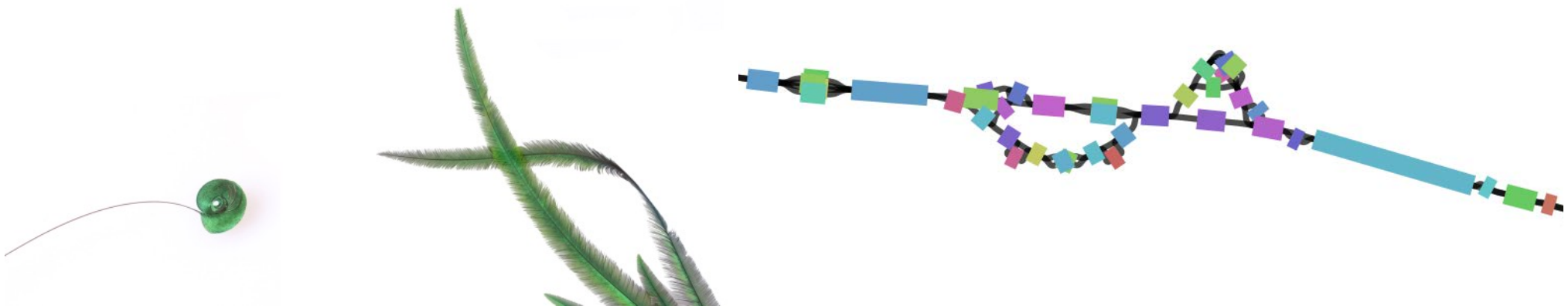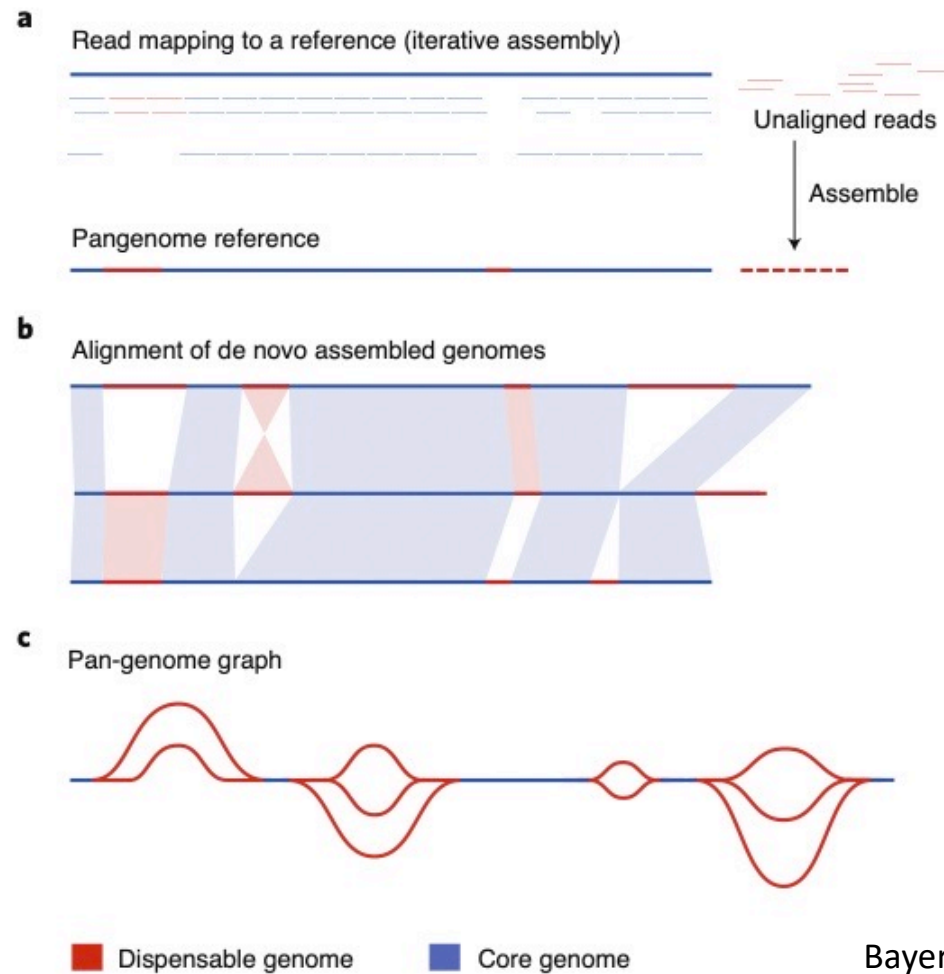# Pangenomes as a new tool for studying ecology and evolution of natural populations

Scott V. Edwards

Museum of Comparative Zoology, Harvard University, Cambridge, USA

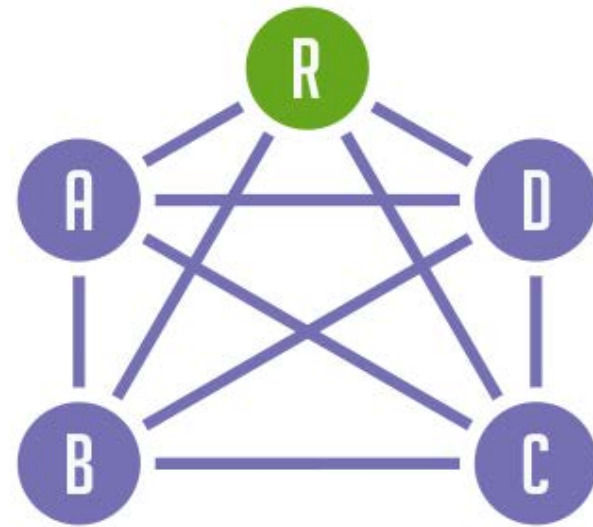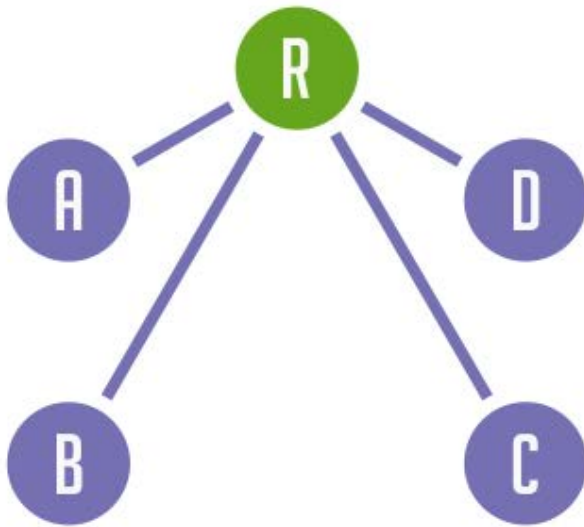# Pangenomes: moving beyond reference-based genomics



**a** Read mapping to a reference (iterative assembly)

Unaligned reads

Assemble

Pangenome reference

**b** Alignment of de novo assembled genomes

**c** Pan-genome graph

■ Dispensable genome   ■ Core genome

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

# Reference-free genomics

Genomic

Pangenomic

Reference model



Eizenga et al. 2021. *Ann. Rev. Genomics and Human Genetics*
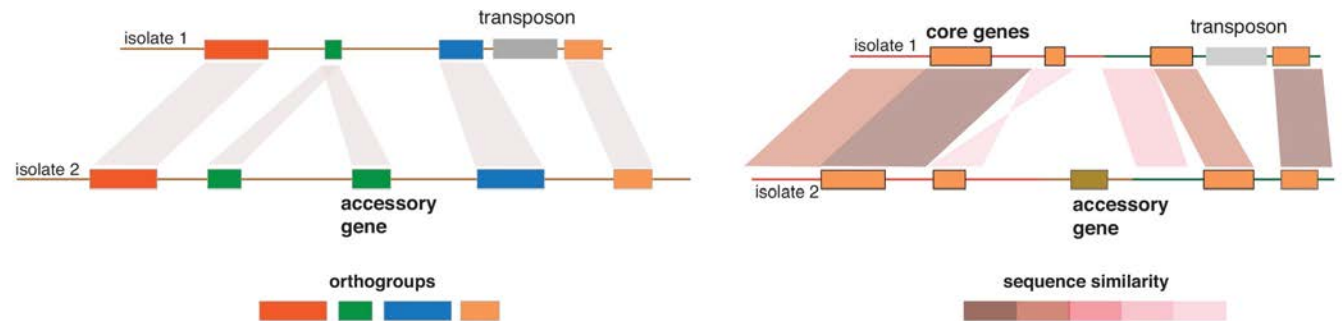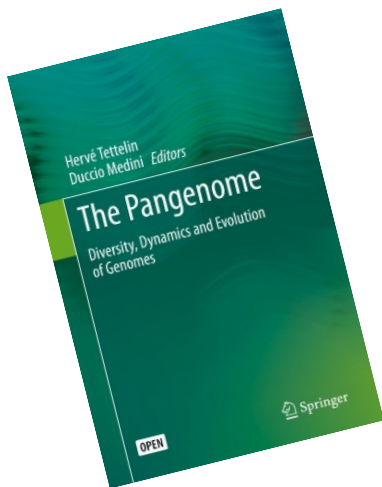
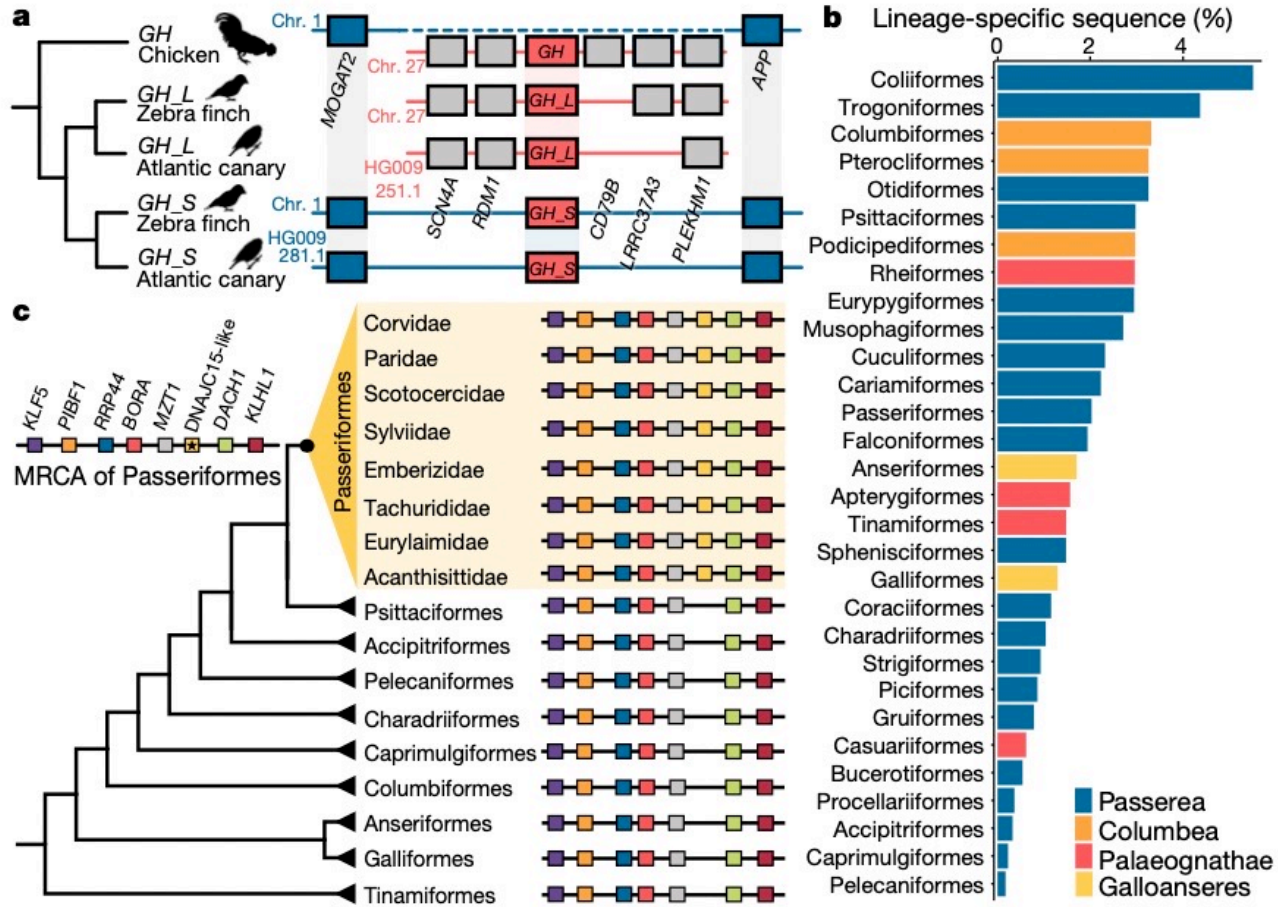# Open and closed pangenomes



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

# 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/

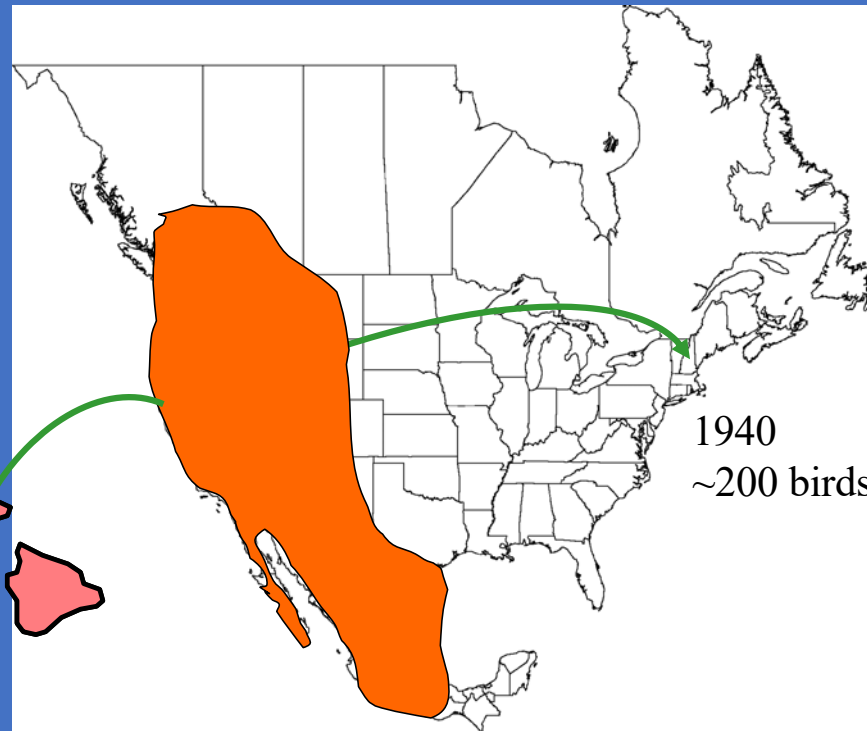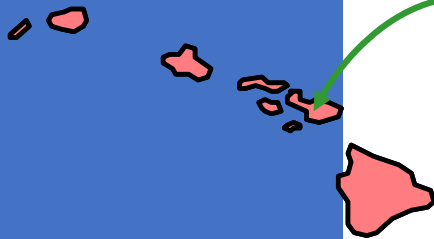# Pangenome approach to comparative genomics



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

# Recent history of House Finch populations

historic range
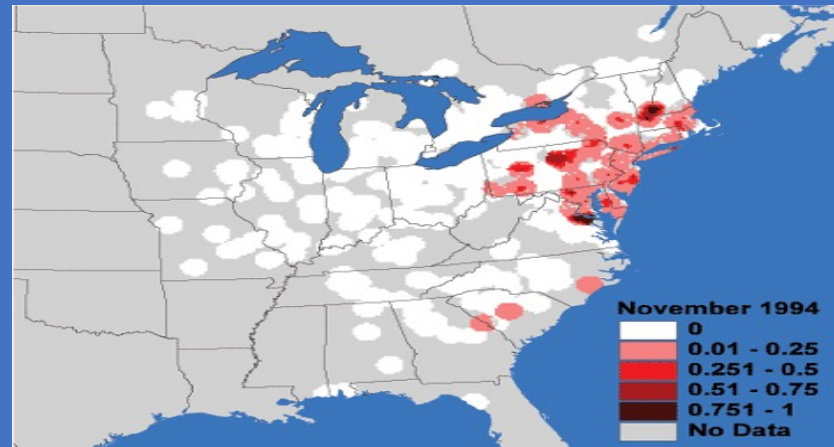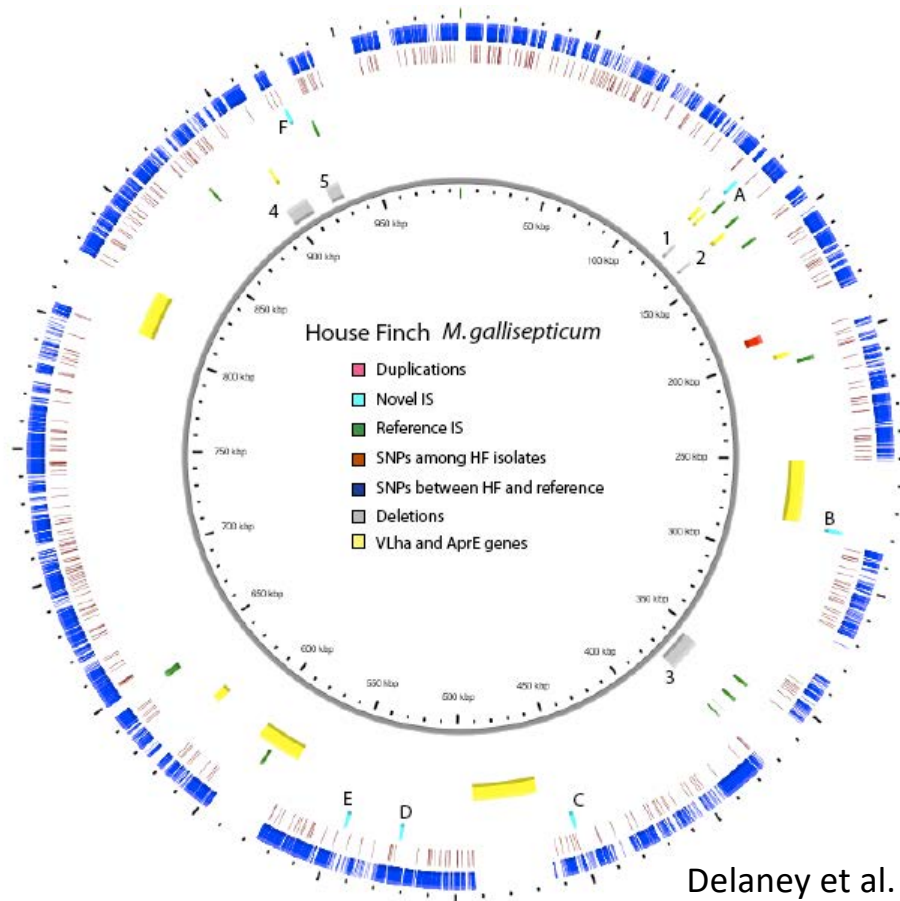
~1870 bottleneck?

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.
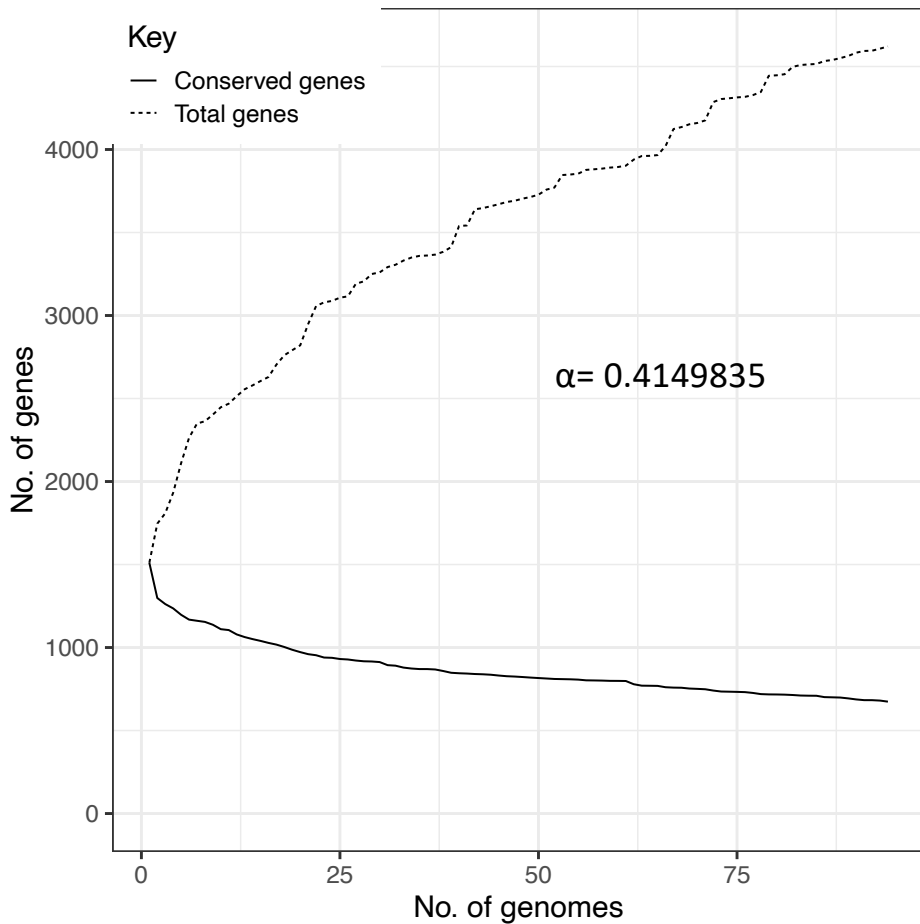
# House Finch *Mycoplasma* genome ~1 Mb



House Finch *M.gallisepticum*

- Duplications
- Novel IS
- Reference IS
- SNPs among HF isolates
- SNPs between HF and reference
- Deletions
- VLha and AprE genes

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

Added 12 new House Finch Mycoplasma strains, sequenced with PacBio

Used

Delaney et al. 2012. *PLoS Genetics*

# Pangenome of *Mycoplasma gallisepticum*



**Key**
— Conserved genes
···· Total genes

α= 0.4149835

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

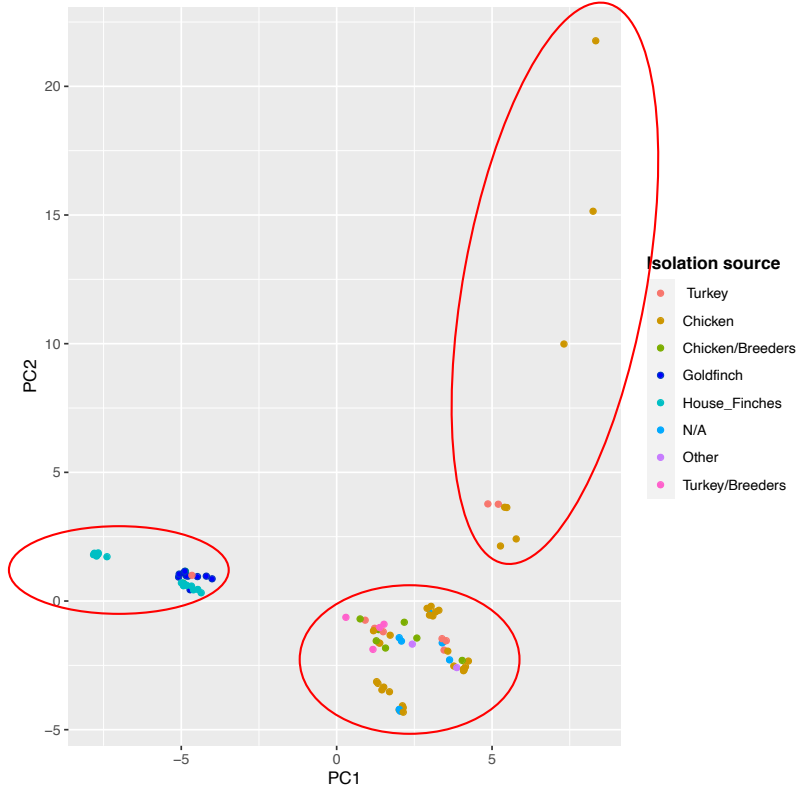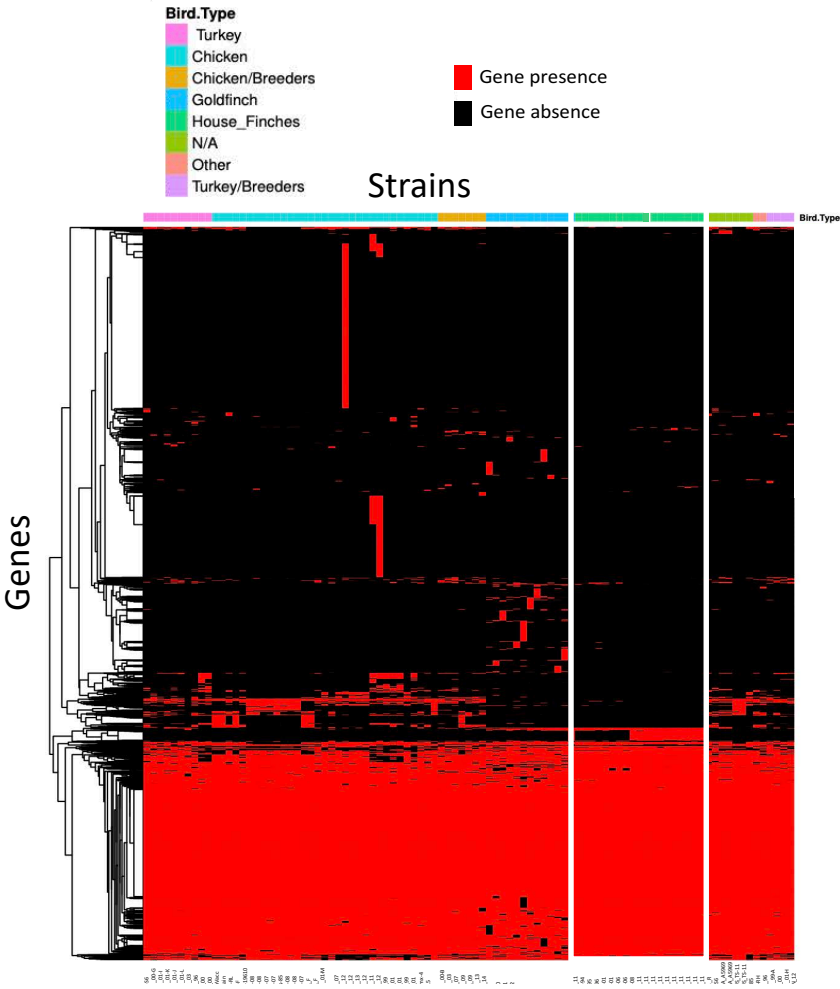| Feature | Info | Number of genes | Percentage |
|---------|------|-----------------|------------|
| Core genes | (99% <= strains <= 100%) | 674 | 14.586 |
| Soft core genes | (95% <= strains < 99%) | 464 | 10.041 |
| Shell genes | (15% <= strains < 95%) | 412 | 8.916 |
| Cloud genes | (0% <= strains < 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
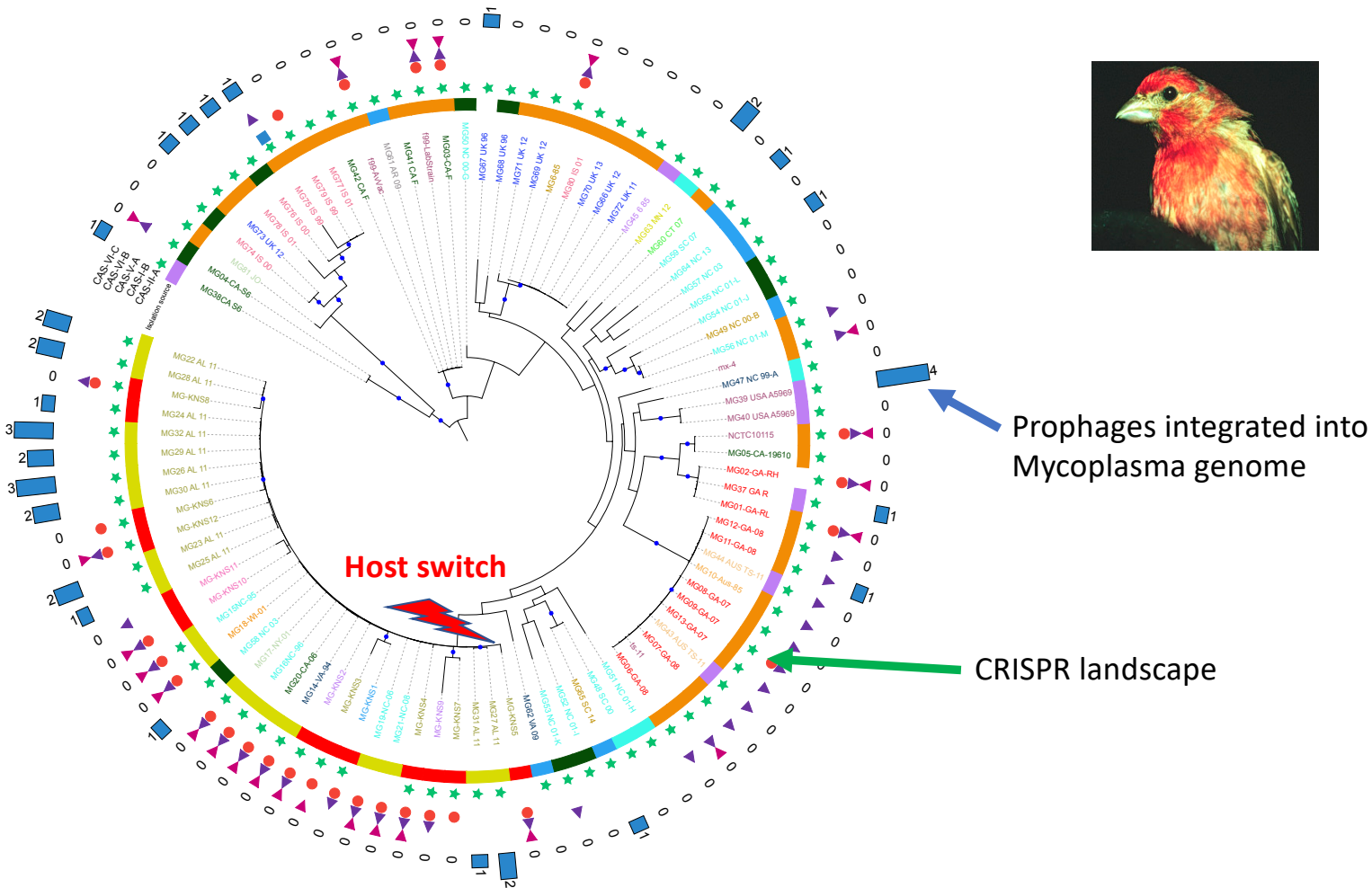if we collected *all* genomes of the species

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

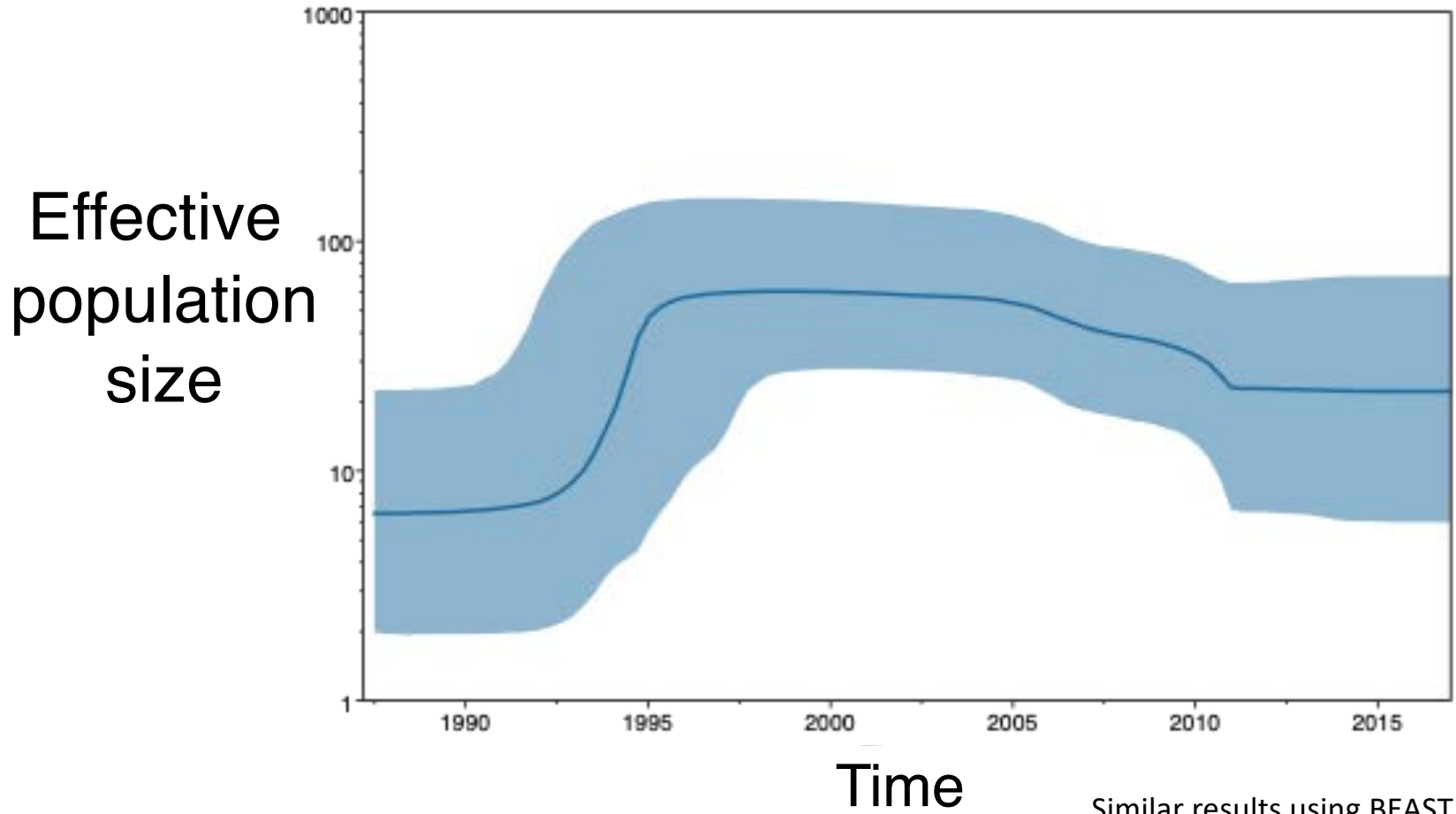*the pan-genome is closed if the estimated alpha is above 1.0

# Mycoplasma pangenome gene repertoire is highly strain-specific

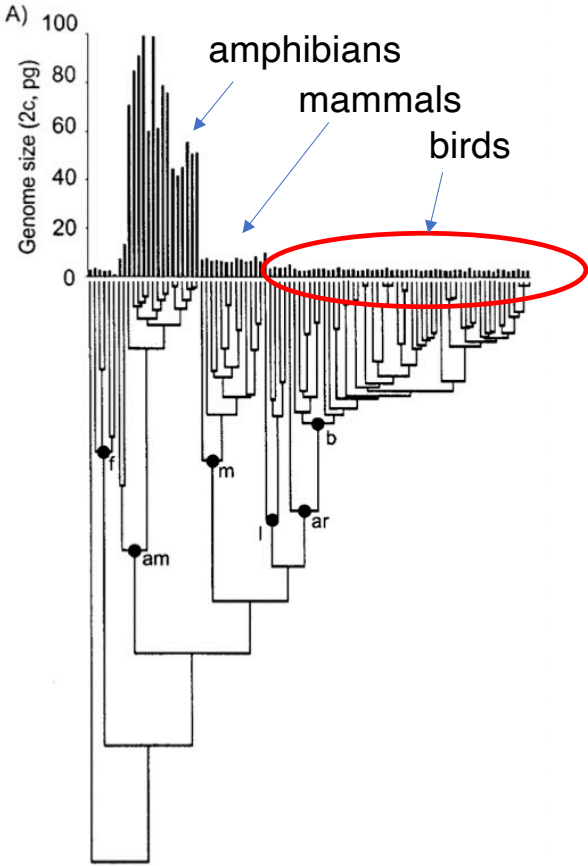# House Finch *Mycoplasma* strains have distinct CRISPR and prophage landscapes



Prophages integrated into Mycoplasma genome

CRISPR landscape

Host switch

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



Effective population size (y-axis, log scale from 1 to 1000)
Time (x-axis, 1990 to 2015)

Similar results using BEAST and Stairway plot

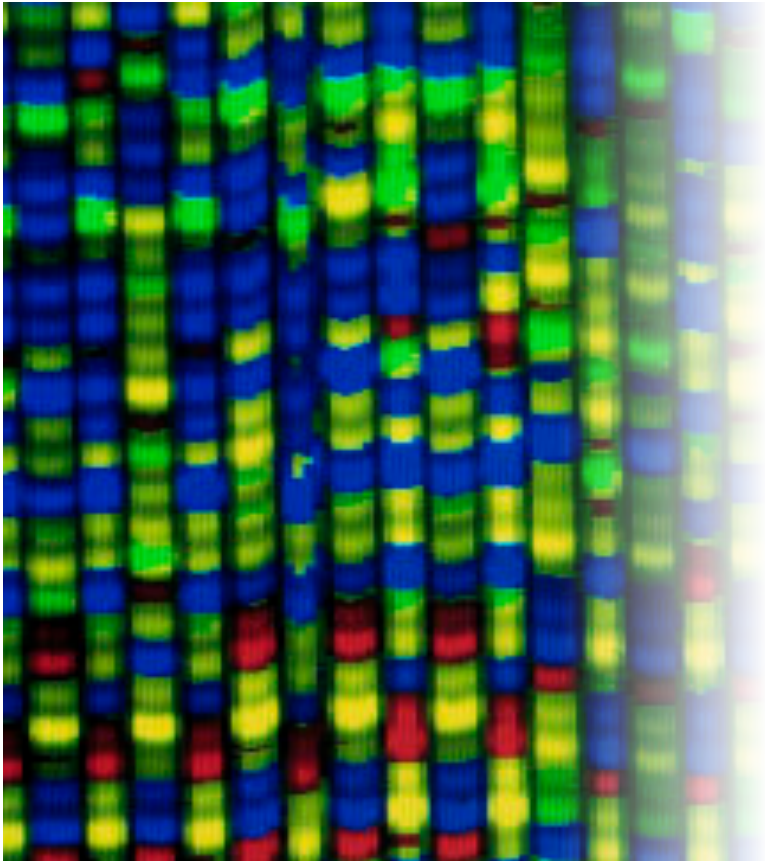# Birds have small, streamlined genomes
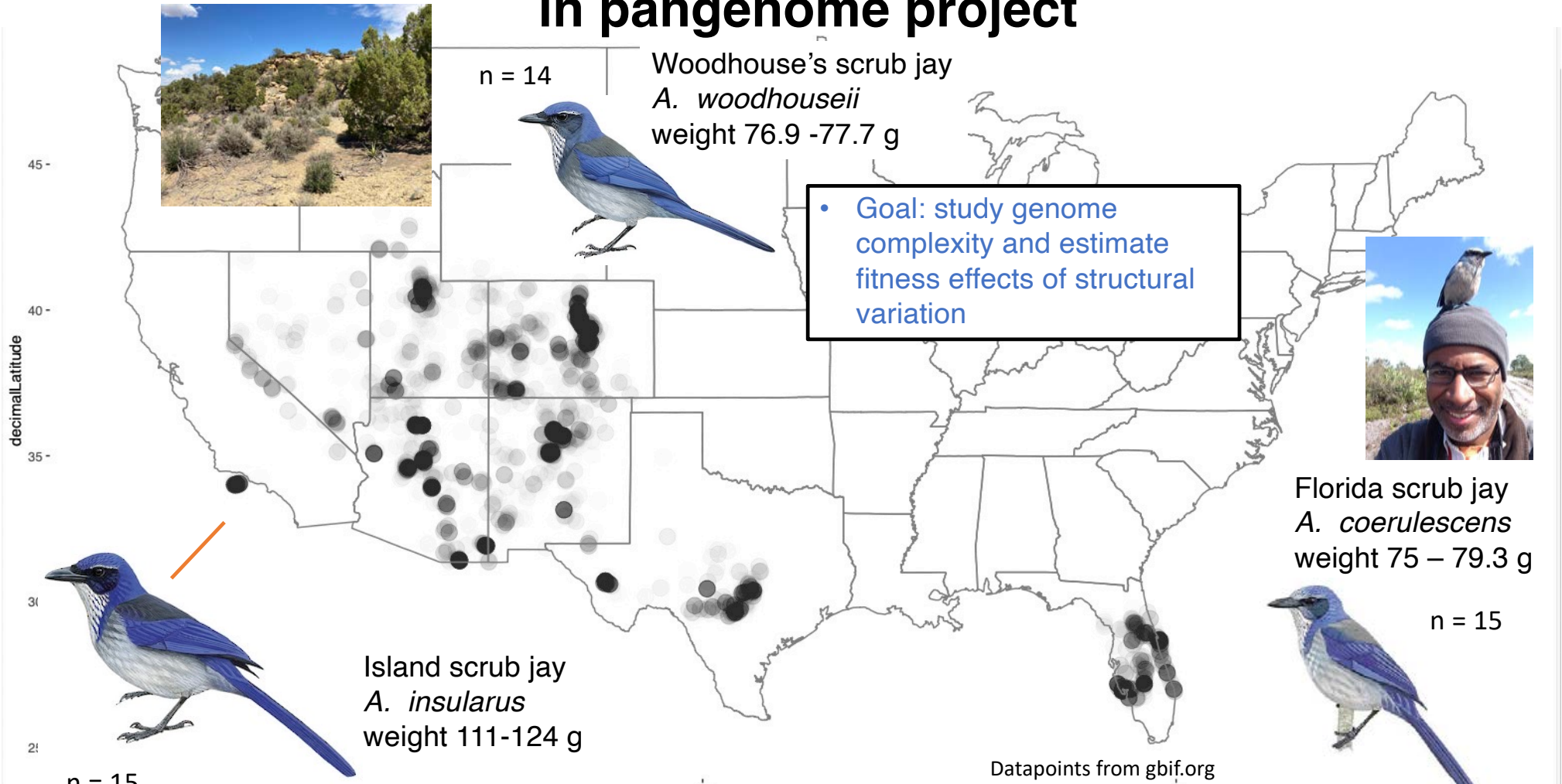


Waltari & Edwards. 2002. *Am. Nat.*



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. insularus*
weight 111-124 g

n = 15

Datapoints from gbif.org

decimalLatitude

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

Scott V. Edwards [iD][1,2,*], V. V. Robin[3], Nuno Ferrand[4], and Craig Moritz[5]

GBE

## Table 1

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

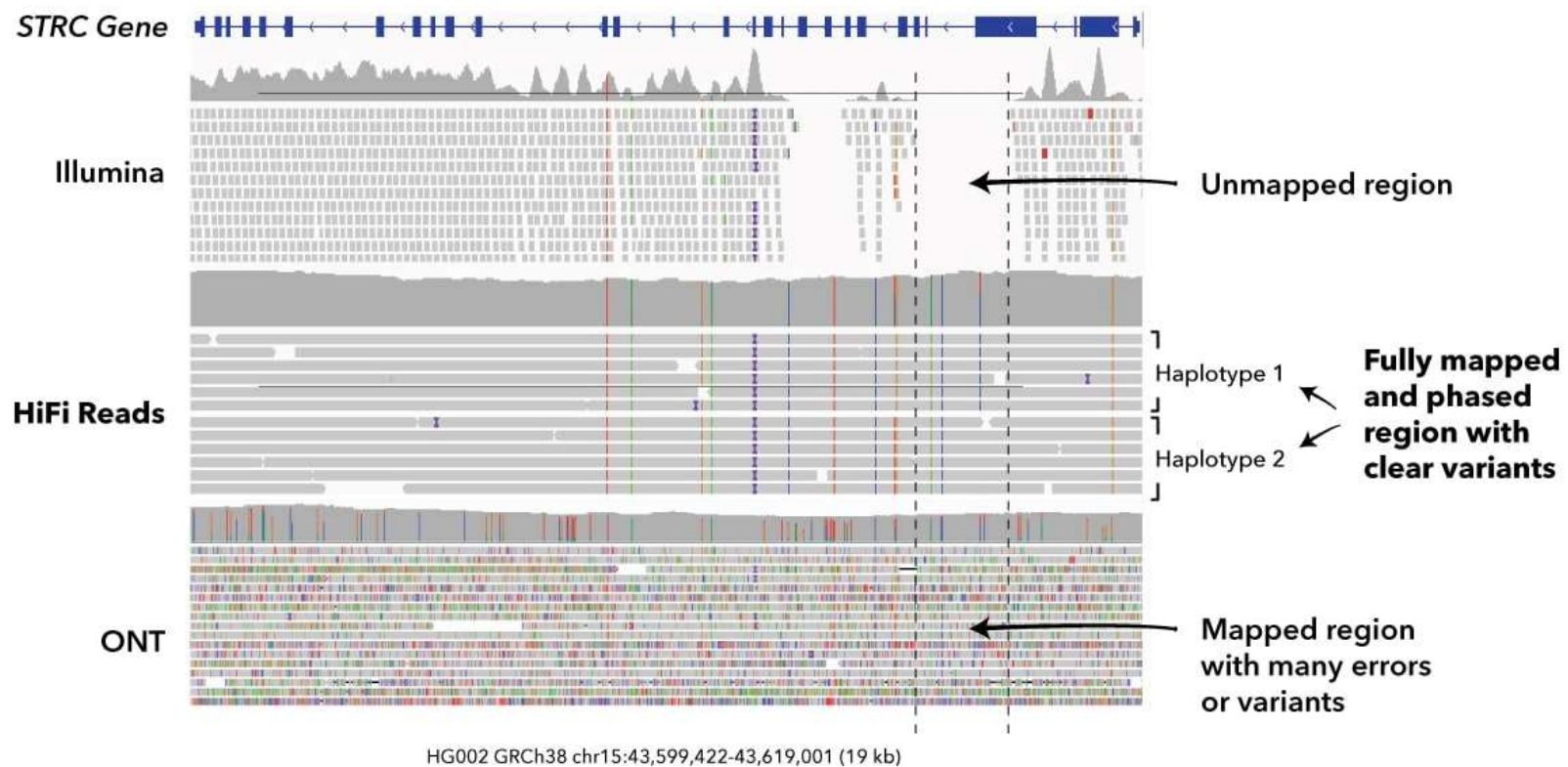| Concept/Parameter | Comparative Population Genomics | Landscape Genomics | Comparative Phylogeography |
|---|---|---|---|
| Comparative perspective | Growing | Nascent | Mature |
| Emphasis on space | No | Yes | Yes |
| Geographic scale | Random mating population | Region | Biome |
| Temporal scale | 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 <u>and</u> accurate

▸ HiFi reads: long & accurate

▸ A breakthrough every ~5 years

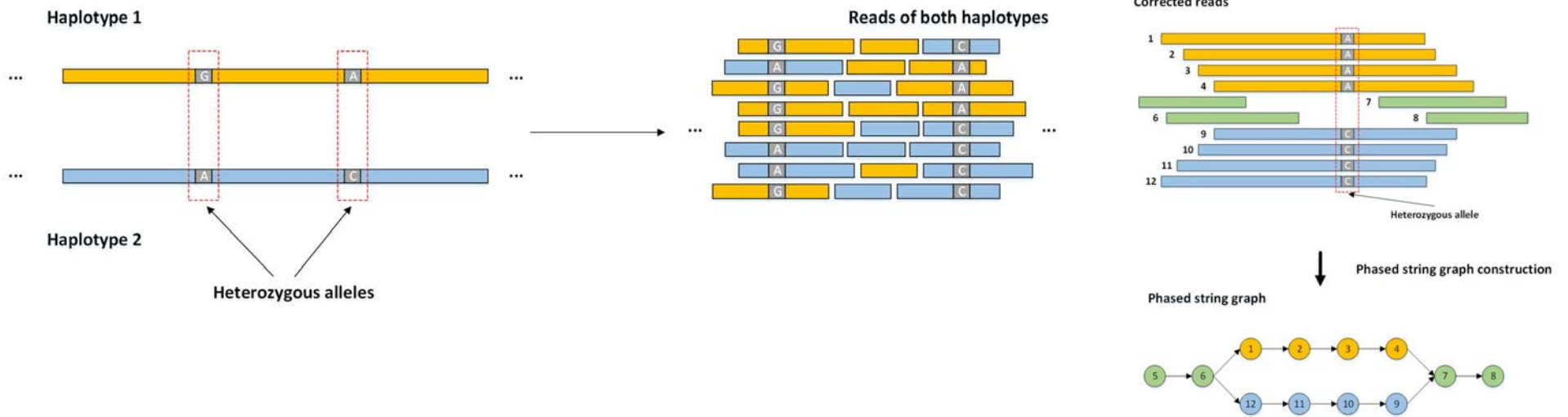▸ Most existing assemblers cannot make full use of the accuracy



Coutesy Haoyu Cheng, Dana Farber Cancer Institute

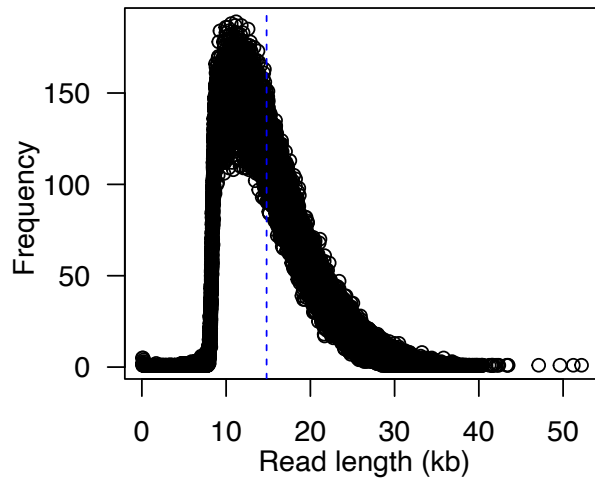# PacBio HiFi reads are long <u>and</u> accurate

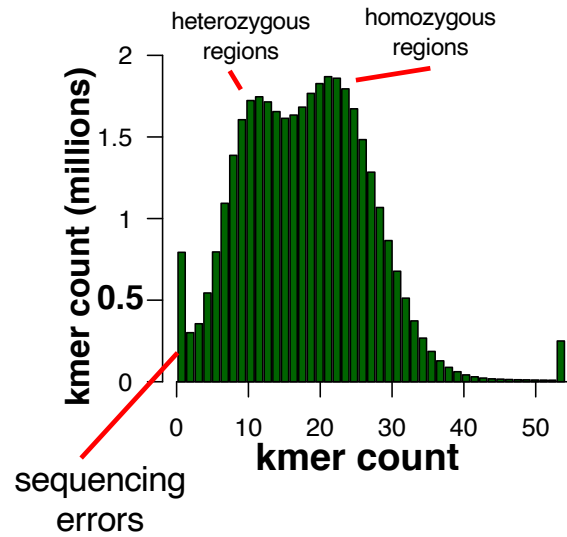# Hifiasm – a HiFi accurate read assembler that resolves haplotypes



Coutesy Haoyu Cheng, Dana Farber Cancer Institute
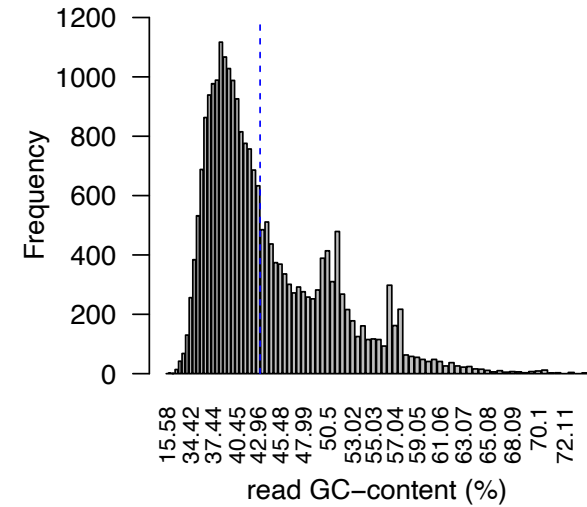
# Scrub-jay PacBio HiFi data characteristics

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



| assembly | Number of contigs | contig N50 (Mb) |
|---|---|---|
| primary | 967 | 14.75 |
| haplotype 1 | 1482 | 7.84 |
| haplotype 2 | 1218 | 7.90 |

*A. woodhouseii*

scaffold length (megabases)

# Quality of assemblies varies by species/tissue type



N50 by sample

# 60-fold range in effective population size across species



**PSMC**
SNP data from genome-wide dipcall regions

**bpp**
~2500 1-kb dipcall regions

Li & Durbin. 2011. *Nature* 2011, 475:493-496; Flouri et al. 2020. *Mol. Biol. Evol.* 37: 1211-1223

RepeatMasker analysis suggests over 25% repeats and transposable elements

Assemblies of Island Scrub Jays are ~100 Mb smaller than Woodhouse's Scrub Jay

N = 30, 30 and 28 haplotypes (AI, AW, AC)

Florida
Island
Woodhouse's

mean within species

total assembly length (Gb)

sum of contigs of primary assembly from hifiasm
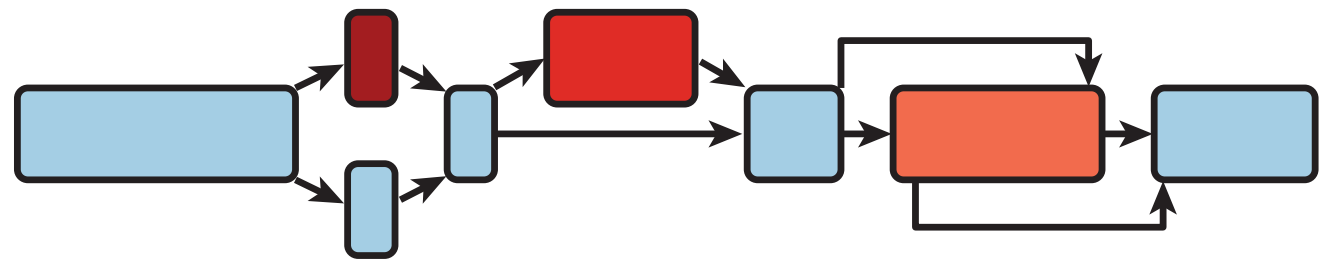
# Estimates of max genome size from Genomescope using k-mers

# Pangenome graphs capture structural variation within species



minigraph

odgi

build genome graph

odgi viz

bandage

visualize genome graph

```
>h1tg000104l
GGCGGGGCCCGGAGGGGCCGGGGCCGCTGAGGGGCCGCGGGTGCGGCAGAGCG
>h1tg000528l
ATGGATACTTTCCAGTCAGAGCTTTATAATAATTTCCATAATTTAAATATTTT
>h1tg000795l
ACTTTGGGGACACCTTTGGGGACACCTCGGGGGACACTTTGGGCCACAAATCC
```

unaligned fasta files

Multiple sequence alignment

Bidirected genome graph

Eizenga et al. 2021. *Ann. Rev. Genomics Hum. Genetics*

# 2D pangenome graph visualizations – PGGB/Odgi

Chr 18 – 12 Mb

Chr 1 – 160 Mb
'telomere kiss'

Chr 23 – 8 Mb

# Variation in depth of a pangenome graph



depth of
MHC
graph (x)

0
2
4
6
8
10

medium
(self-)node-depth regions
depth = ~ n haplotypes

high node-depth regions –
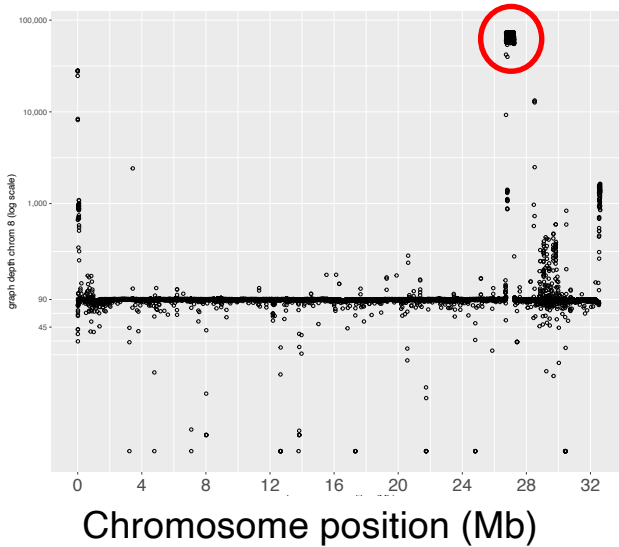large structural variants, satellites

low node-depth
regions – SNPs
or small
indels

# Graph depth of microchromosome 27 correlates with LTRs and satellites

graph depth chromosome 27

chromosome position

bases in 100 kb window (bp)

satellite / repeat
- rnd-5_family-33_rnd-5_family-92_ltr-1_family-44
- rnd-6_family-14455
- sj_sat_2216_2221_2218
- sj_sat_circ144_rnd-1_family-59_ltr-1_family-129
- sj_sat_circ153-5084

# Further examples of graph depth scans

## Chr 8



## Chr 1A



## Chr 20

# Graph depth of W chromosome correlates with LTRs and satellites

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



Pangenome graphs
generated with odgi
and visualized with Bandage

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

7.5 kb polymorphic indel

2.5 kb polymorphic indel

# Smaller regions of complexity in hox1a region



depth of
*hox1a* region
graph (x)

| | |
|---|---|
| | 0 |
| | 2 |
| | 4 |
| | 6 |
| | 8 |
| | 10 |

White region
Gray regions
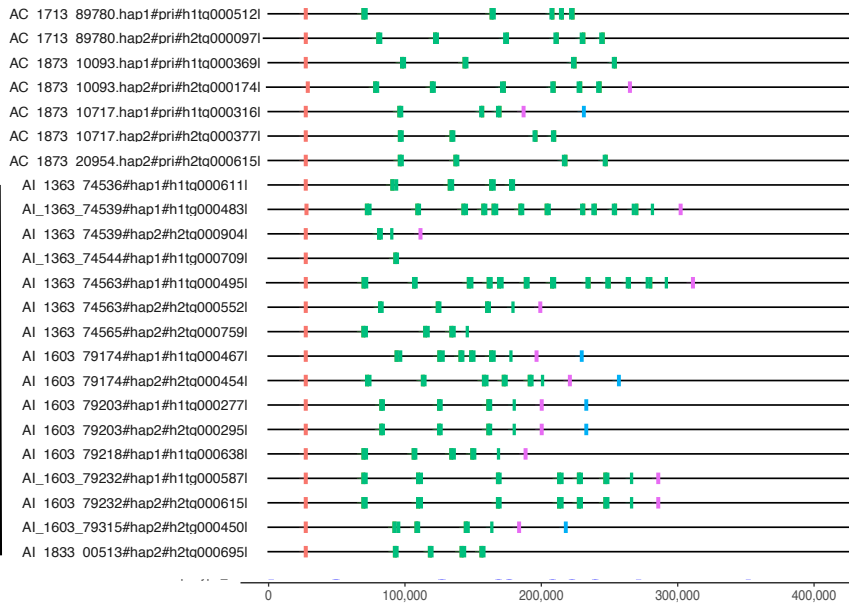Red regions a                                    gions

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



(not to scale)

# Extraordinary haplotype diversity in MHC class II region in Scrub Jays
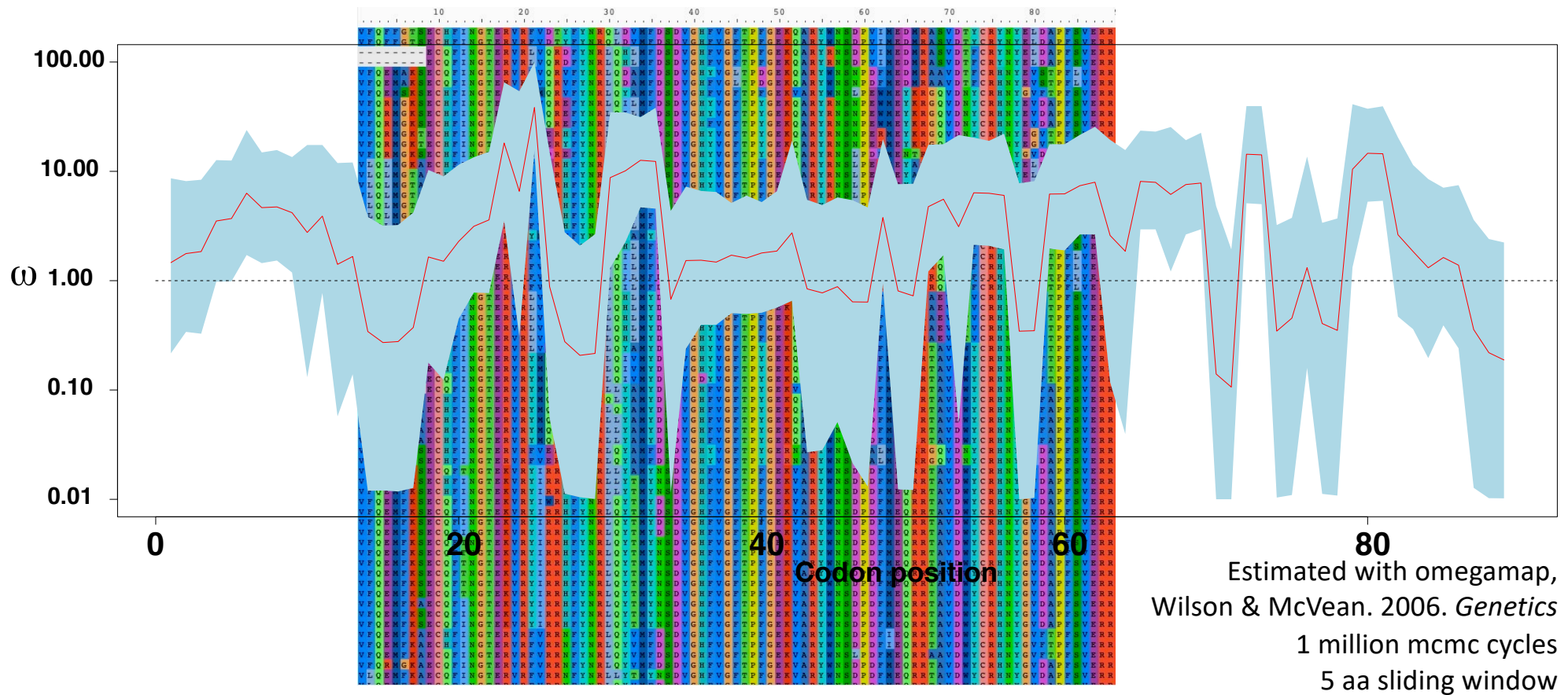
BRD2  MHC-class II  RXRB_exon4  SLC39A7

Florida

Island

Woodhouse's

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

Estimated with omegamap, Wilson & McVean. 2006. *Genetics*
1 million mcmc cycles
5 aa sliding window

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



MCZ_orn_366490.hap1.h1tg000470l

MCZ_Orn_366498.hap2.h2tg000086l

MCZ_orn_366494.hap1.h1tg000827l

MCZ_Orn_366498.hap2.h2tg000086l

MCZ_orn_365338.hap1.h1tg000795l

MCZ_Orn_365327.hap2.h2tg000373l

MCZ_Orn_366487.hap2.h2tg000648l
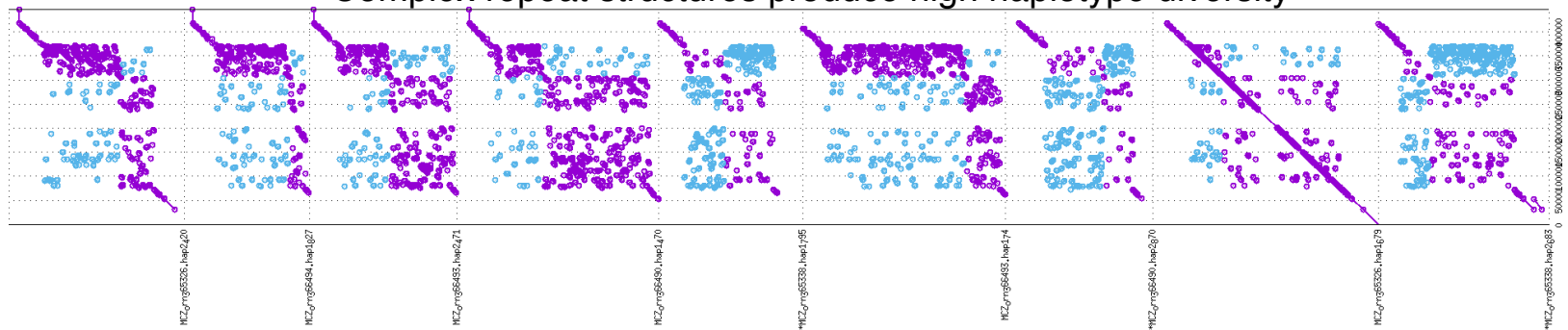
MCZ_orn_365326.hap1.h1tg000679l

Phylogenetic paths of
Mhc exon2 alleles
on individual haplotypes

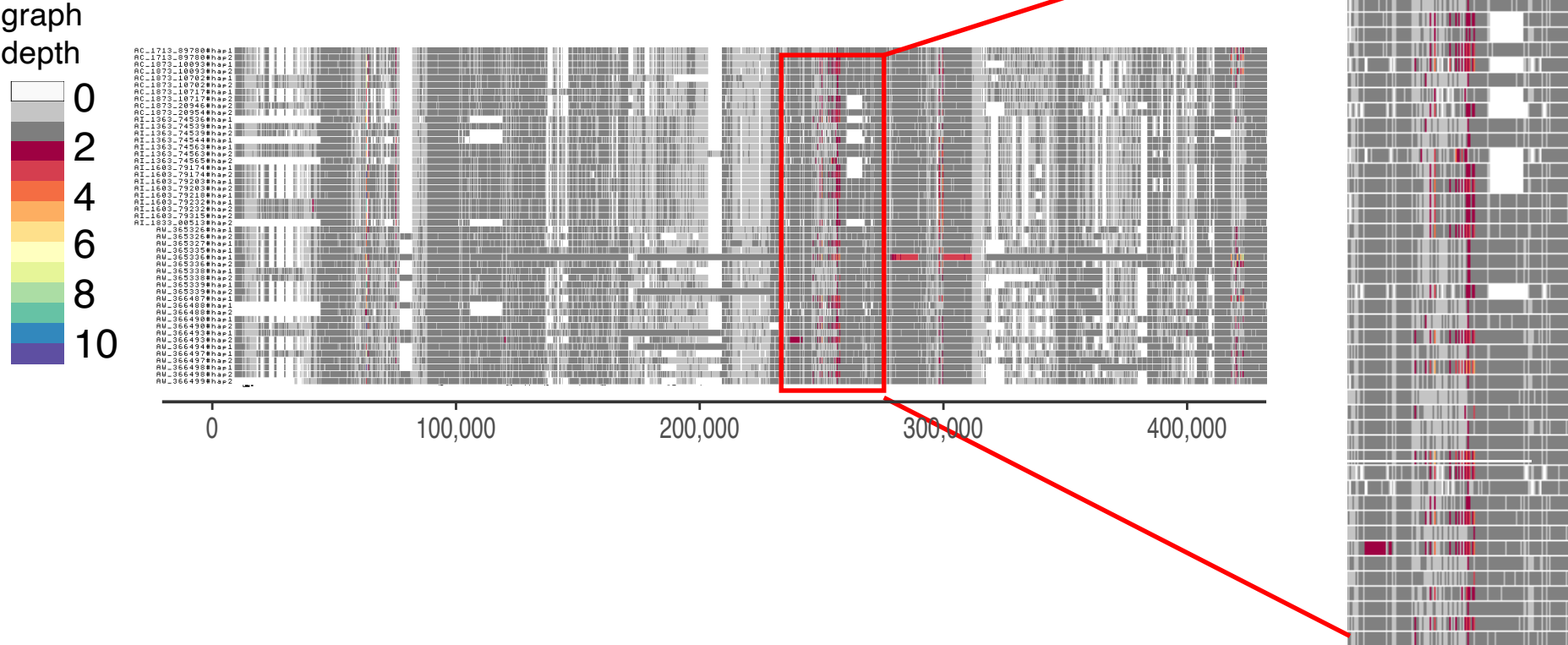# Visualization of MHC class II region in 22 haplotypes of Woodhouse's scrub-jays with odgi



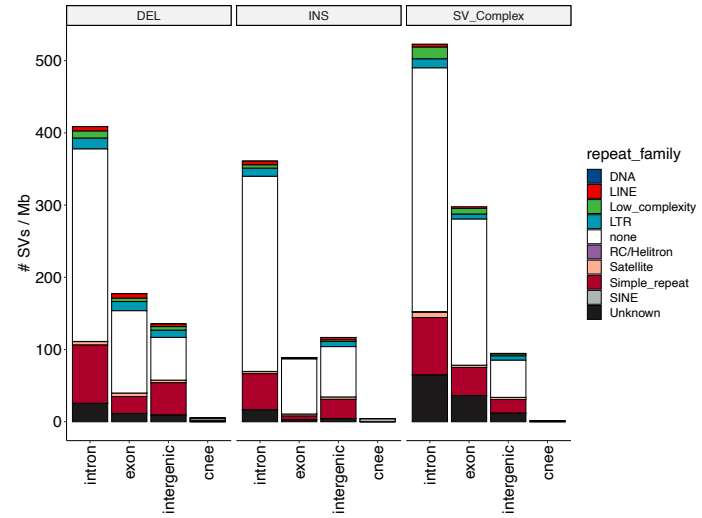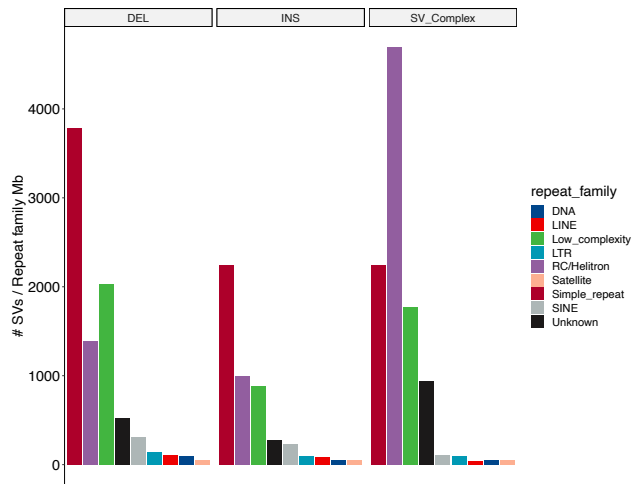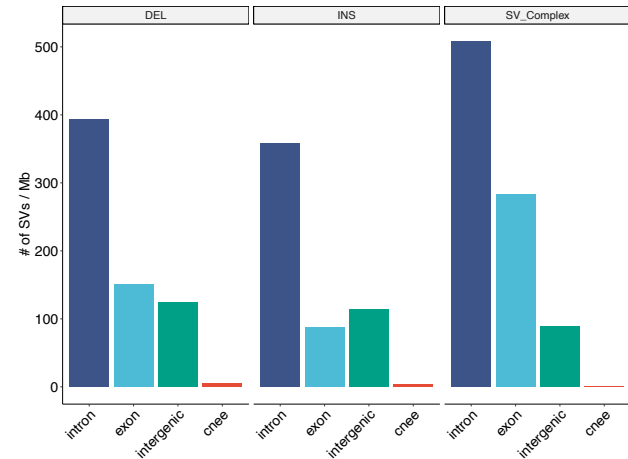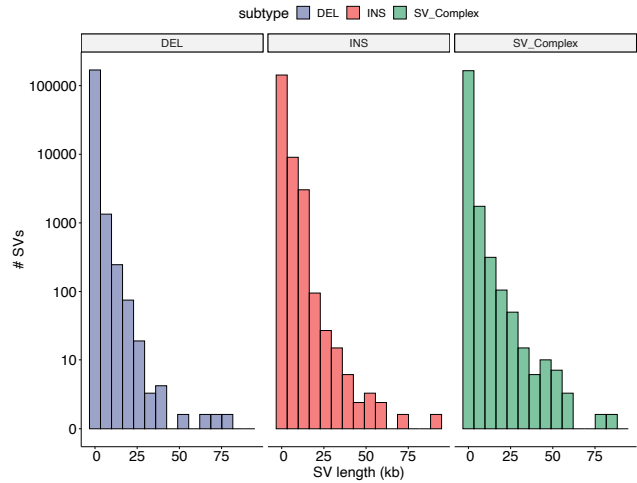Complex repeat structures produce high haplotype diversity



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

MHC graph depth shows single-copy MHC regions surrounded by complex VNTRs

# SVs from PGGB: Danielle's plots

# inversion lengths longer in heterozygotes

## Florida          Island          Woodhouse's



inversion length (bp – log scale)

327 inversions across data set; range 233 bp – 99.5 kb (hard limit); 16-53 per haplotype
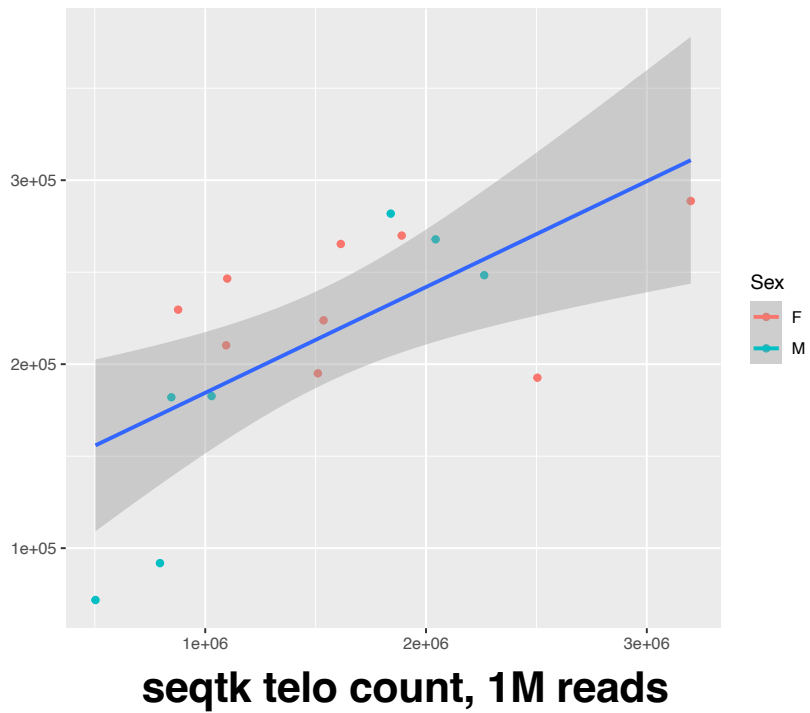
genotype
het
hom

SVs counted in diploid mode with svim-asm: Heller et al. 2021. *Bioinformatics* **36:** 5519-5521.

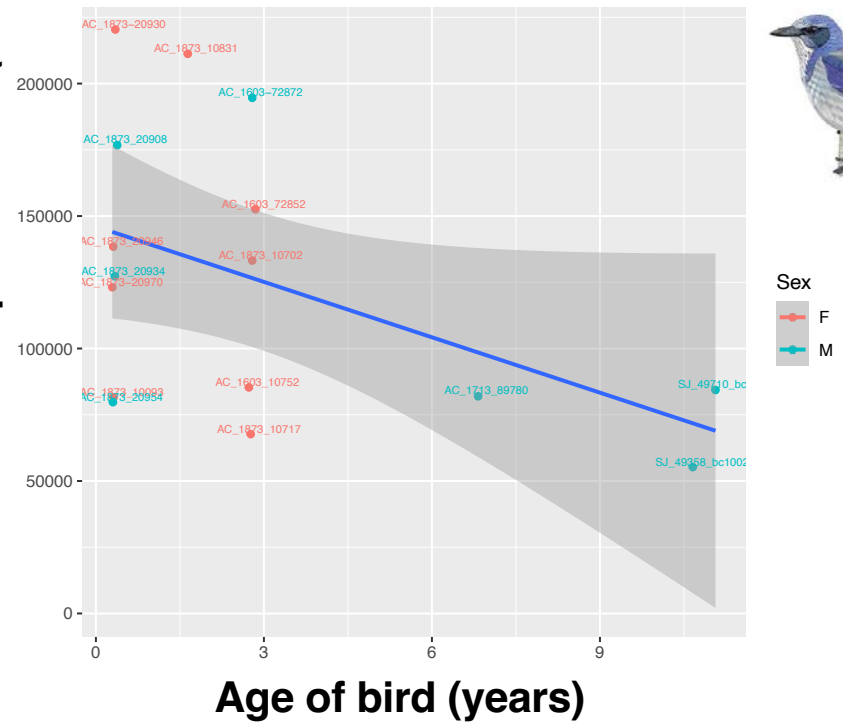# Telomere sequences are generally found at chromosome ends

# Telomere abundance declines with age in Florida birds



using seqtk telo, H. Li unpubl.

# Structural variants affect gene expression
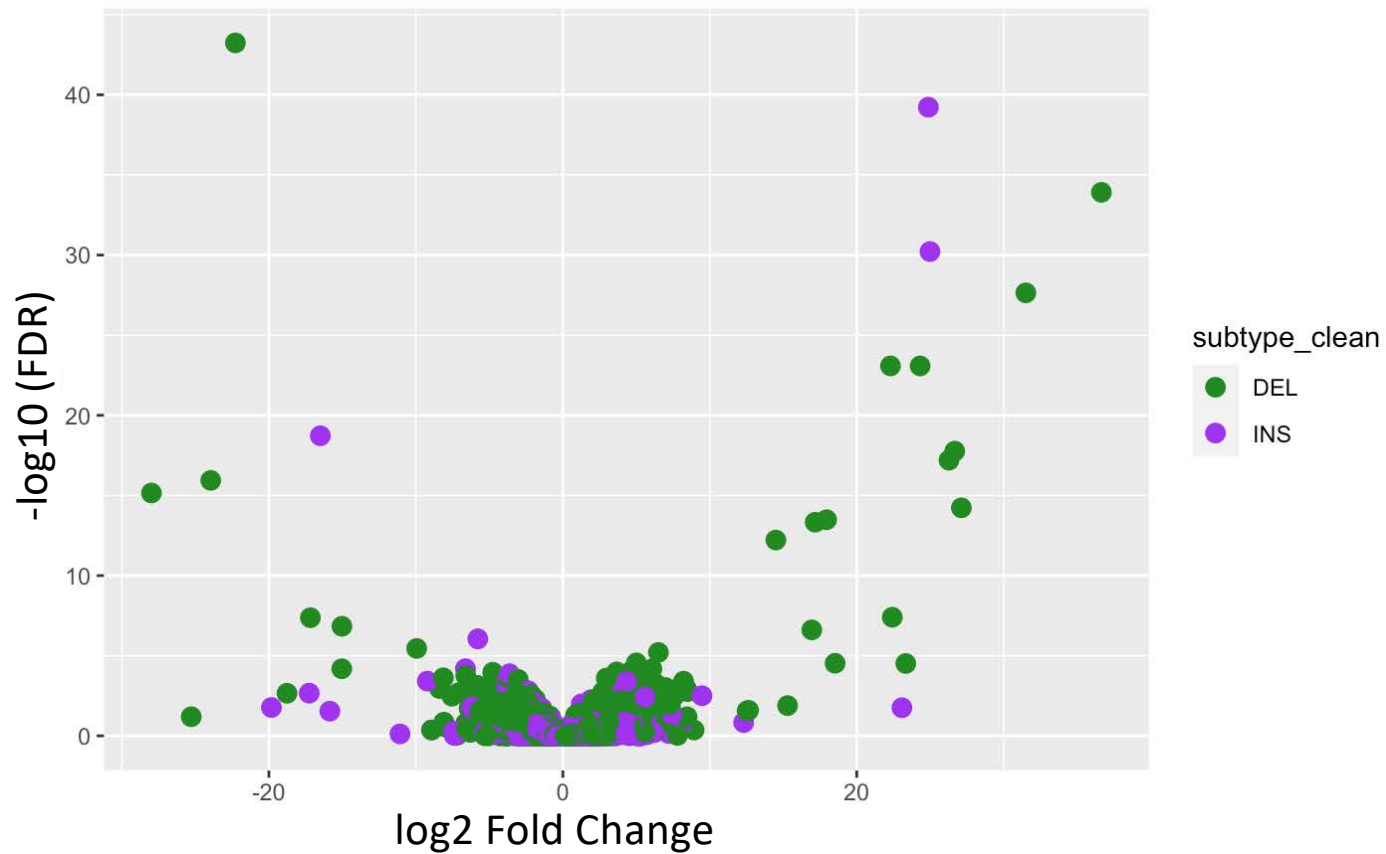


Significant (5% FDR)

● TRUE
● FALSE

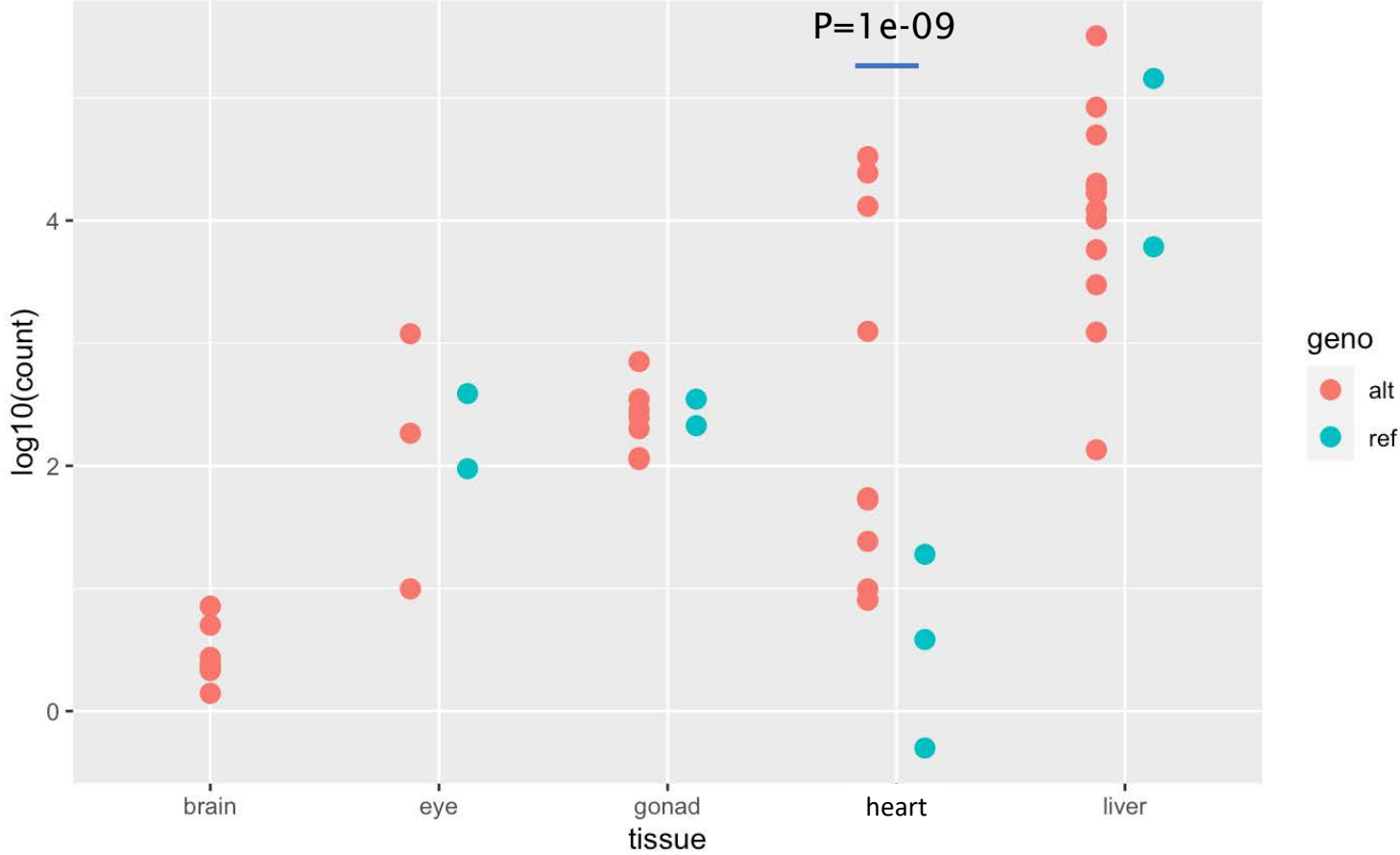58 SVs that are associated with expression differences of nearest gene, at 5% FDR

# SVs influence tissue-specific gene expression
## 16 bp deletion in CNEE near CDH4 gene

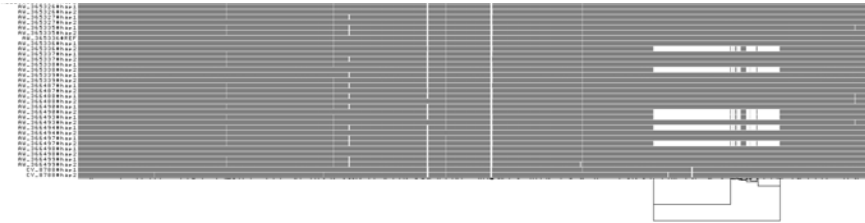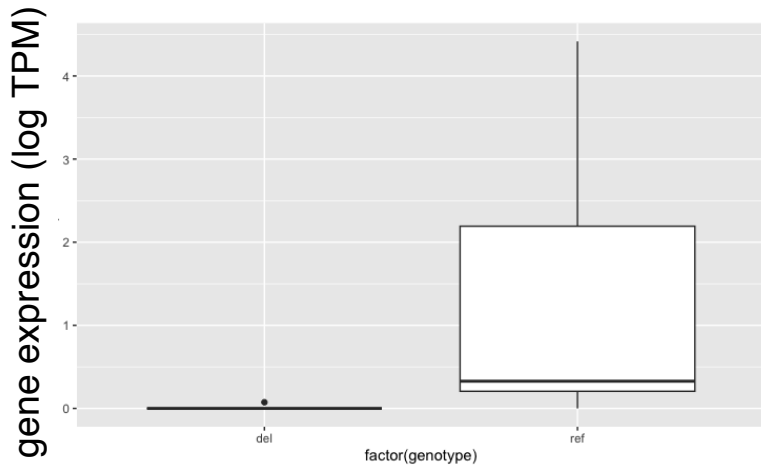SVs influence tissue-specific gene expression
43 bp deletion near THRSPB

# Visualizing examples of SVs influencing gene expression
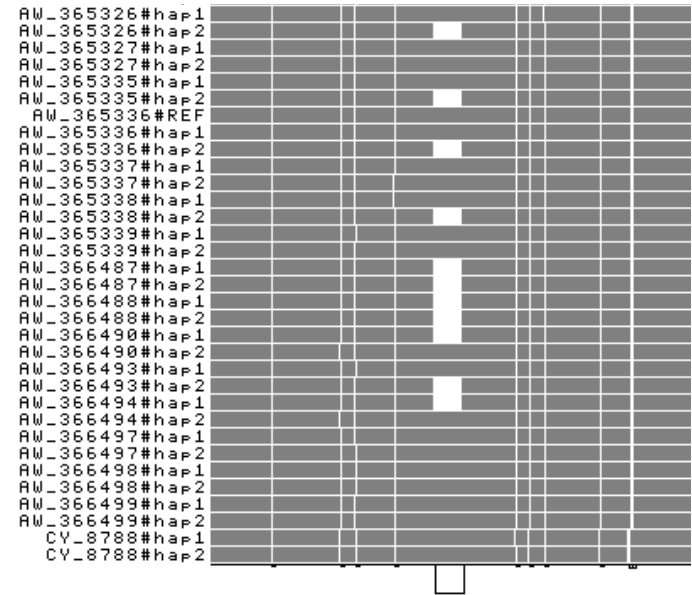
gene-LOC116806707 – some AW birds with
0 and some with > 100 tpm



odgi viz 1D visualization of region around
gene-LOC116806707

# Another example of SV and gene expression

17bp deletion in a CNEE near the GTPBP2 gene
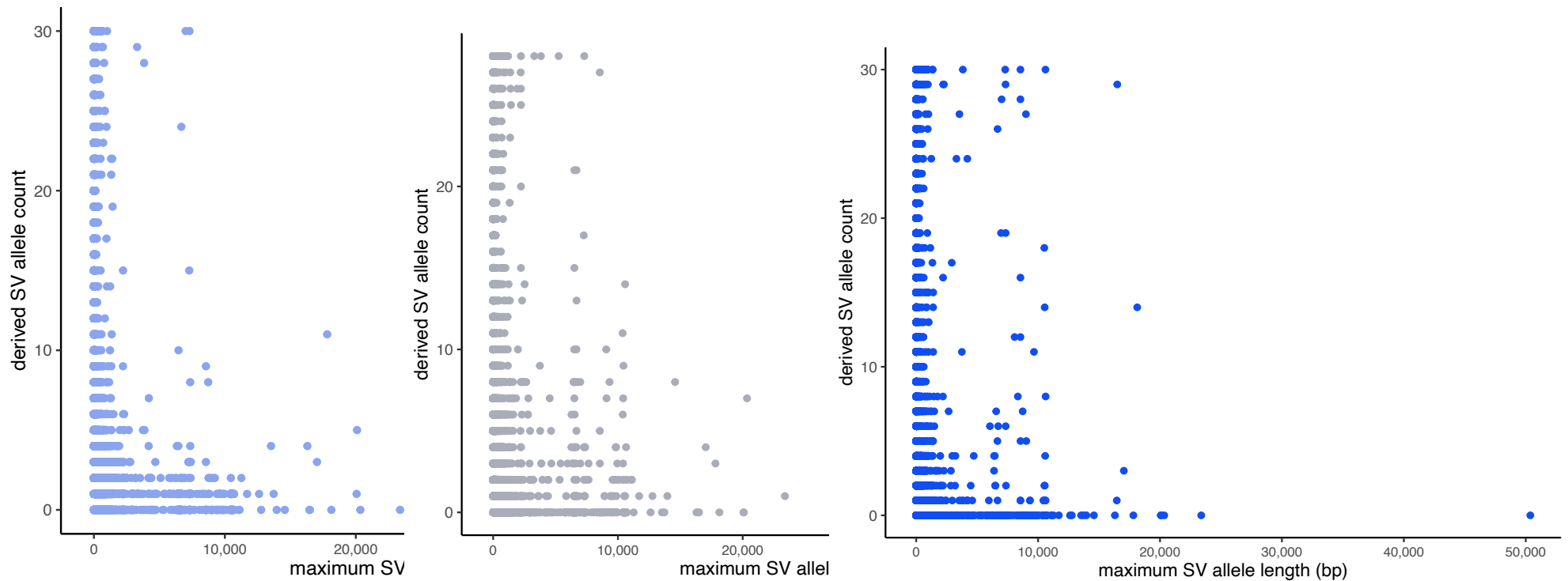
# Long SVs increase in frequency in small populations
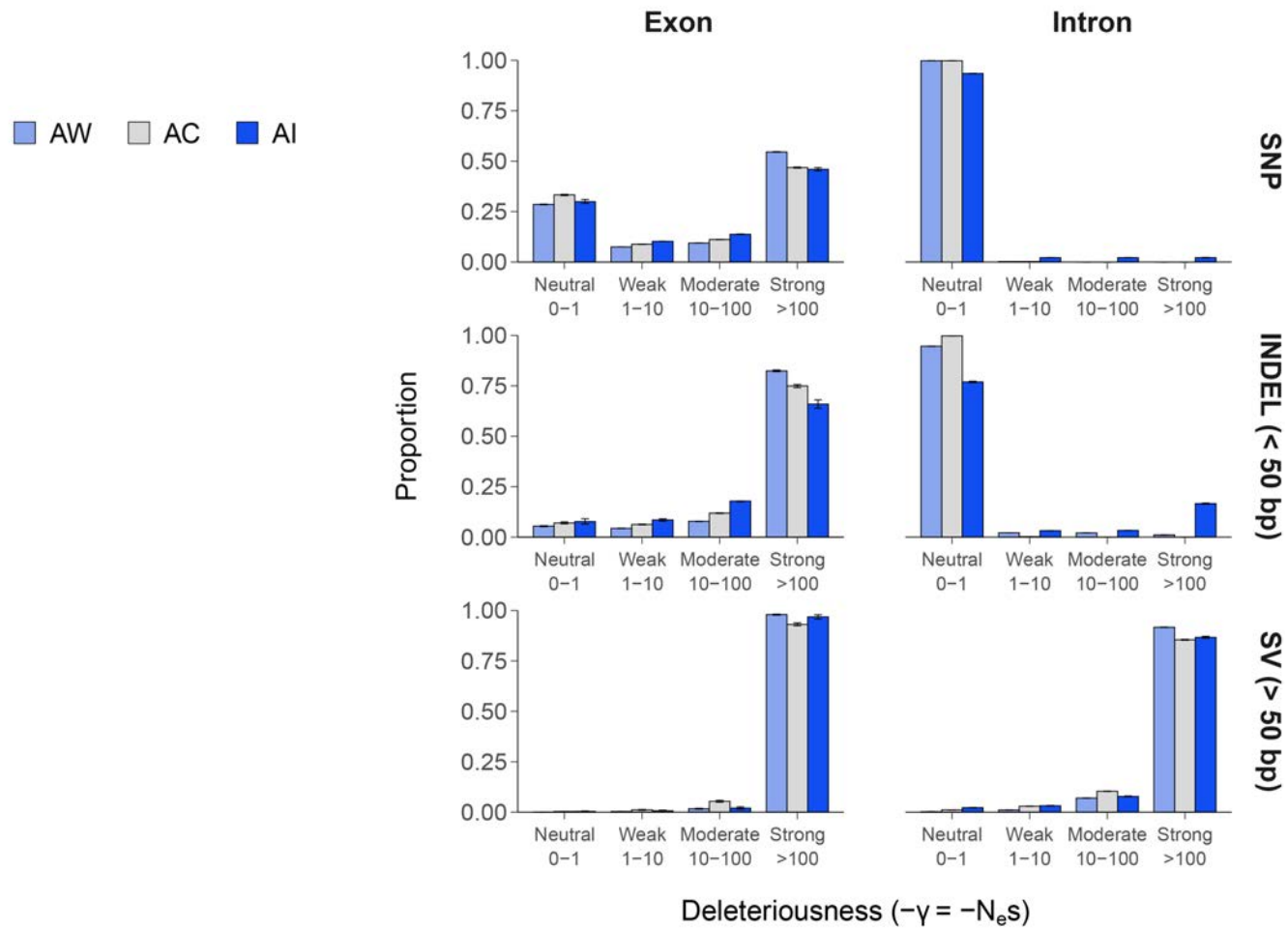
**Woodhouse's Scrub Jay**   **Florida Scrub Jay**   **Island Scrub Jay**



4,032,715 SVs from PGGB vcf downsampled to 100k SVs

# Distribution of fitness effects  of SVs

# 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

**Pangenome informatics**
Erik Garrison
Andrea Guarracino

**Harvard team - Woodhouse's Scrub Jay and Informatics**
Tim Sackton
Danielle Khost
Heng Li
Bohao Fang
George Kolyfetis

**Fieldwork**
Greg and Donna