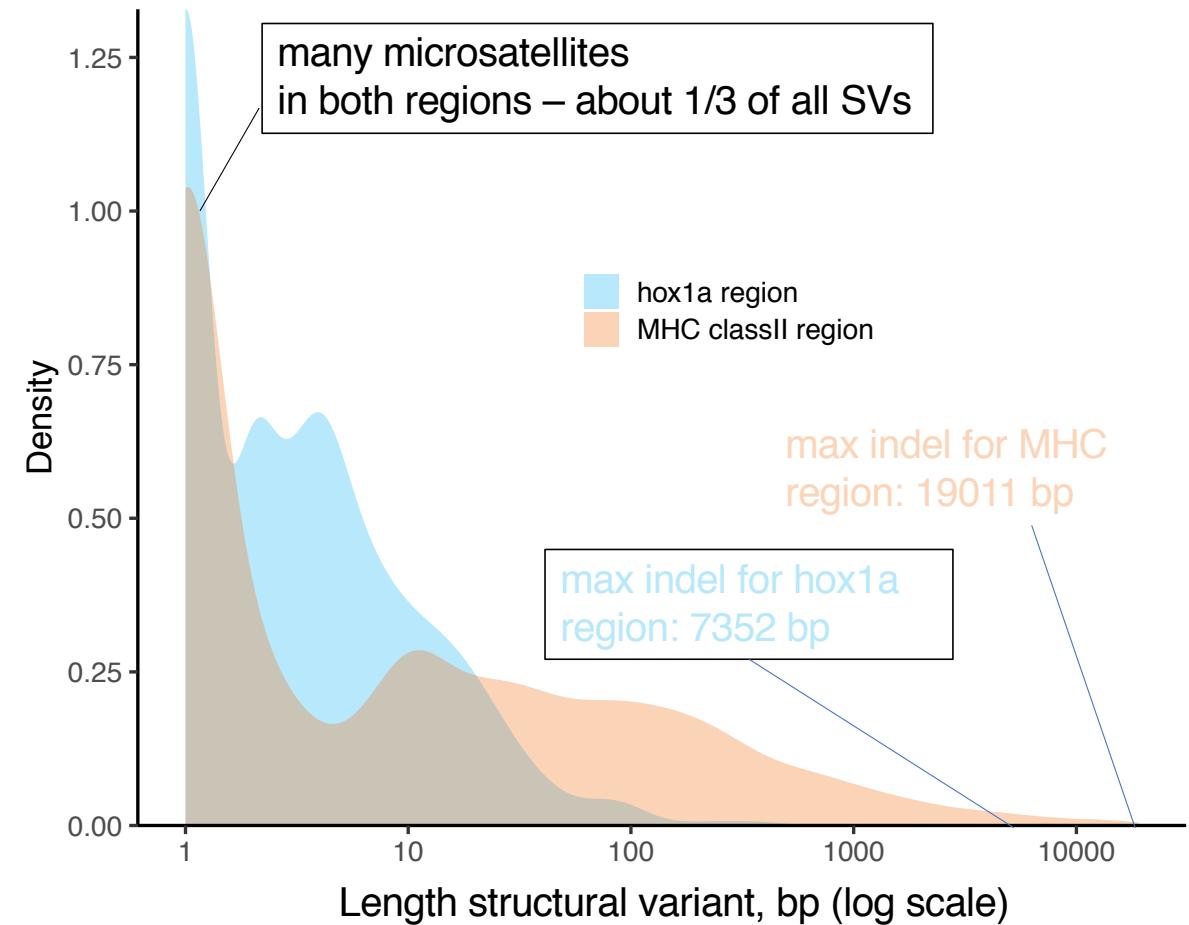
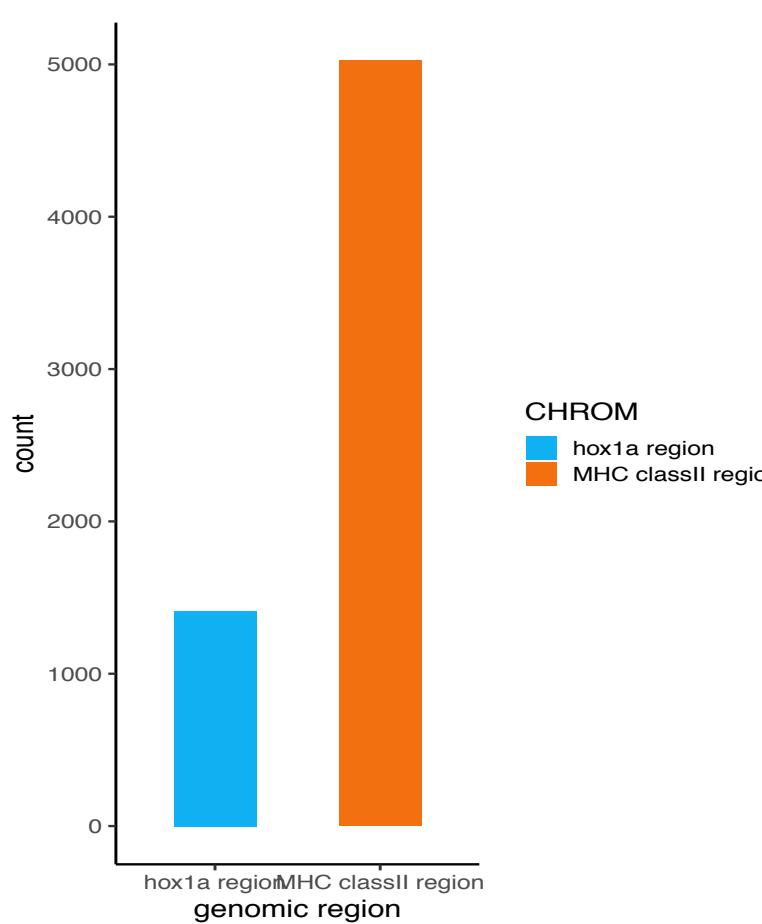


made with odgi and pangenome graph builder pipeline  
Guarracino et al. 2021. *Bioinformatics*, in press.

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

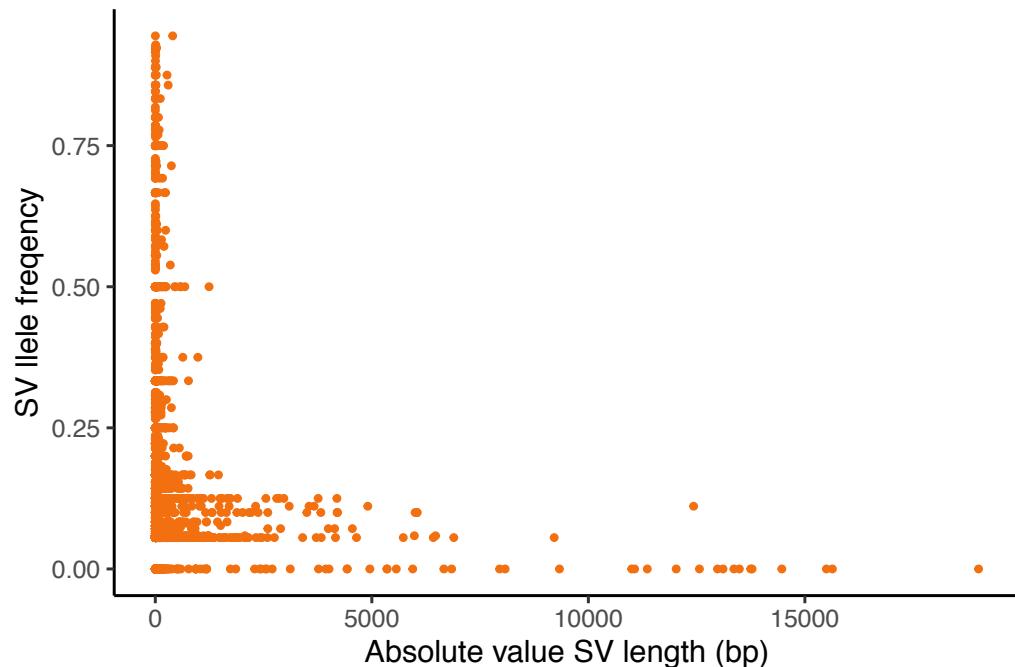


# MHC region has more numerous and longer structural variants than hox1a region

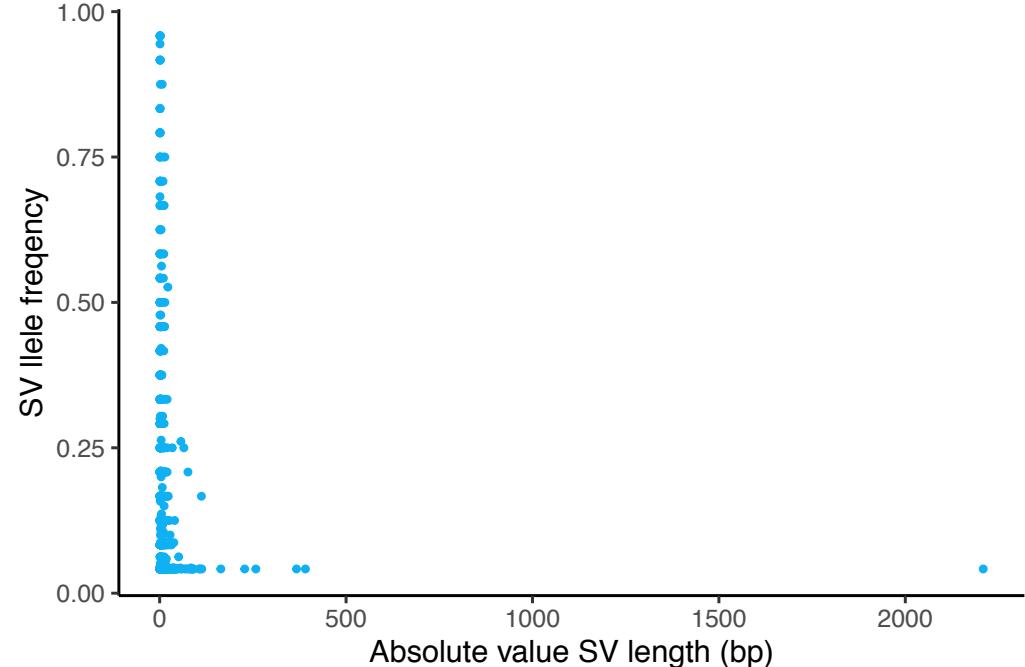


# Low frequency of large SVs in both MHC and hox1a regions (~400-kb)

Structural variants in MHC class II region



Structural variants in hox1a region

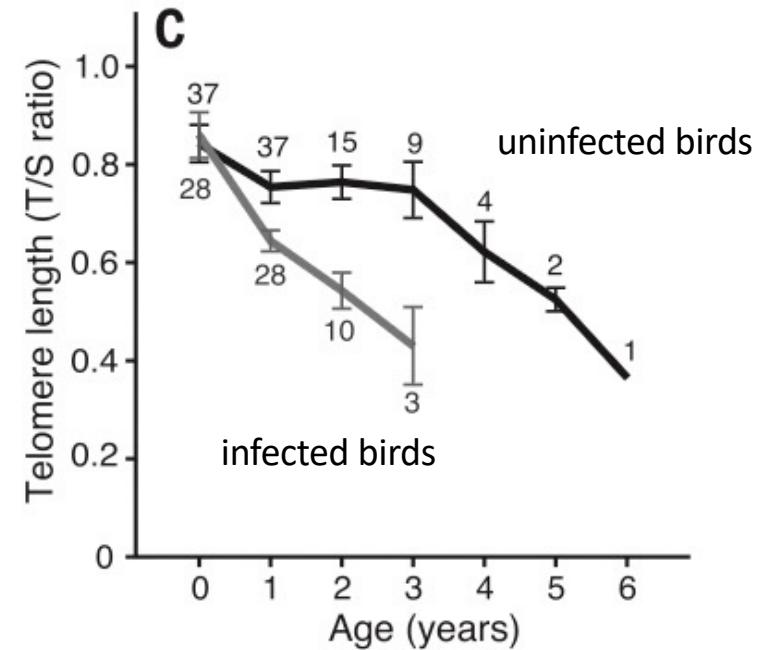
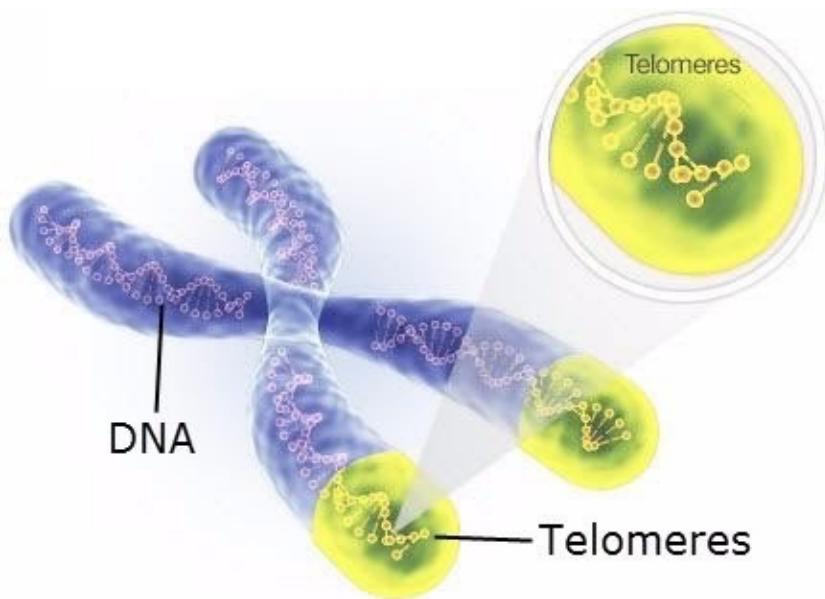


# Telomeres – barometers of age and stress in birds

RESEARCH | REPORTS

CHRONIC INFECTION

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

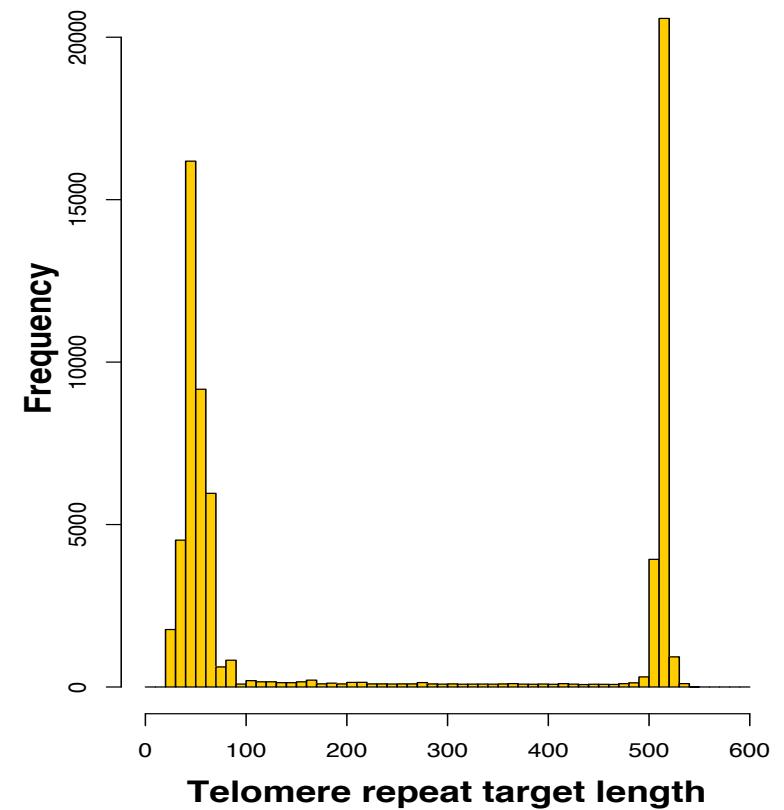


<https://medibalans.com/telomere/>

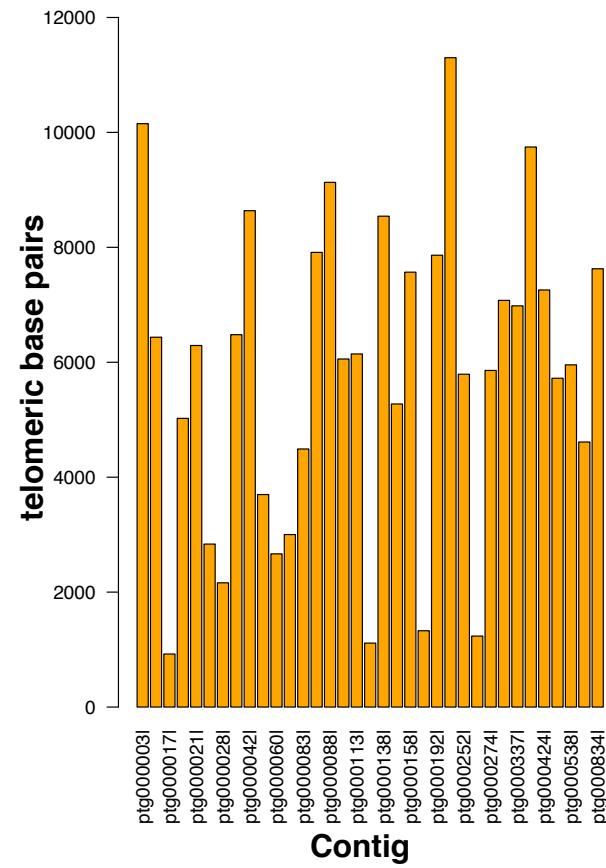
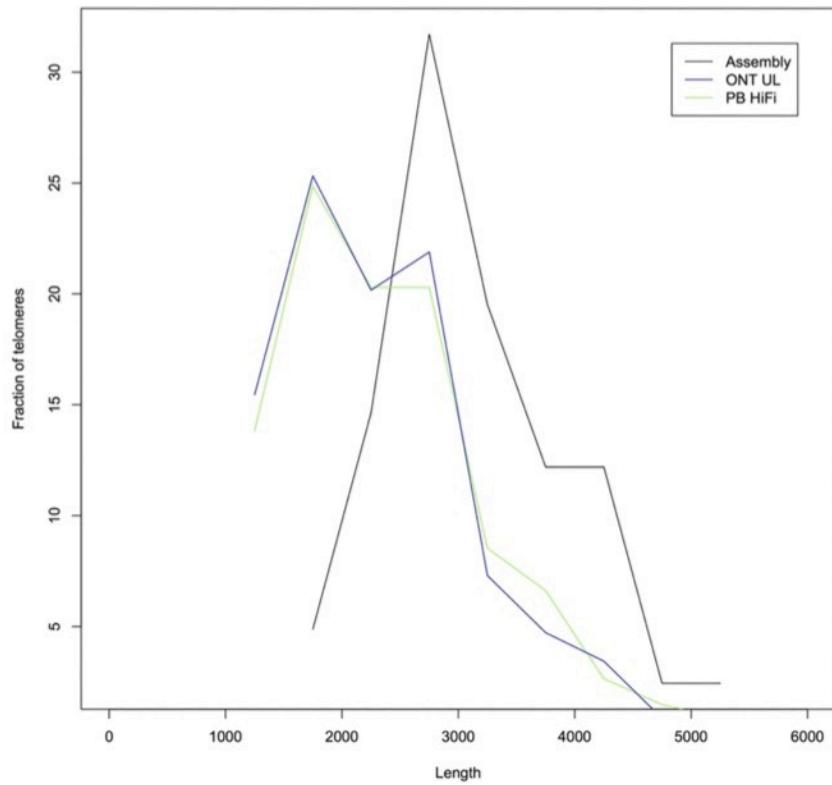
Ashgar et al. 2015. Science 347:436-438



# Blasting scrub jay contigs with telomere probe



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



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

## Telomeric sequences in scrub jays often occur near ends of contigs



# 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

## Pangenome informatics

Erik Garrison  
Andrea Guaracino



**Fieldwork**  
Greg and Donna

# 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

