

Session 4:

Challenging genomes to curate and strategies to work with them

Genome Reference Informatics Team (GRIT)
Wellcome Sanger Institute - Tree of Life

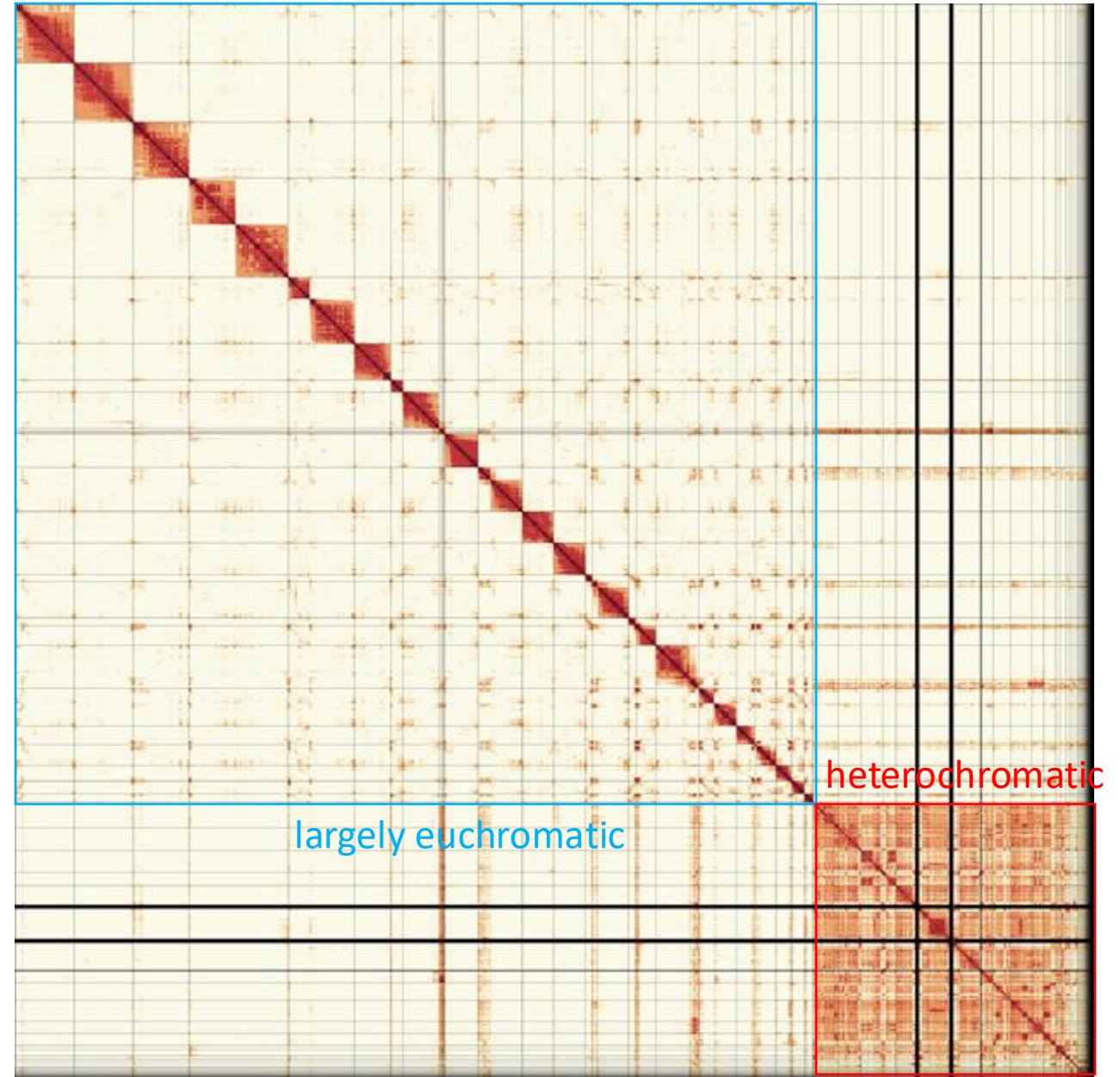
High repeat and heterochromatin content

Contrast between **euchromatic** and **heterochromatic** portion of the genome

Non-repetitive HiC signal can be seen for 26 chromosomal entities, in stark contrast to the heterochromatic portion of the genome (centromeric and short-arm sequences which in the case of this wasp do not have enough specific association with a particular chromosome to enable them to be placed.



iyNysSpin1_1



High repeat and heterochromatin content

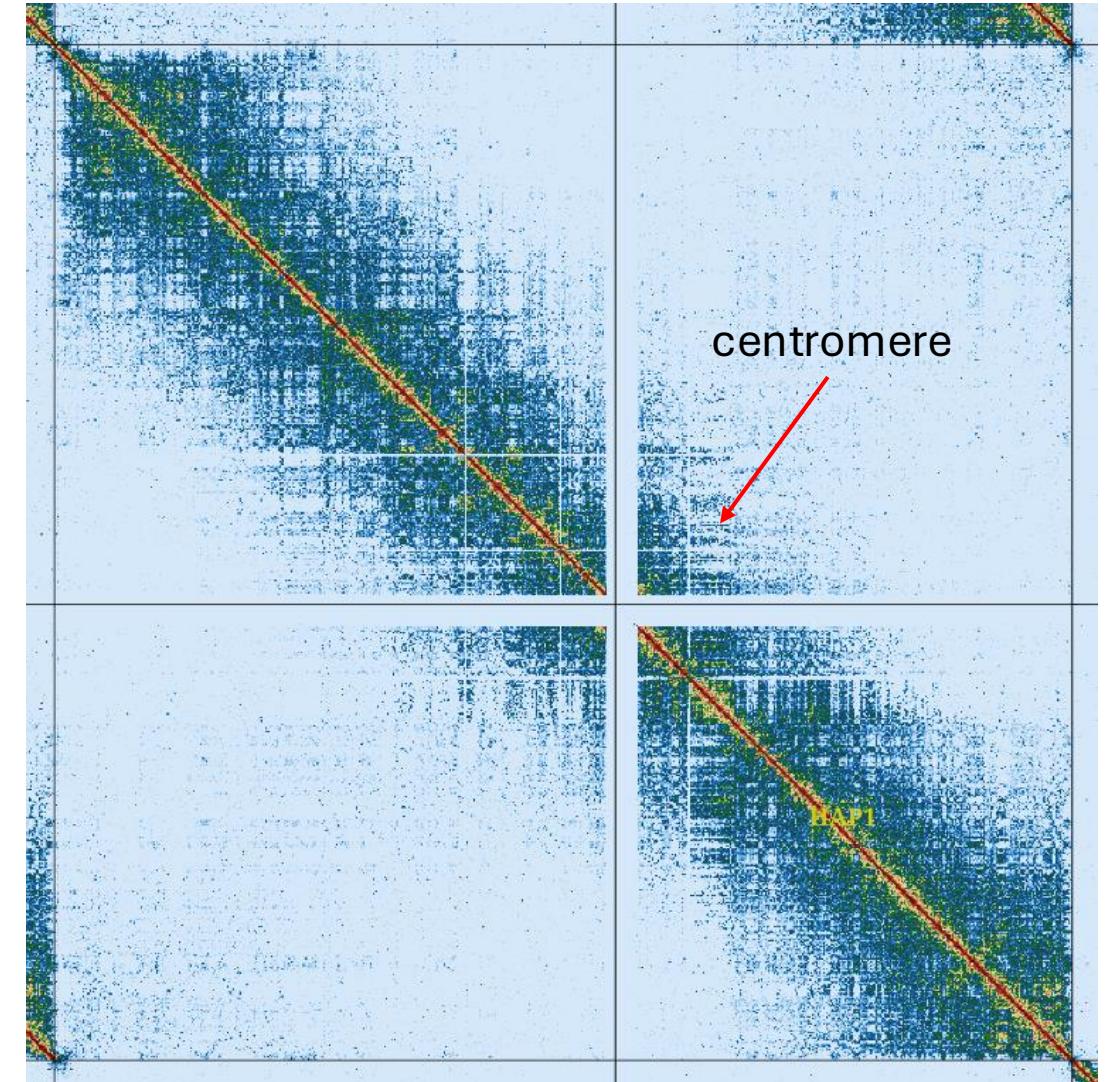
HiC bias

- More represented:
- High GC content: due to PCR and sequencing bias

- Less represented (Low mappability):
Repetitive regions: Centromers, telomeres and repetitive regions

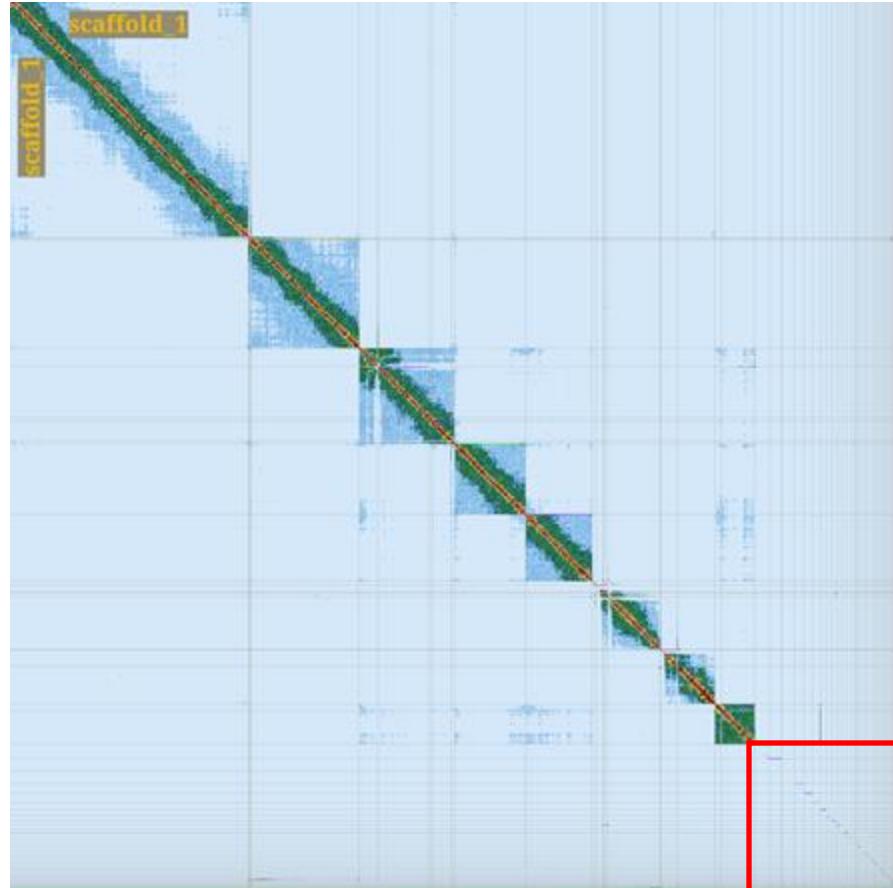
HiC reads align to many regions

Worse when you have HiC low coverage

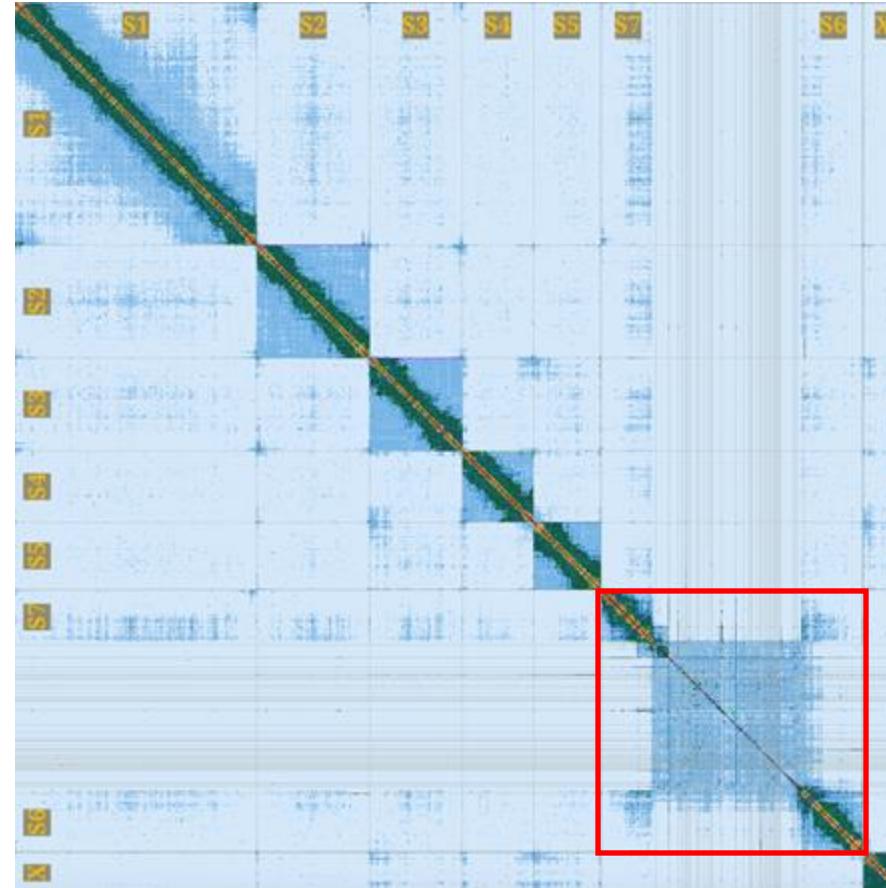


Multi-mapping + karyotype

Multi mapping reads reveal hidden linkage between ‘separate’ chromosome scaffolds and blank repetitive scaffolds



multi-mapping ‘off’



multi-mapping ‘on’



Rhagonycha fulva

Karyotype image confirms presence of large heterochromatic chromosome



Microchromosomes

Birds, sharks and reptiles (only?)

Bird genomes are organized in macro- and micro-chromosomes

(By Tom Mathers)



Chicken chromosomes (n = 39)

Masabanda et. al. 2004, *Genetics*

The chicken genome is typically divided into 10 **macrochromosomes** (>23 Mb in length) and 29 **microchromosomes** (22 – 2.5 Mb in length).

Microchromosomes can be further divided into **micro** (22 – 5 Mb, n = 19) and “**dot**” (< 5 Mb, n = 10) chromosomes.

Dot chromosomes represent a major challenge for assembly and curation.

* Sizes based on latest chicken near-T2T assembly (Huang et. al. 2023, *PNAS*).

In cuckoo, the 10 smallest chromosomes represent 1.2% of the genome assembly but contain 10% of the protein coding genes.

Microchromosomes

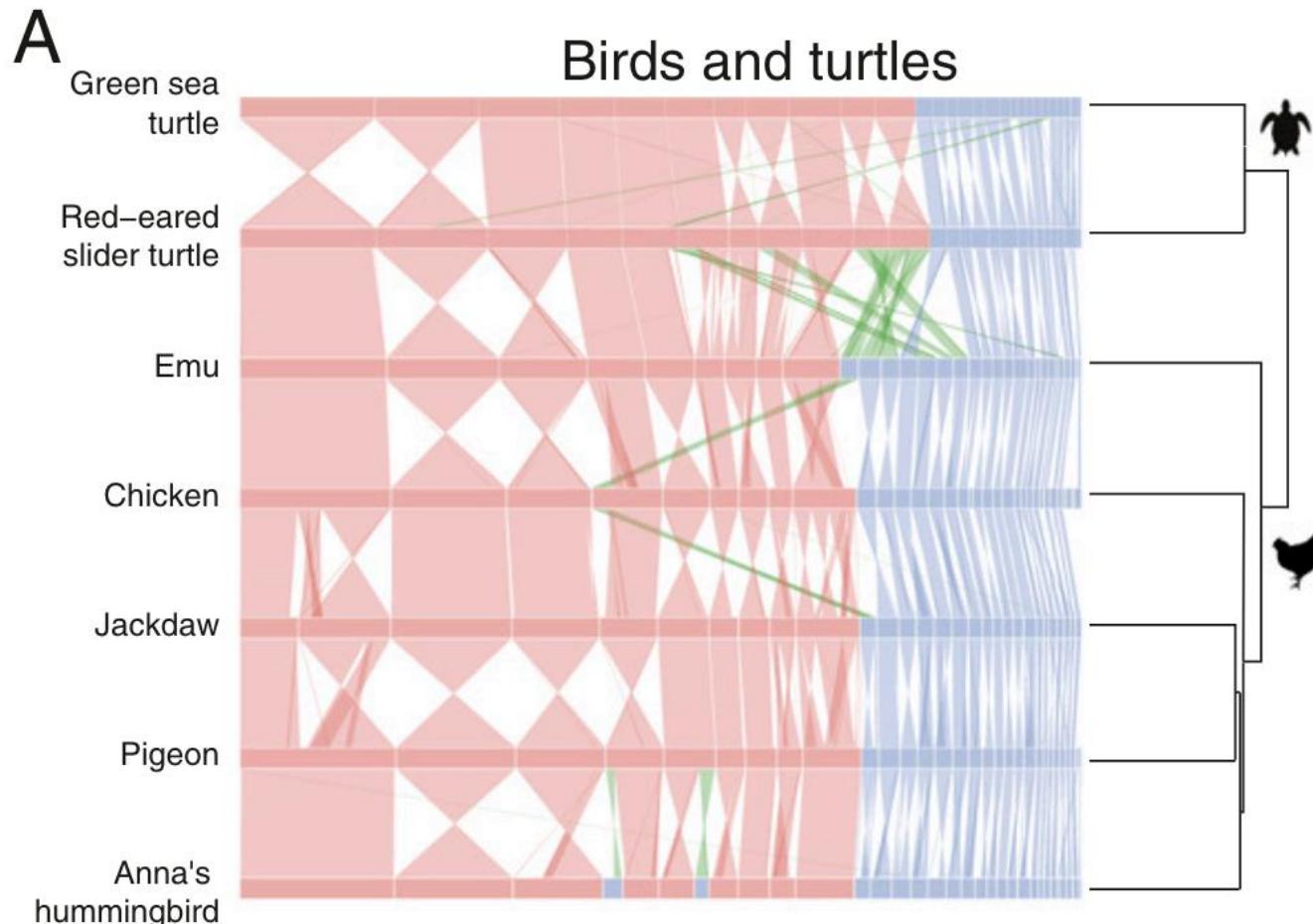
Less than 2% of the assembly and 99% of the effort....

Microchromosomes are often highly fragmented and mixed in with repeat scaffolds in assembly “shrapnel”

Considerable gene content

Microchromosomes are highly conserved across birds and reptiles

(By Tom Mathers)

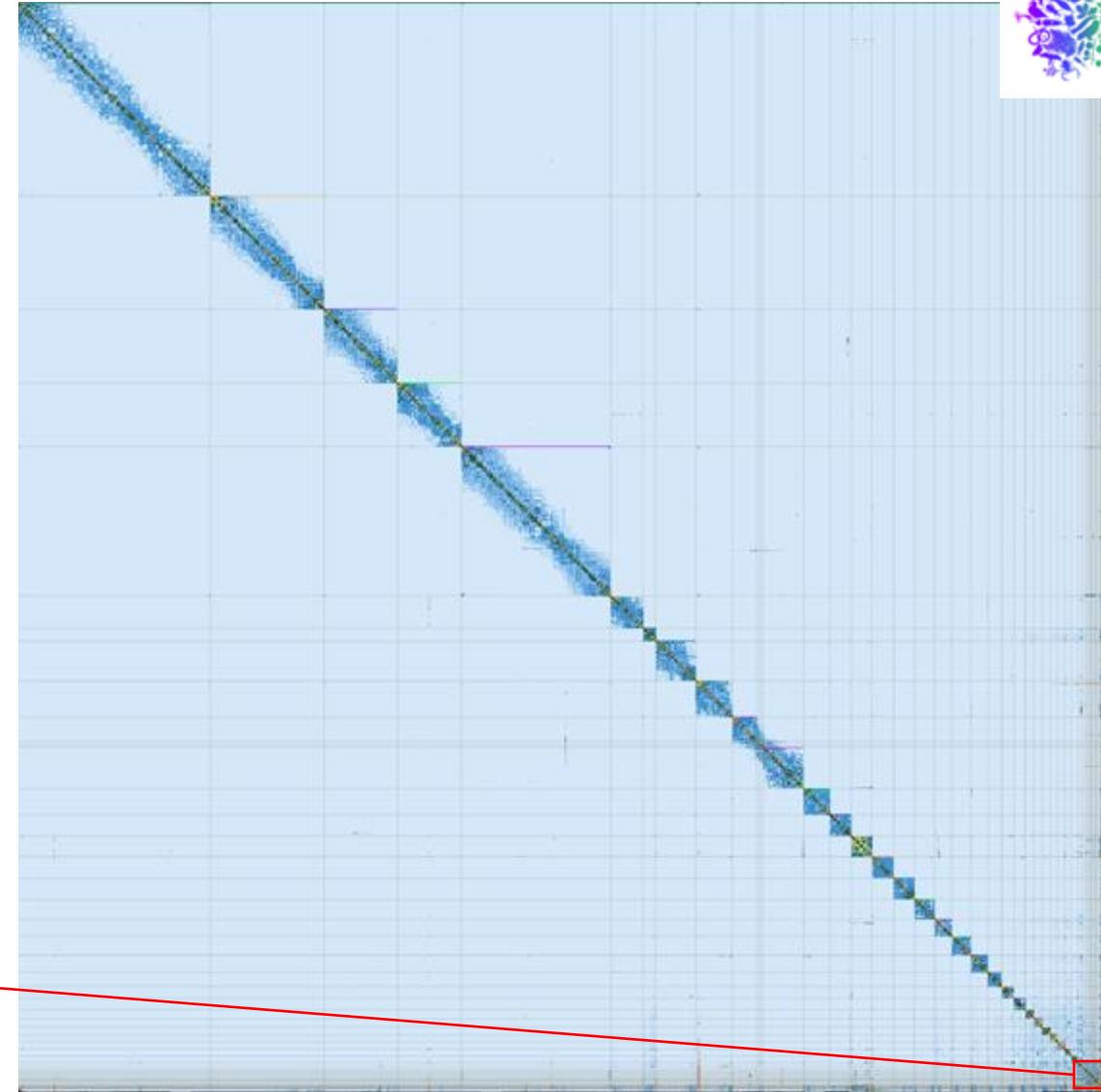
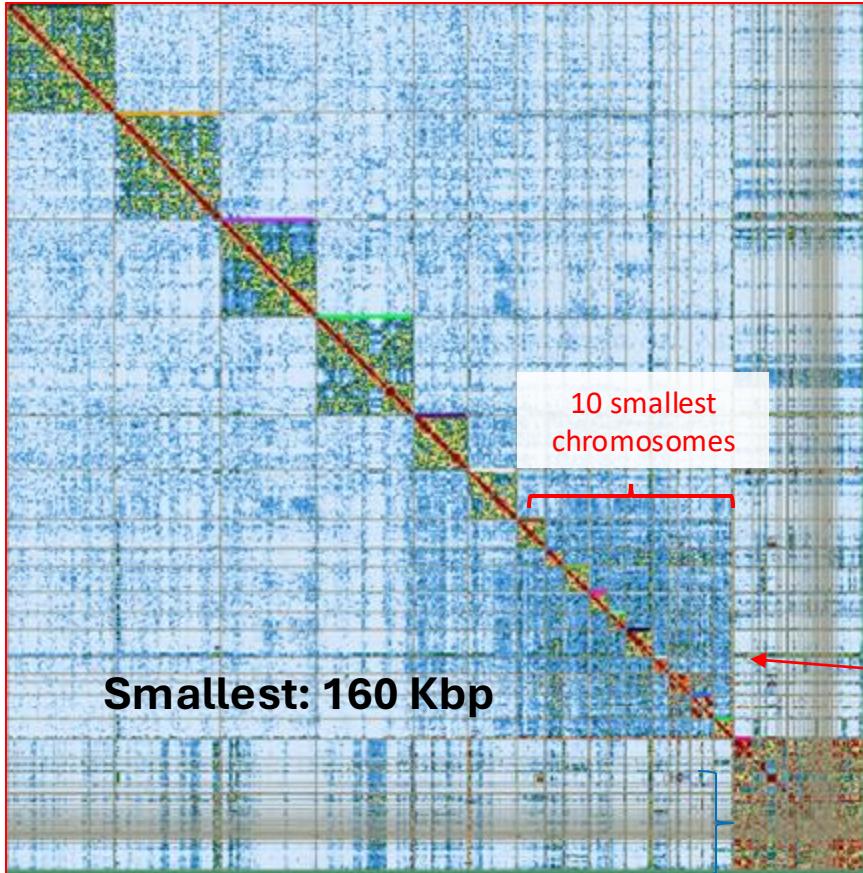


Micro-chromosomes - Birds

(bCucCan1)

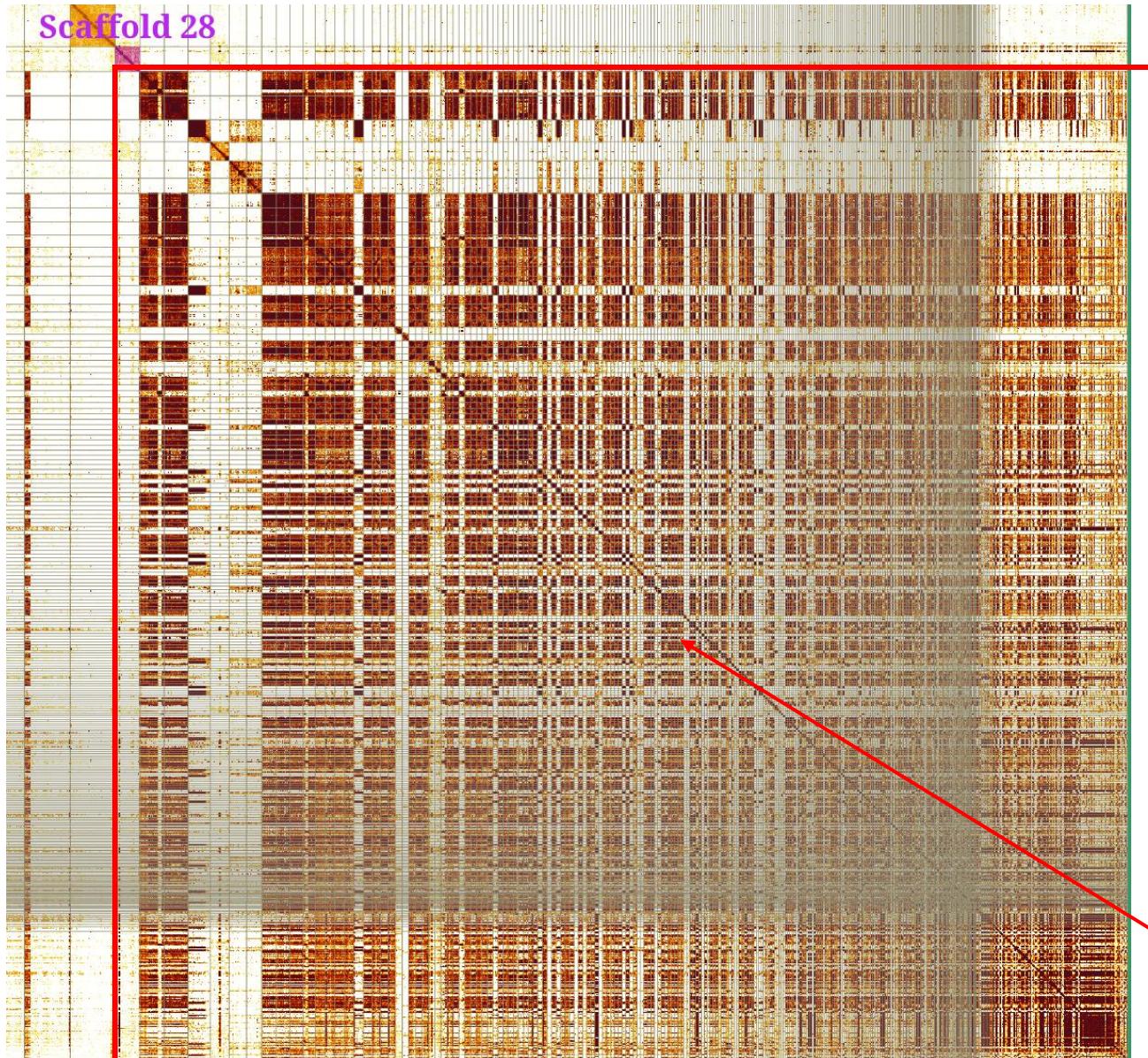
HiFi data

- Disproportionate amount of time curating the **smallest 10 micro-chromosomes** (<1.2% of the assembly)....



Microchromosomes

(By Tom Mathers)



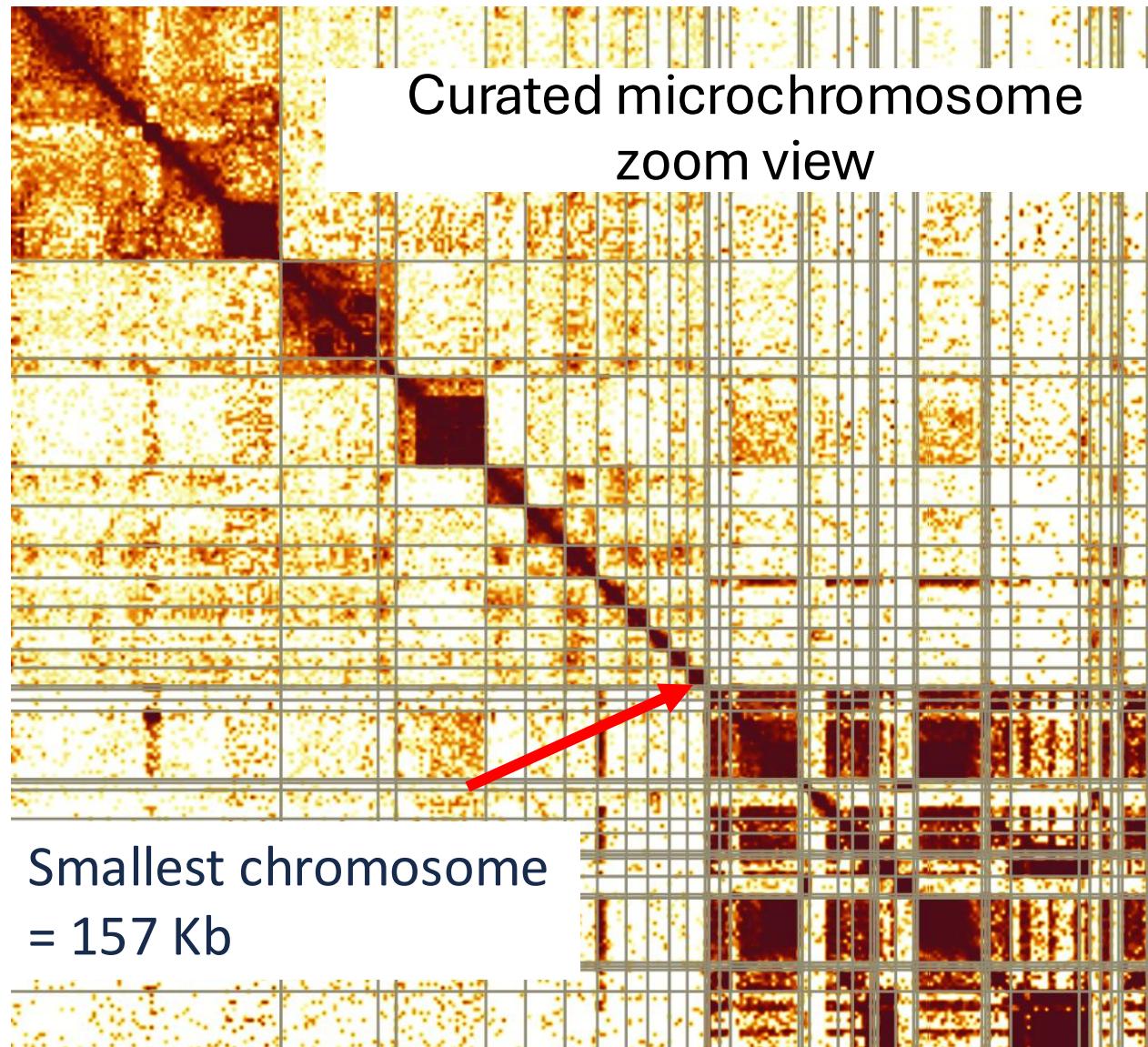
Quick curation of larger scaffolds only recovers 28 chromosomes.

Expected karyotype is 39 autosomes + Z + W

Remaining 13 chromosomes are somewhere in here!

Micros ???

Microchromosomes



To find these missing chromosomes we **rely on elevated background HIC signal between micros.**

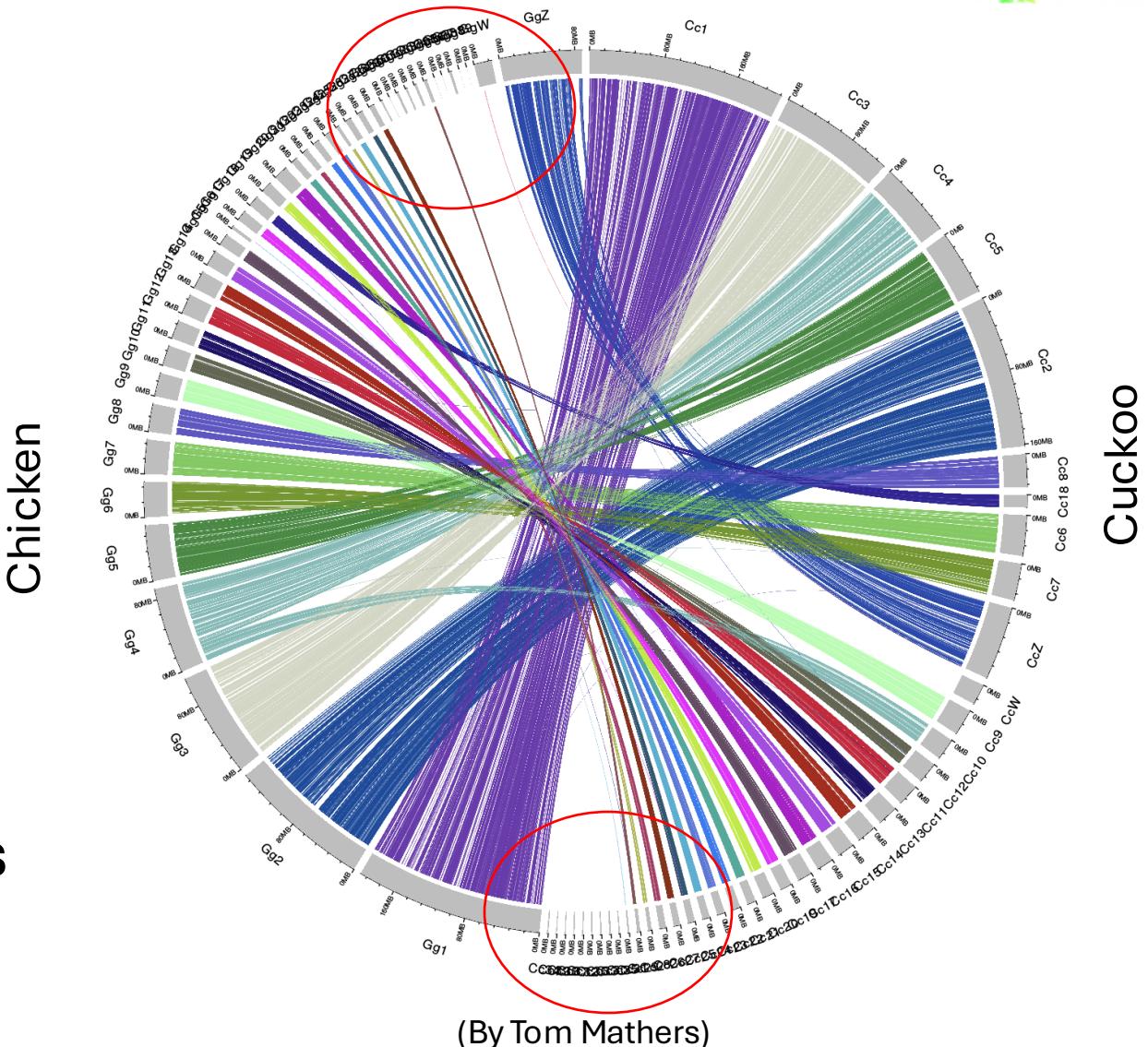
Additionally,
JBrowse looking for gene-rich small contigs

Manually inspect whole genome **alignments** vs well assembled relatives (if present).

Very time consuming.

Most challenging group: birds

- Very time consuming
- Highly repetitive and fragmented (ideal is a hybrid assembly using HiFi and ONT reads)
- Really tiny
- BUSCO Aves ODB 10 gene set does not cover all microchromosomes



How do we fish out the micros?

Main approaches we use for birds:

1. MicroFinder script for birds (HiFi/ ONT)

Miniprot to **map a set of conserved microchromosome-associated proteins** to a draft assembly and then **counts the resulting hits and orders the input assembly by the number hits**

2. Nucmer (HiFi/ONT – sometimes)

3. Gene content in Jbrowse (HiFi only)



How do we fish out the micros? (Birds)

MicroFinder script for birds:

<https://github.com/sanger-tol/MicroFinder>

Recommended:

16 cores

24 Gb RAM

Scaffolds > 5Mbp will not be ordered

The script should be run for each haplotype separately:

```
MicroFinder.sh <hap1_fasta> scaffold_length_cutoff
```

```
MicroFinder.sh <hap2_fasta> scaffold_length_cutoff
```

scaffold_length_cutoff (Kbp)

It will:

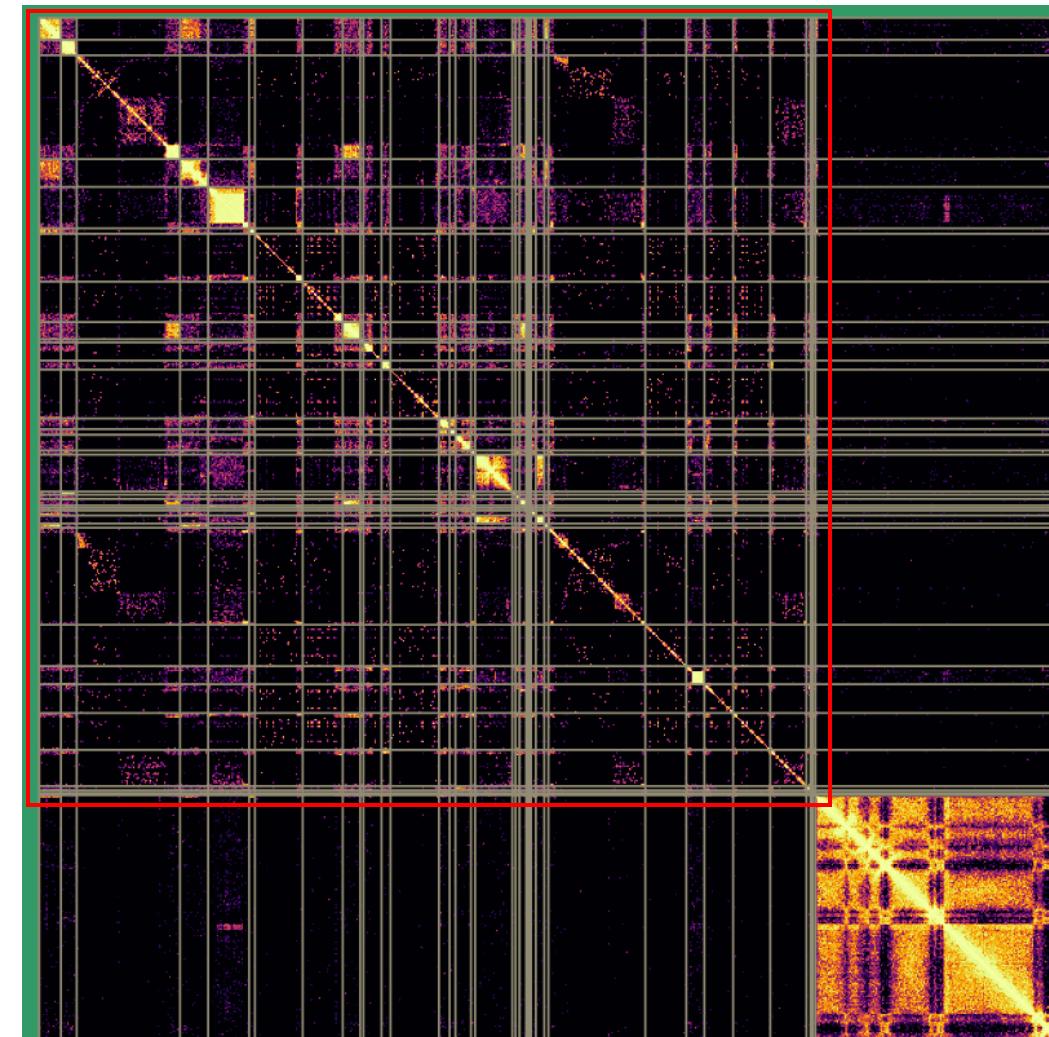
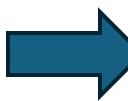
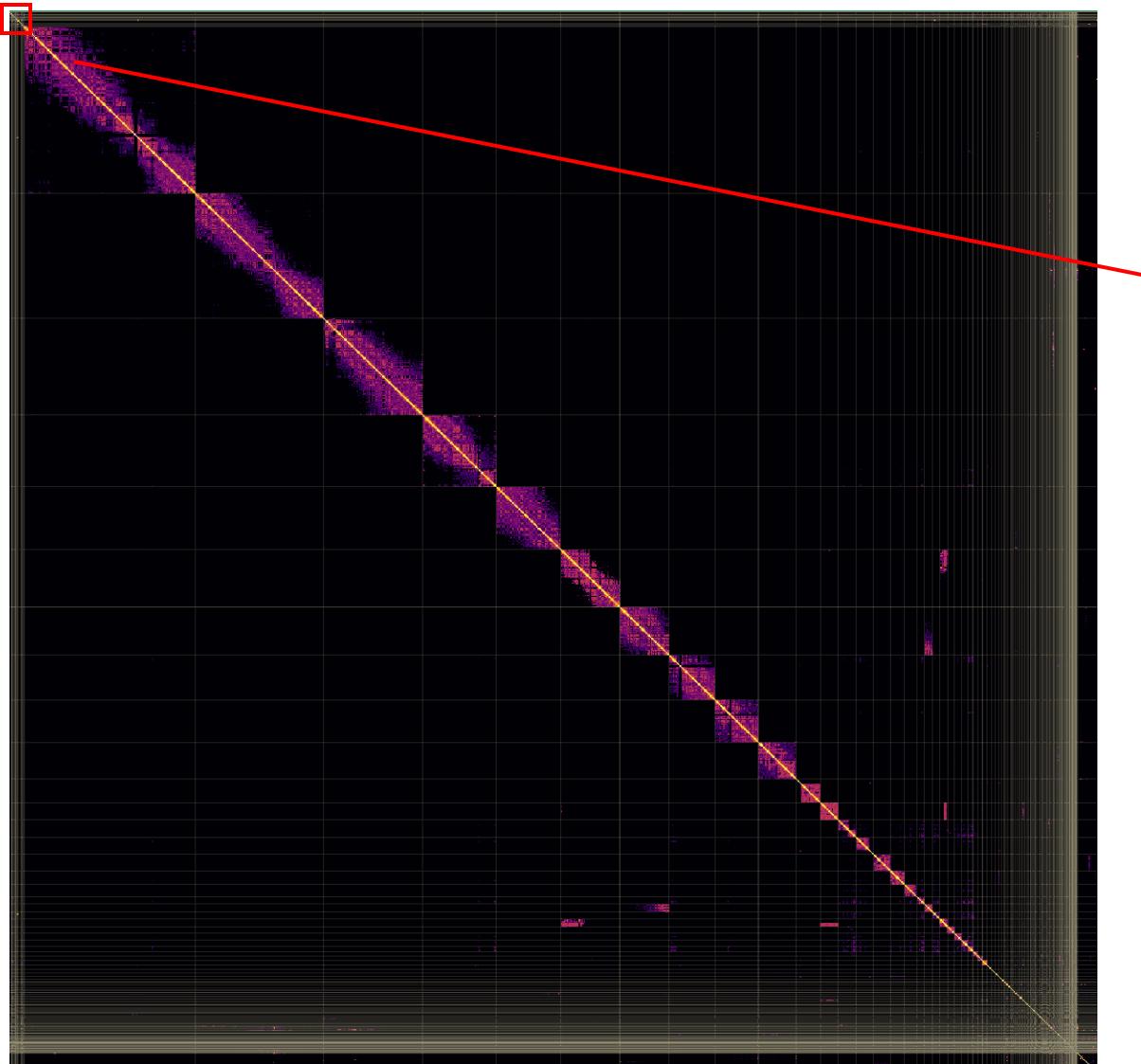
Align your genome to a conserved database of bird microchromosomes and look for gene content

Sort by number of gene hits and then by size (< 5Mbp only) and move them to the beginning of the fasta file

Generate a new fasta file

How do we fish out the micros? (Birds)

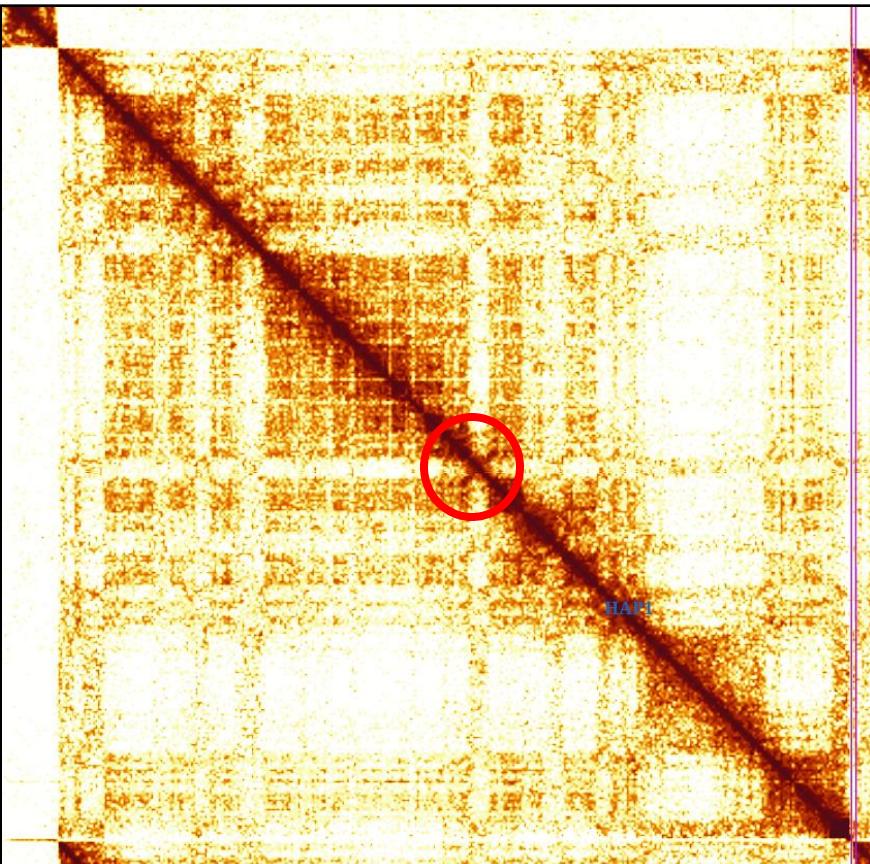
Potential micros will appear on the top left of hap1 and hap2 new Pretext maps



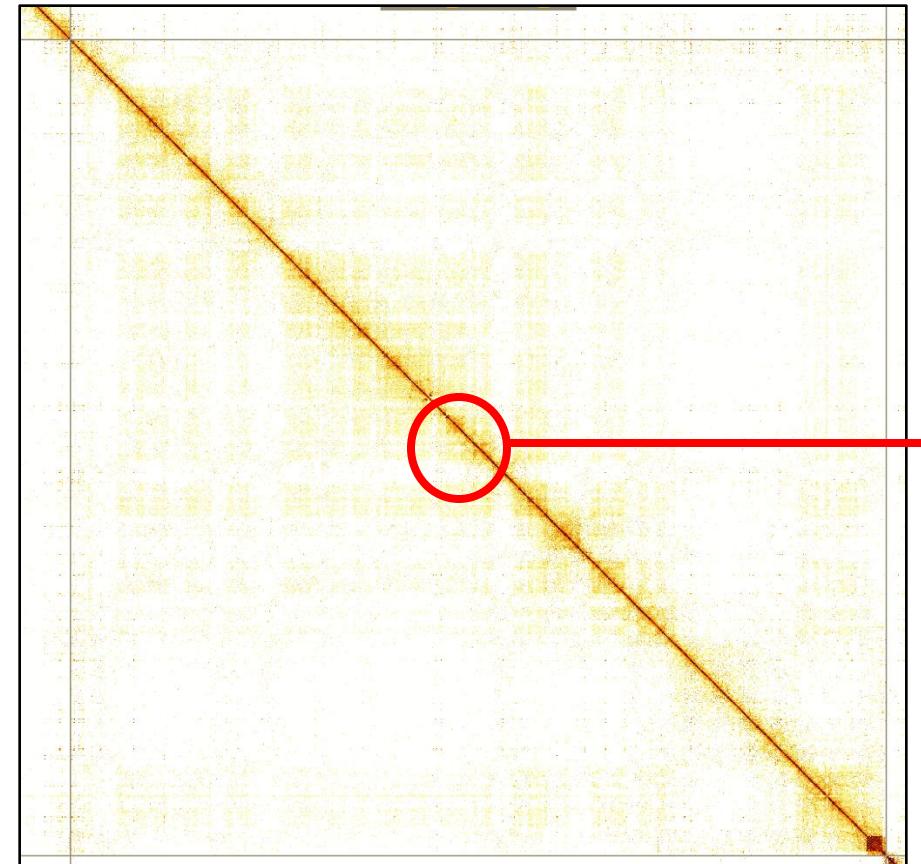
Ordered by gene count

Two pass assembly curation enables high resolution curation of microchromosomes

Full map

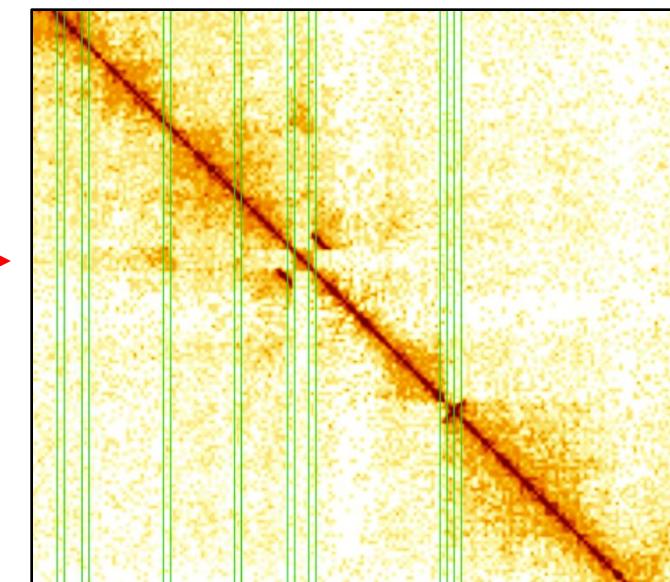


Second pass map
(scaffolds < 20 Mb)



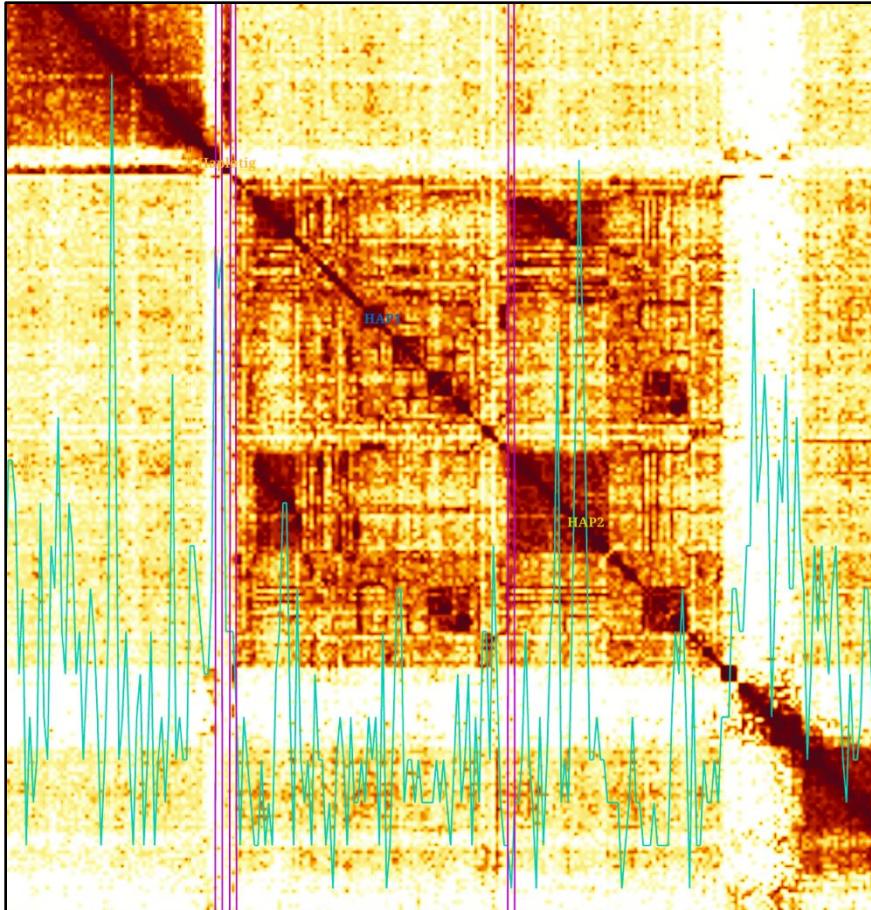
High resolution maps

Zoomed in second pass map
(scaffolds < 20 Mb)

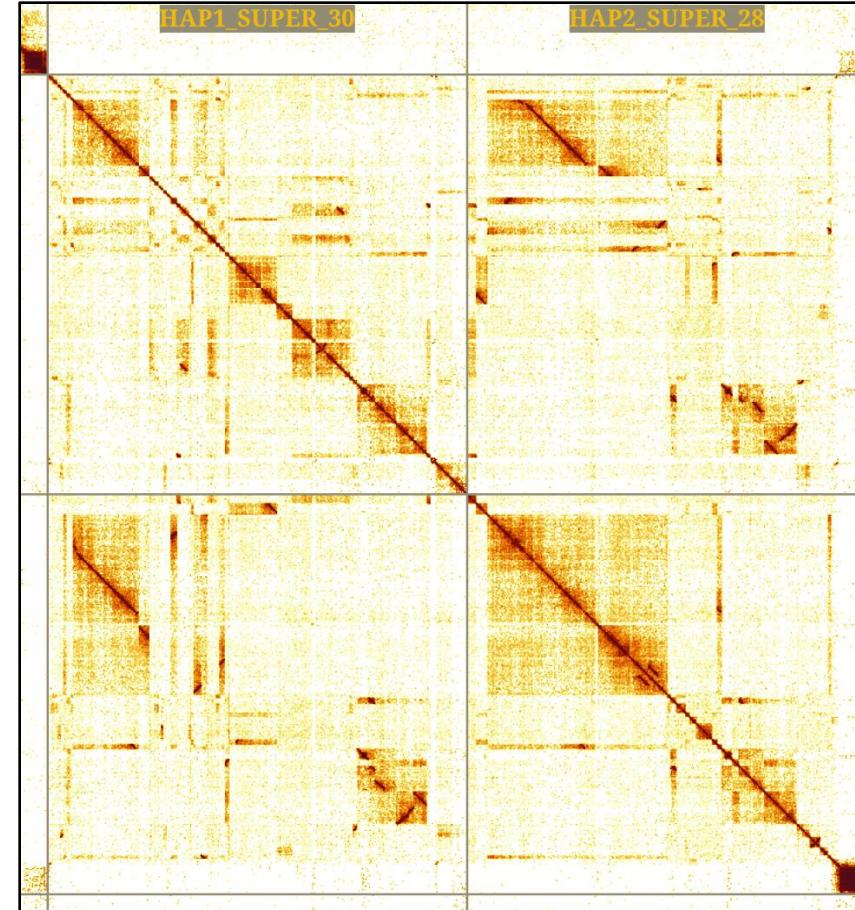


Second pass curation can reveal some horrors!

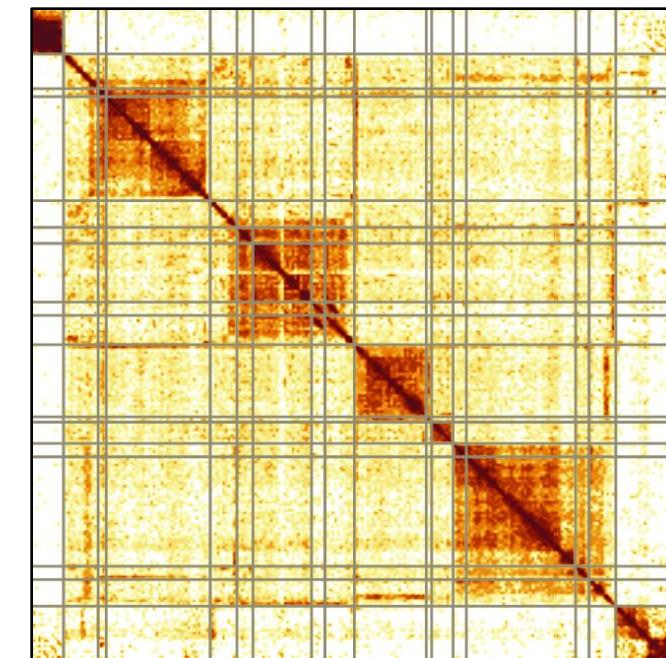
Full zoom of Hap1 and Hap2 of
chr 30 in full map



Second pass map
(scaffolds < 20 Mb)



Rearrangement shows
chr 30 is actually 5
chromosomes!
(Hap1 shown)

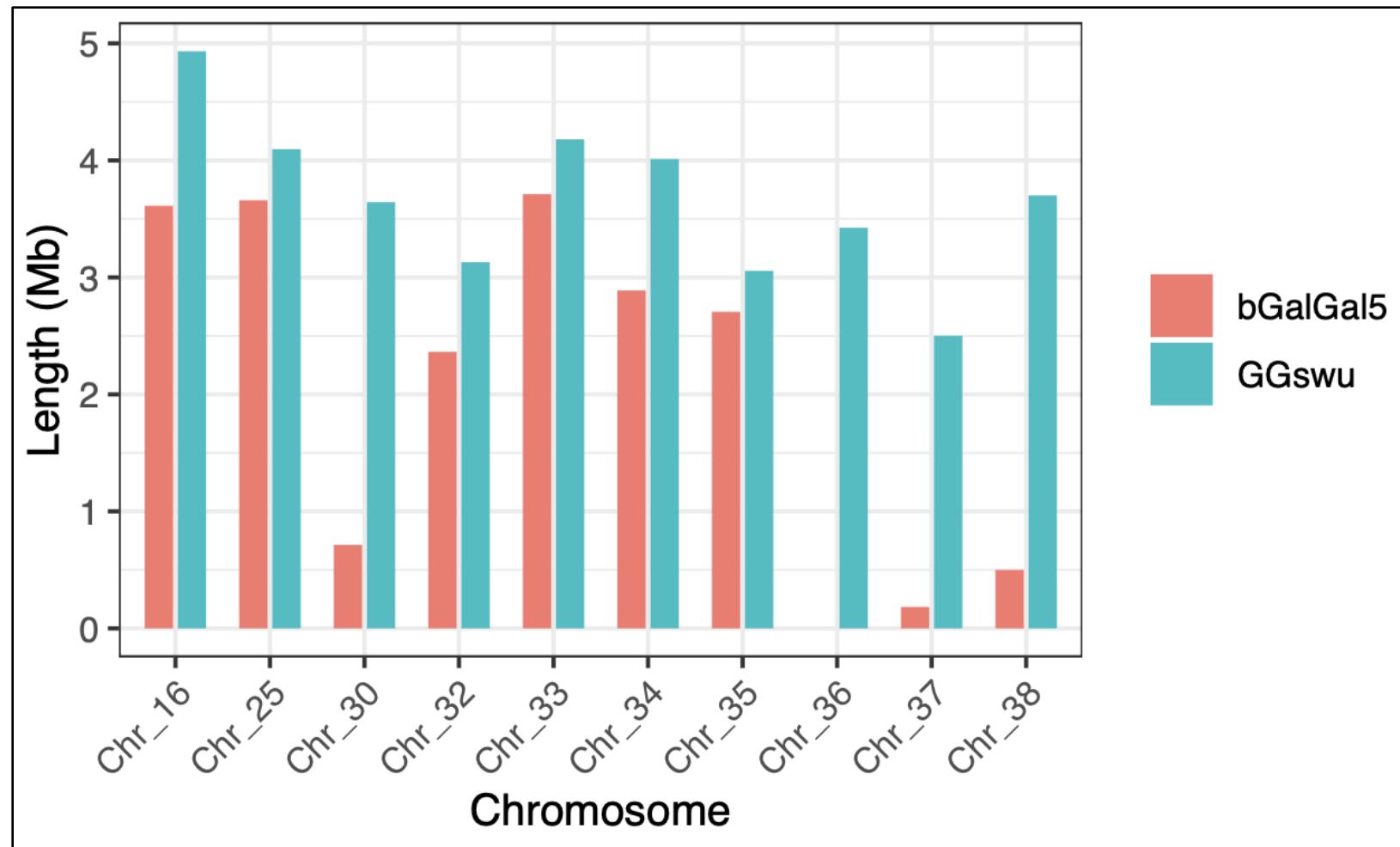


Is it possible to have more complete bird
microchromosomes?

Nanopore data looks to be promising!!!

Bird and fish gene-rich regions

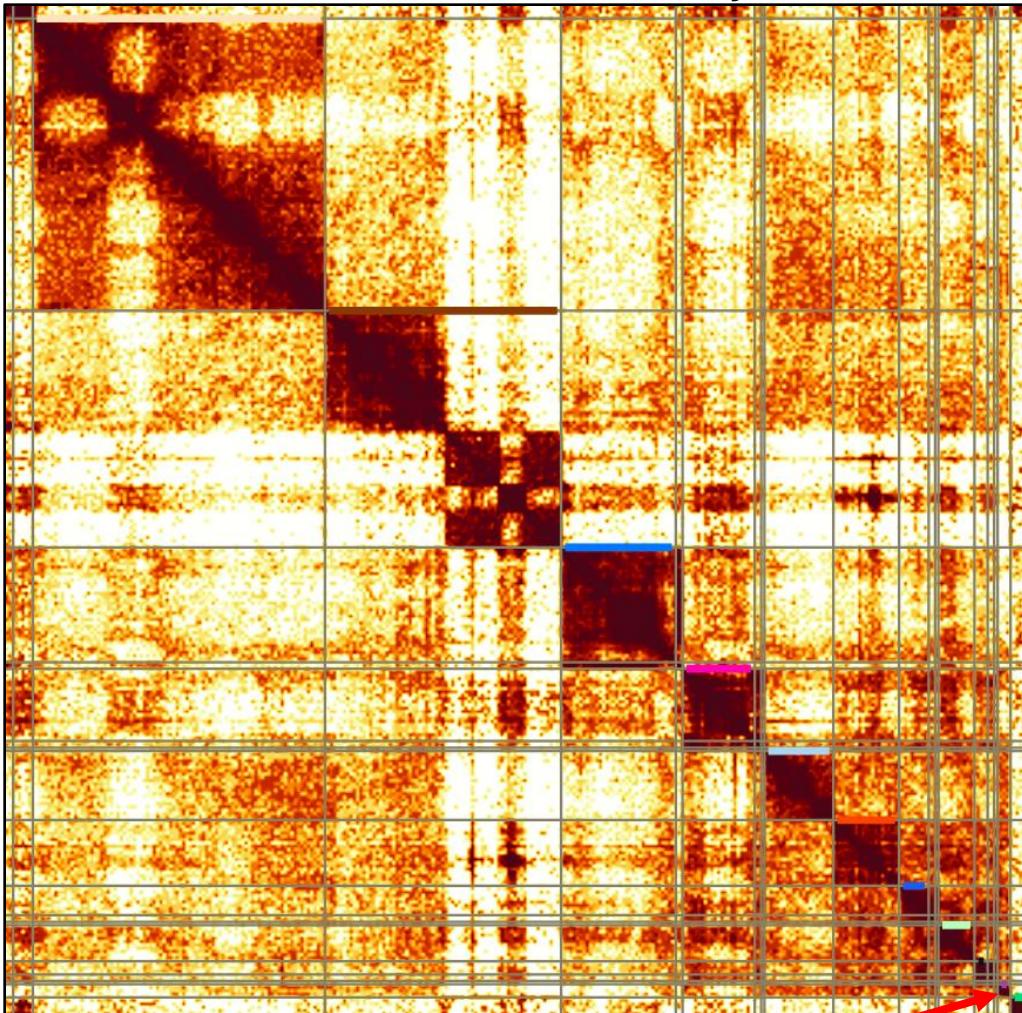
“Dot” chromosomes are substantially shorter in bGalGal5 (HiFi) than Ggswu (HiFi + ONT)



*In contrast, the size difference of the macrochromosomes is < 5%

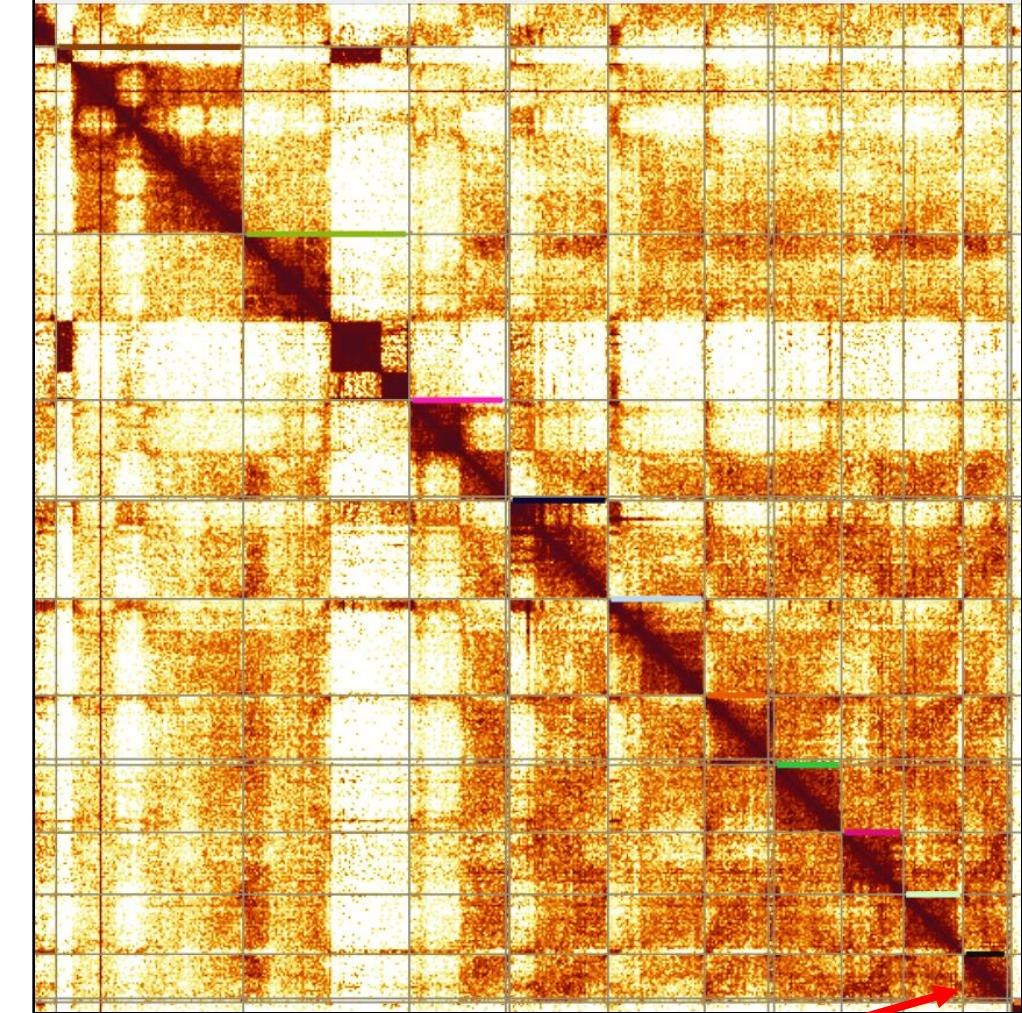
bAytFul3 *Aythya fuligula* (tufted duck) smallest 10 chromosomes

PacBio HiFi assembly



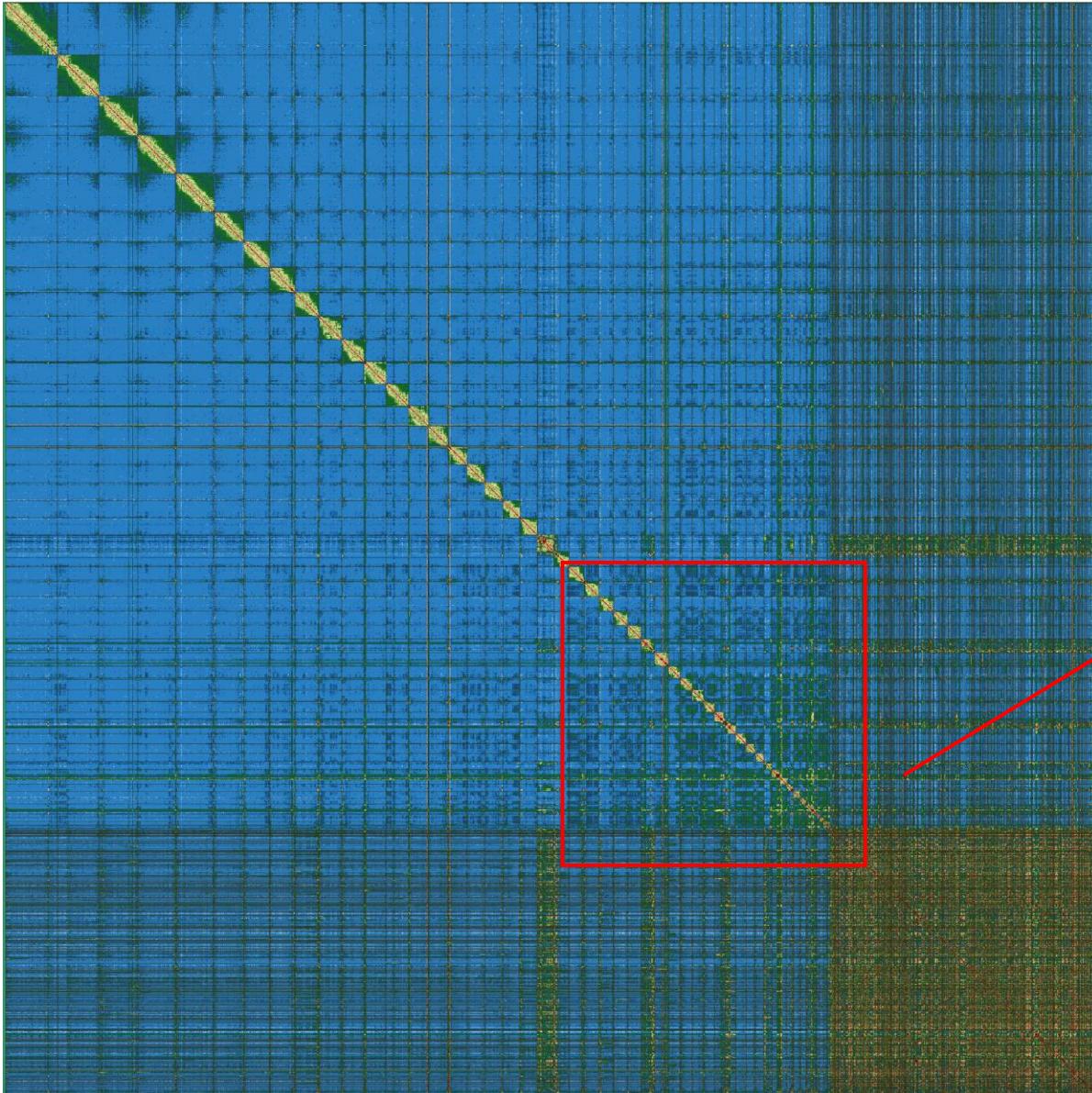
Chr 39 = 125 kb

Nanopore assembly

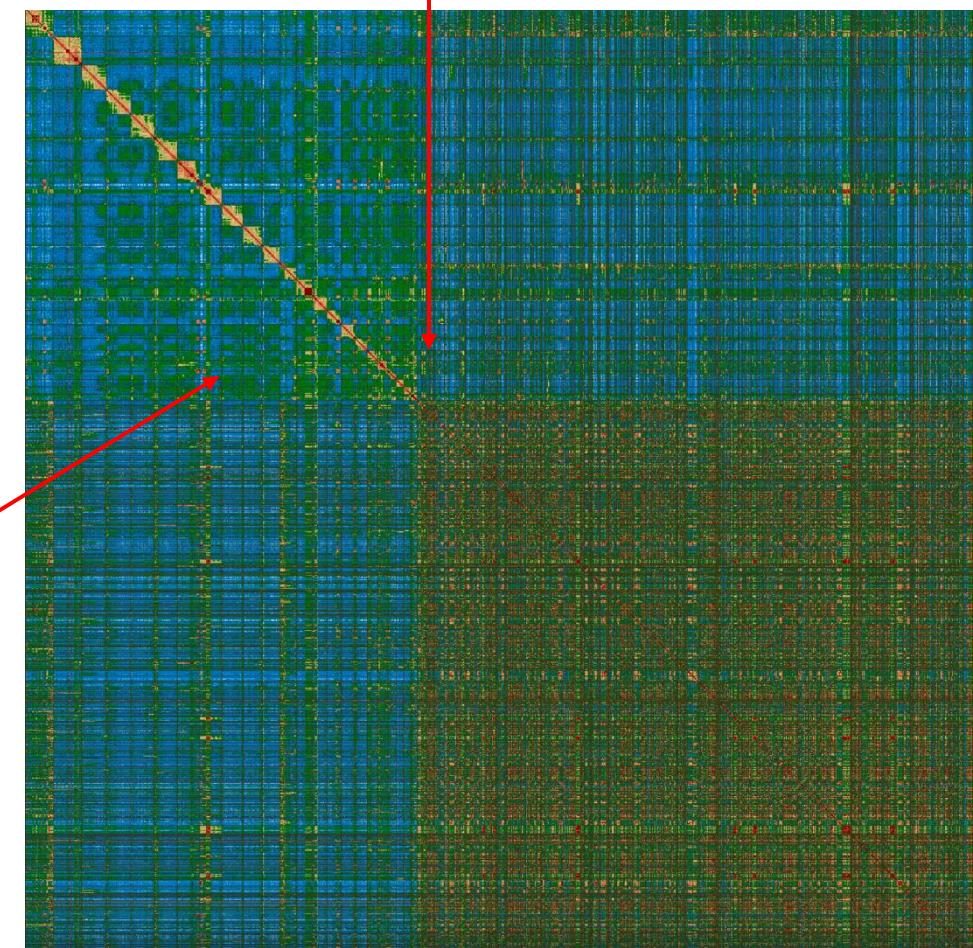


Chr 39 = 1.15Mb

Micro-chromosomes - Sharks

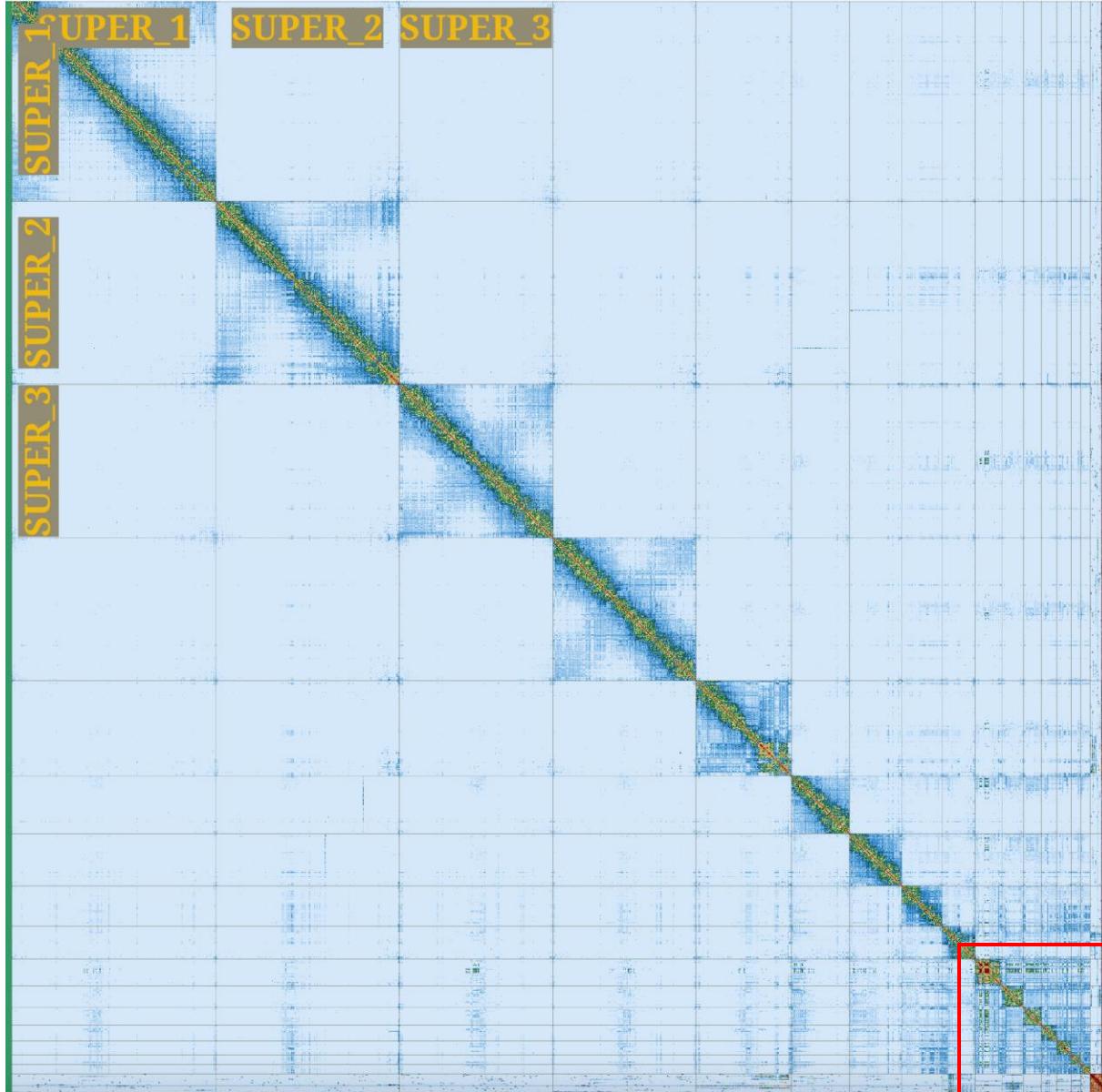


16 micros
Smallest one: < 9 Mbp

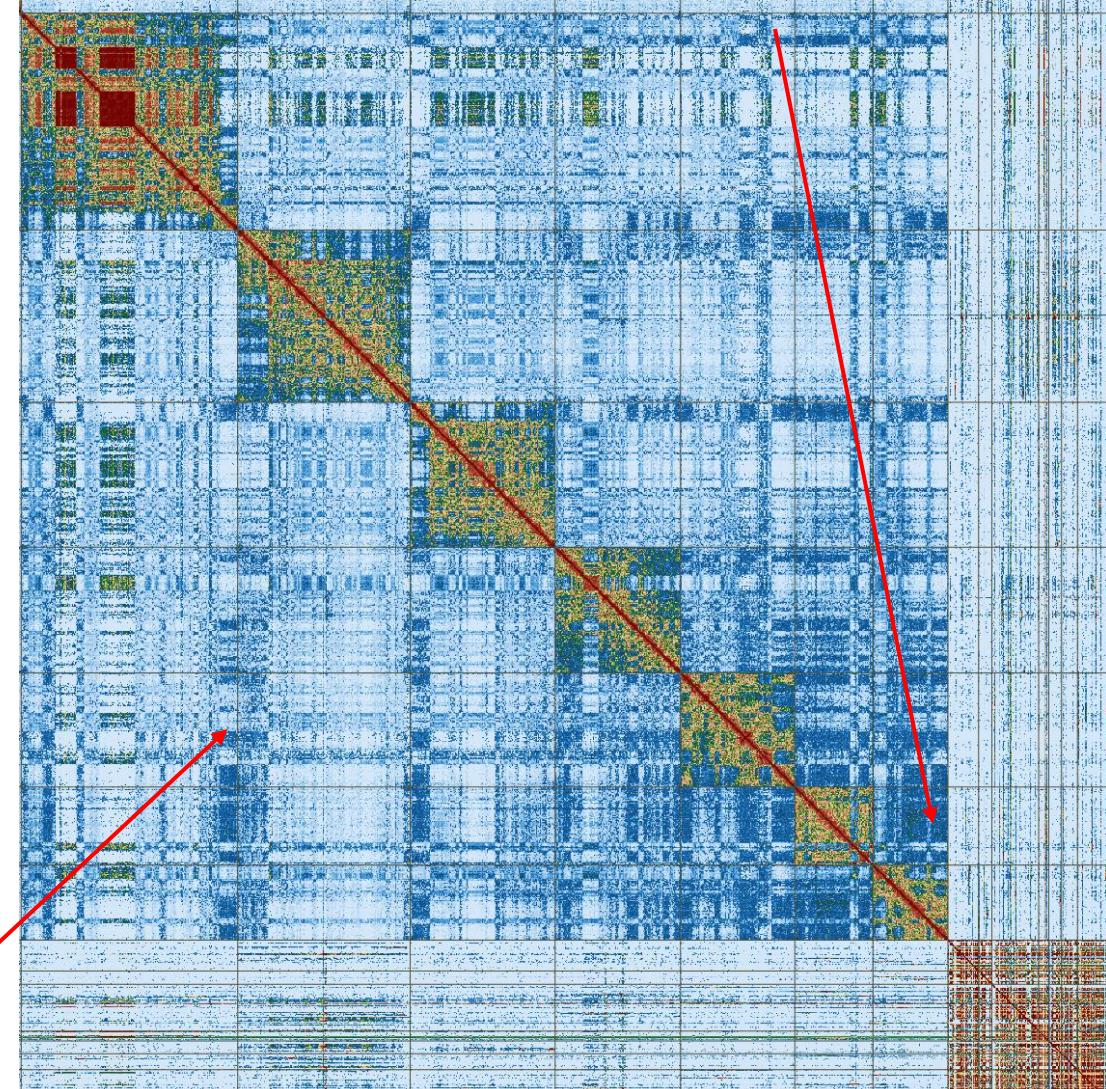


Main challenge:
Usually the most repetitive ones

Micro-chromosomes - Reptiles



7 micros
Smallest one: 12 Mbp





**High chromosome number
Poor HiC**

High chromosome number + poor HiC + no telo information



Symbiodinium – 92 chroms protist, no karyotype available or close species with genomic data available

Some approaches:

1. Adjust gamma contrast in Pretext



2. Higher possible contrast in Pretext colours, dark background

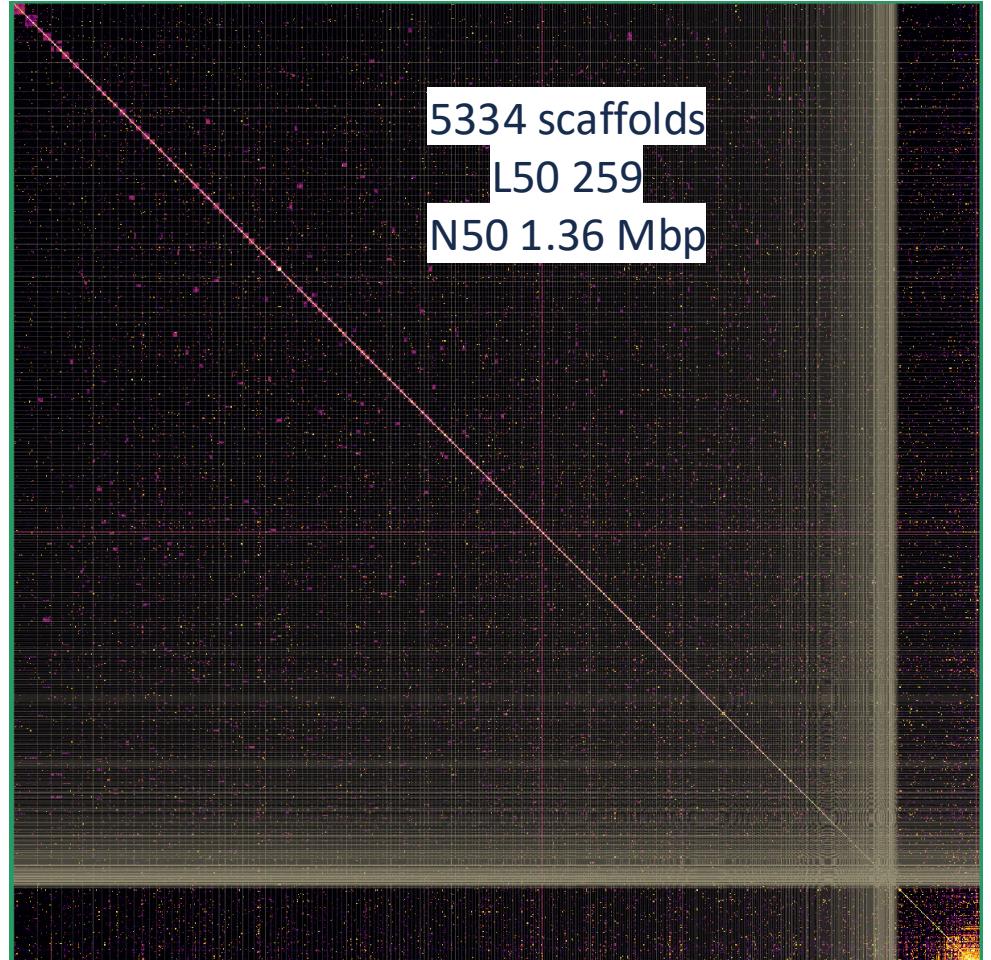
3. Normal resolution maps

4. Zoom in

5. Use a comparator (when available)

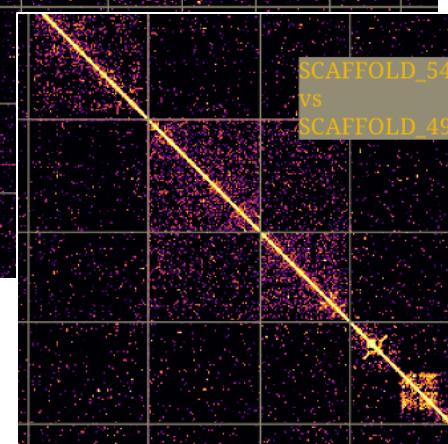
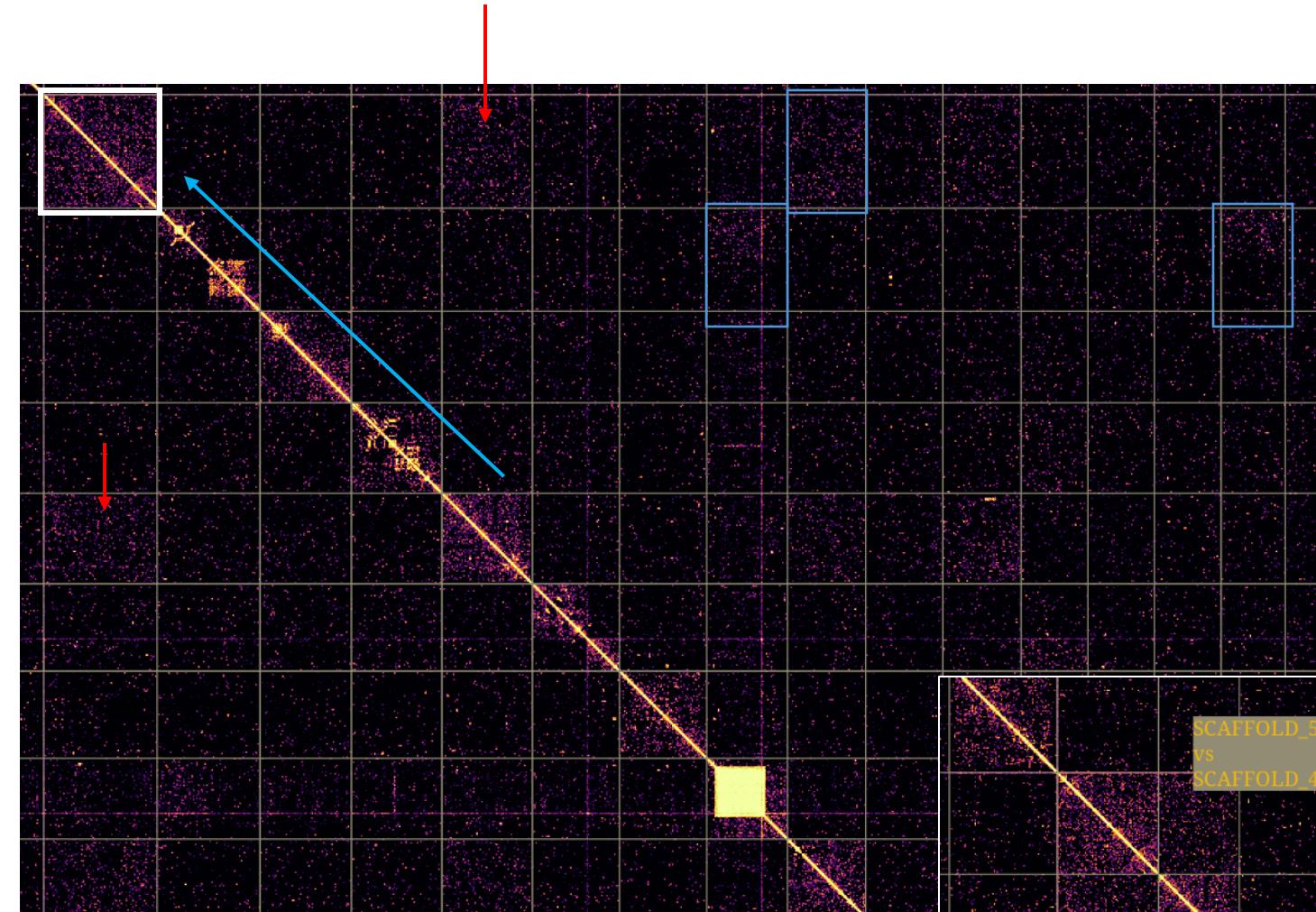
6. Telomere track

7. Top up

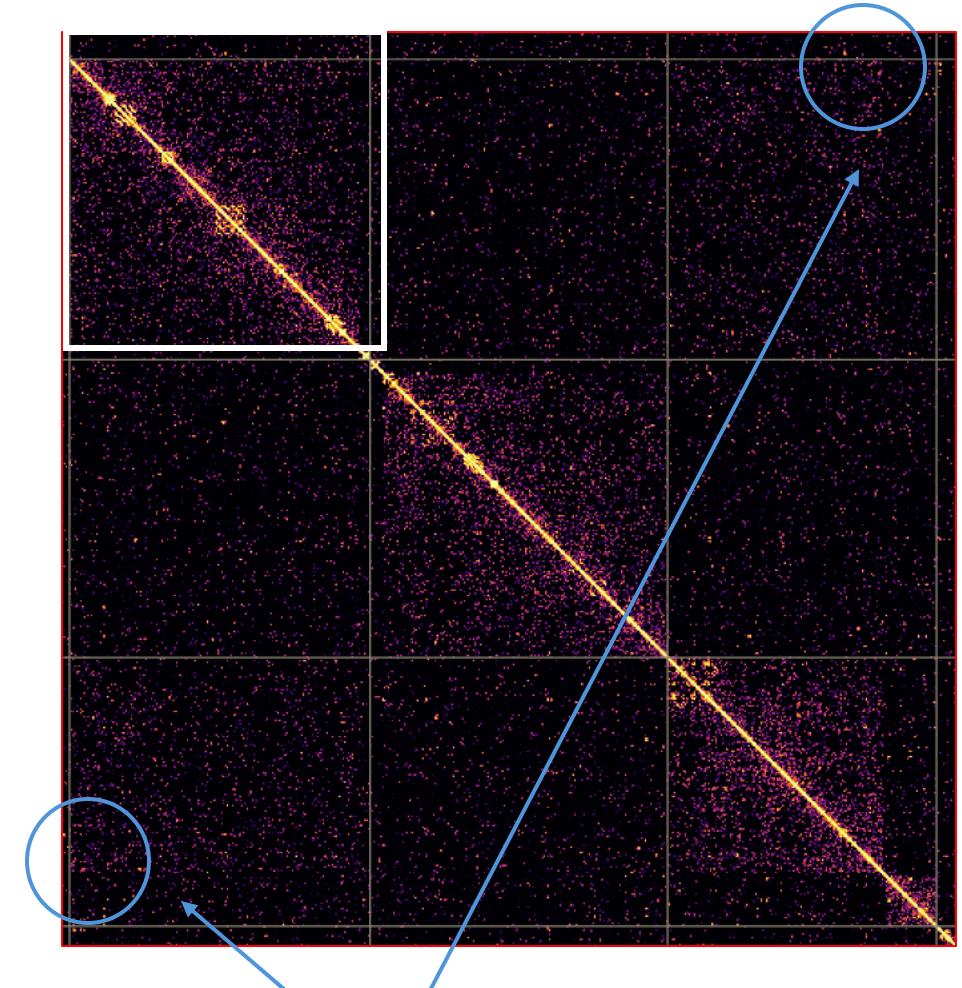




High chromosome number + Bad HiC + no telo information



After zoom in

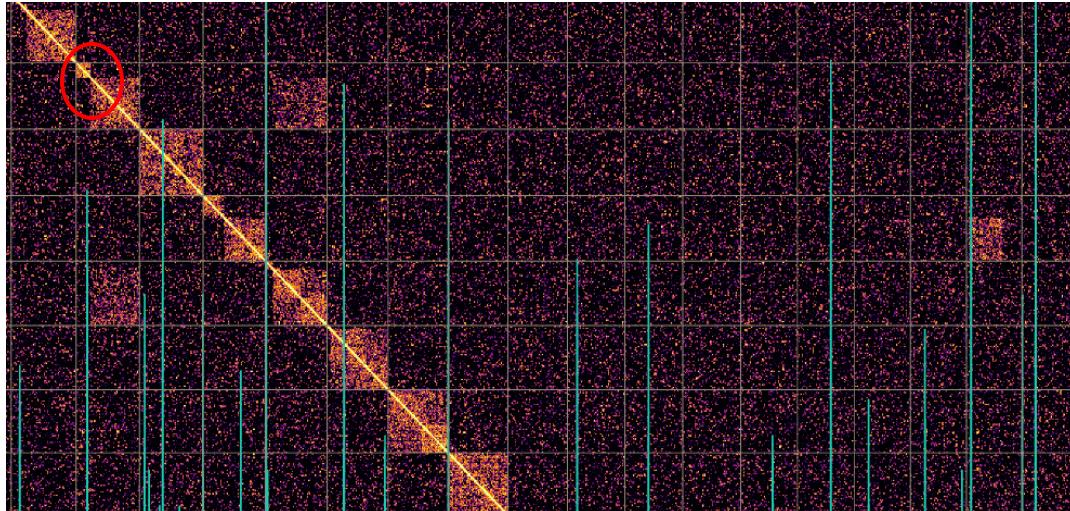


Strongest HiC signal

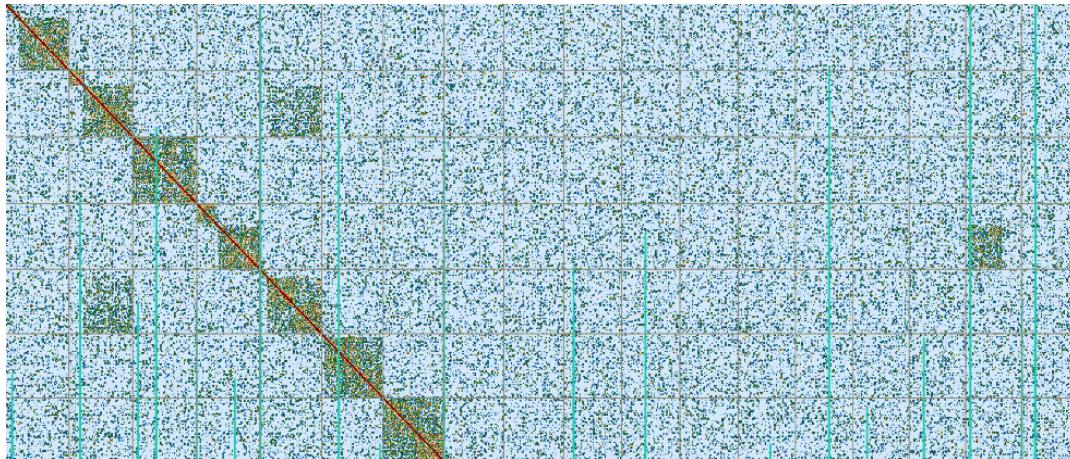


High chromosome number + Bad HiC + no telo information

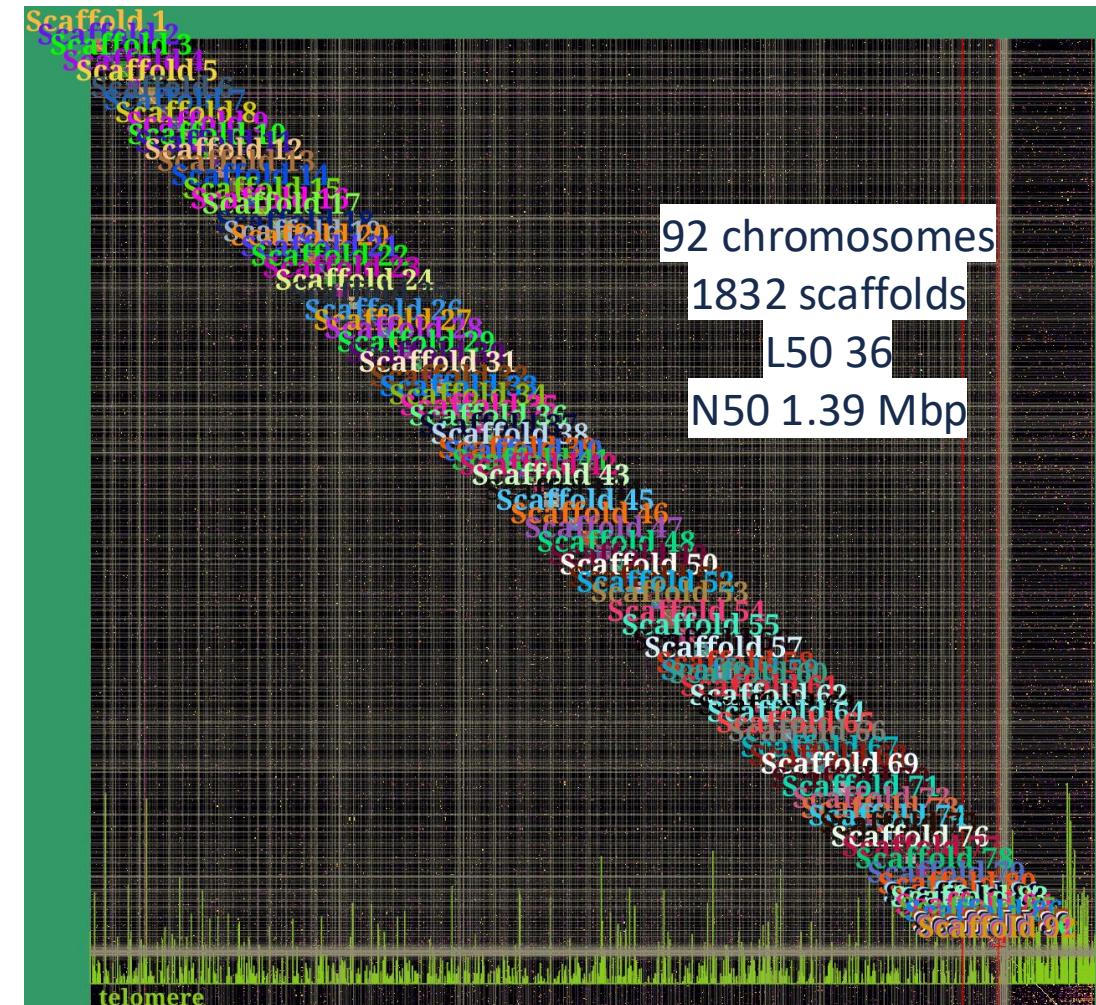
Gap track may also be helpful for breaks and cuts



Difference in the HiC signal visualization



Main source: HiC signal
High contrast colours sets (darker background helps)





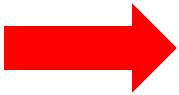
Haplotype phasing

Haplotype phasing

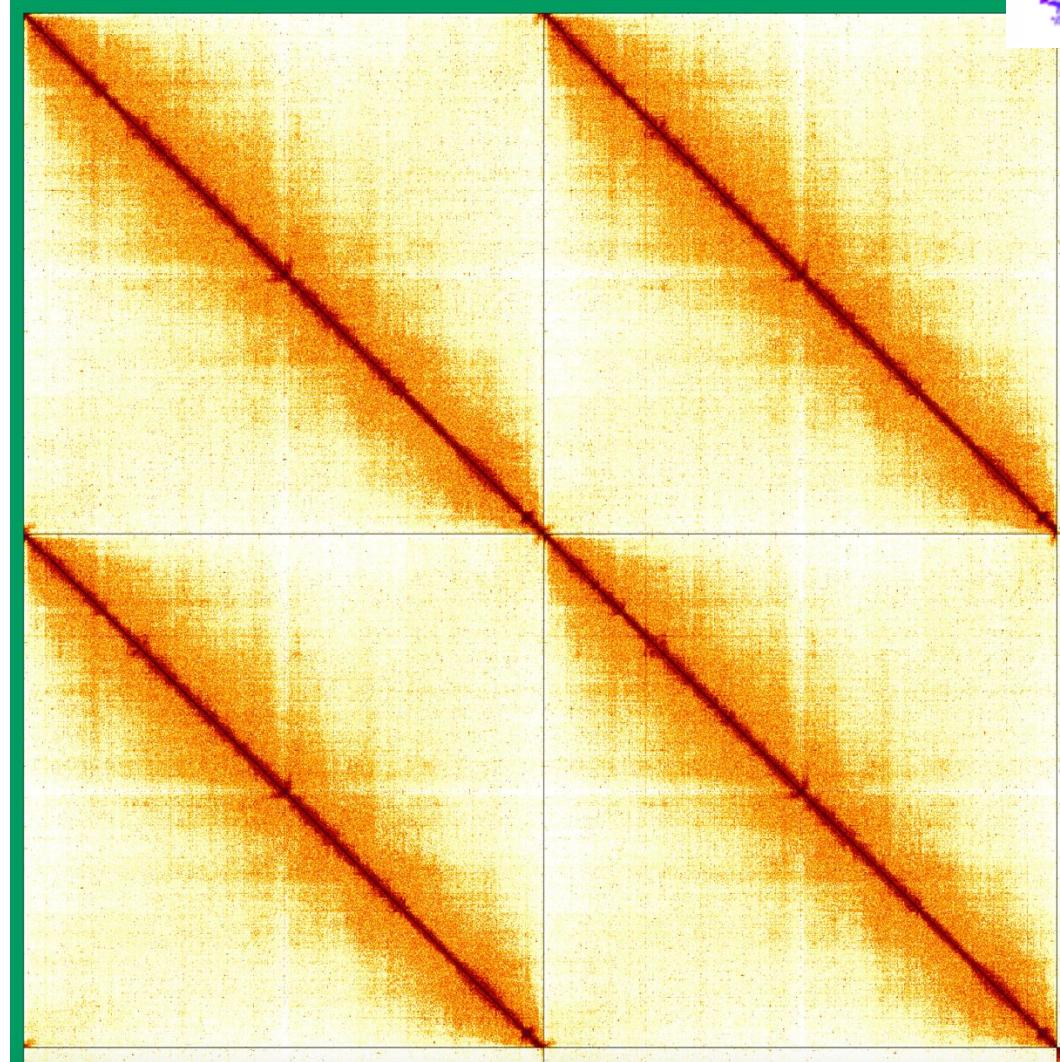


Good phasing

This is how the map should look
like when phasing worked well

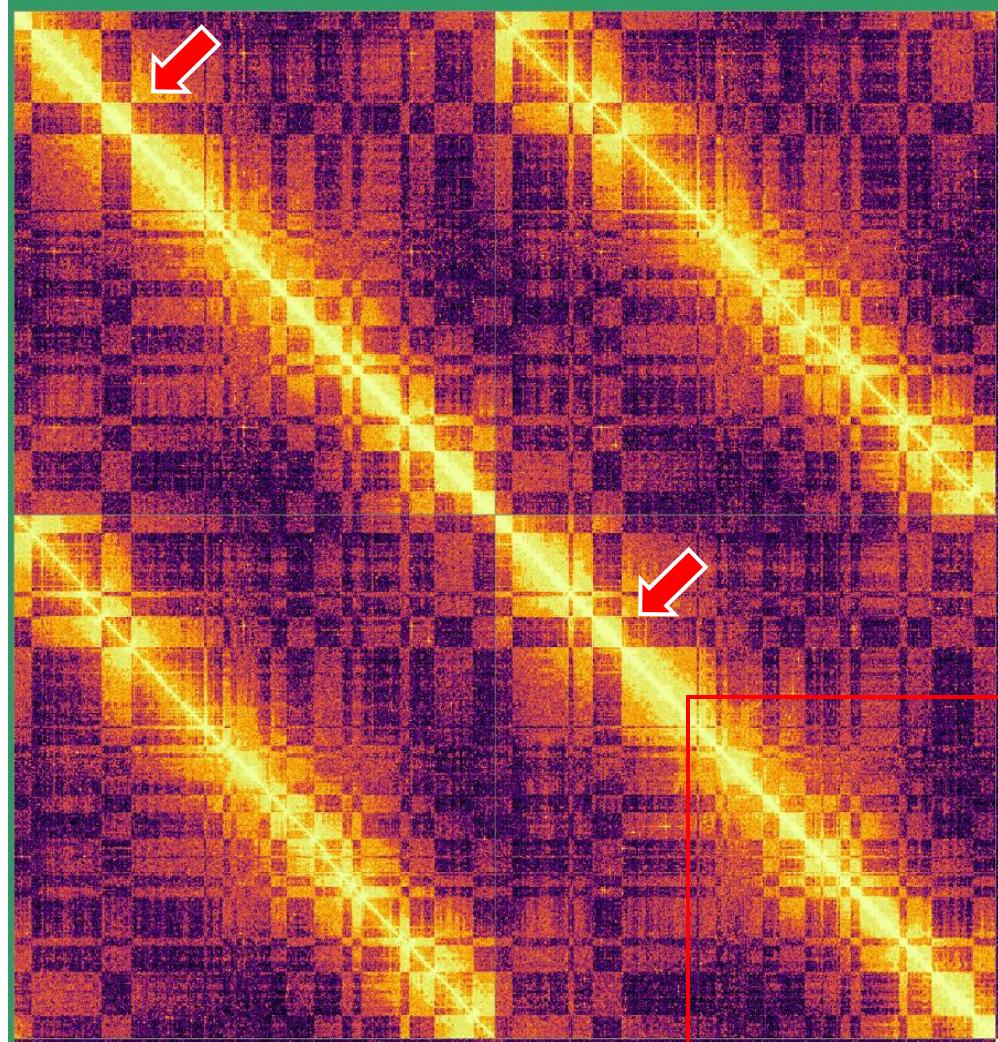


Contiguous HiC signal
No blanks or weak signal regions



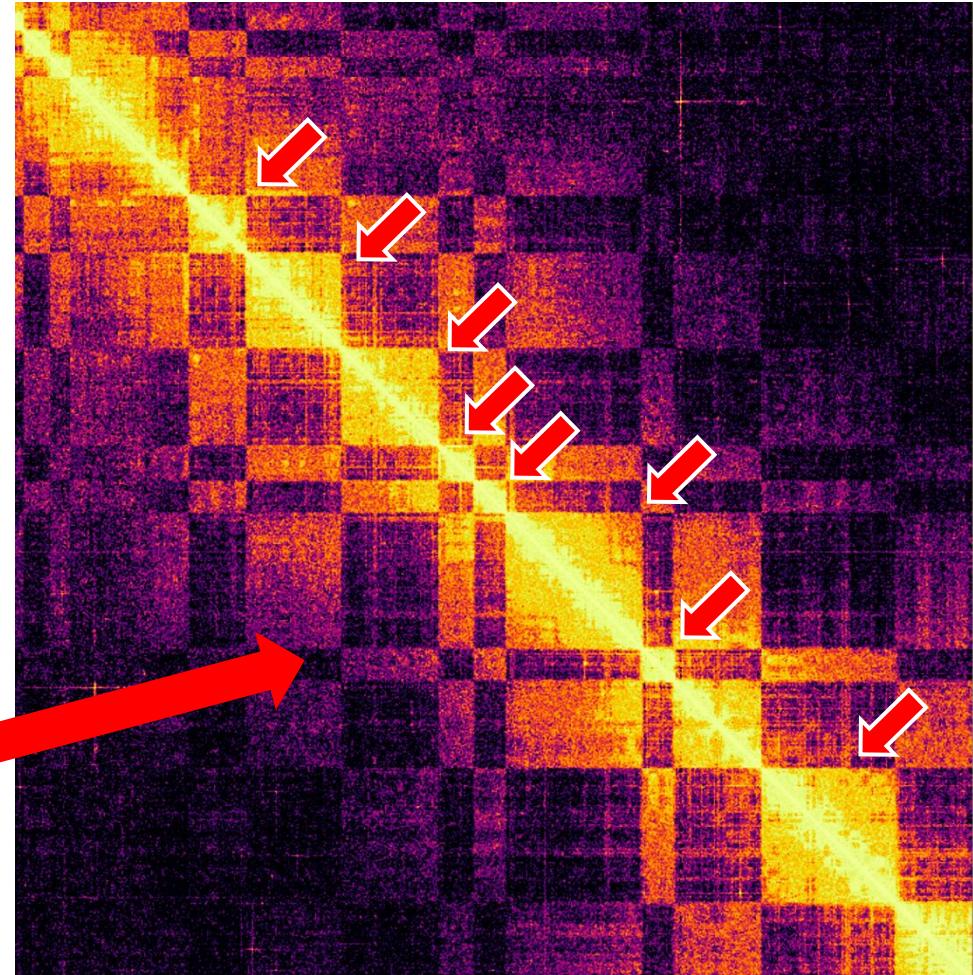
dcOxyDigi1

Haplotype bad phasing



xbMysUnda1

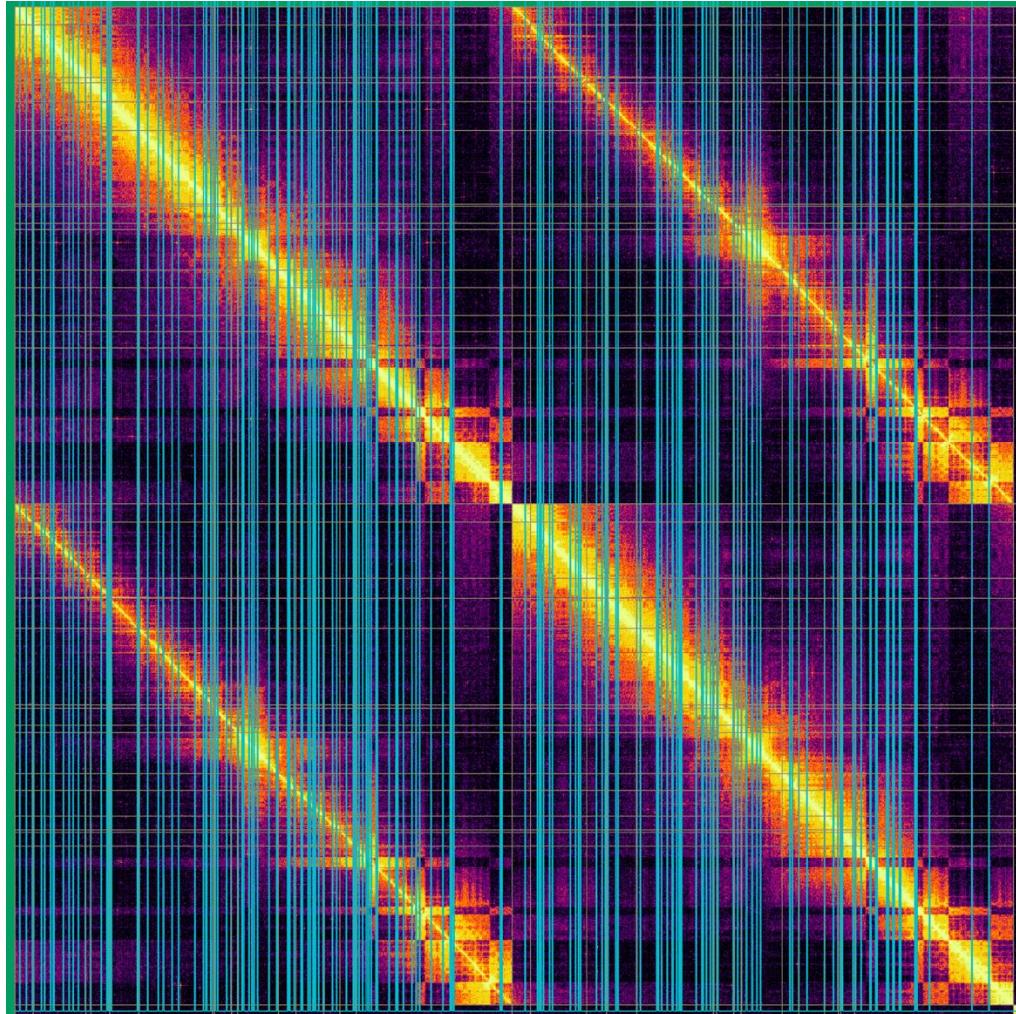
But when it doesn't work it is a source of confusion...



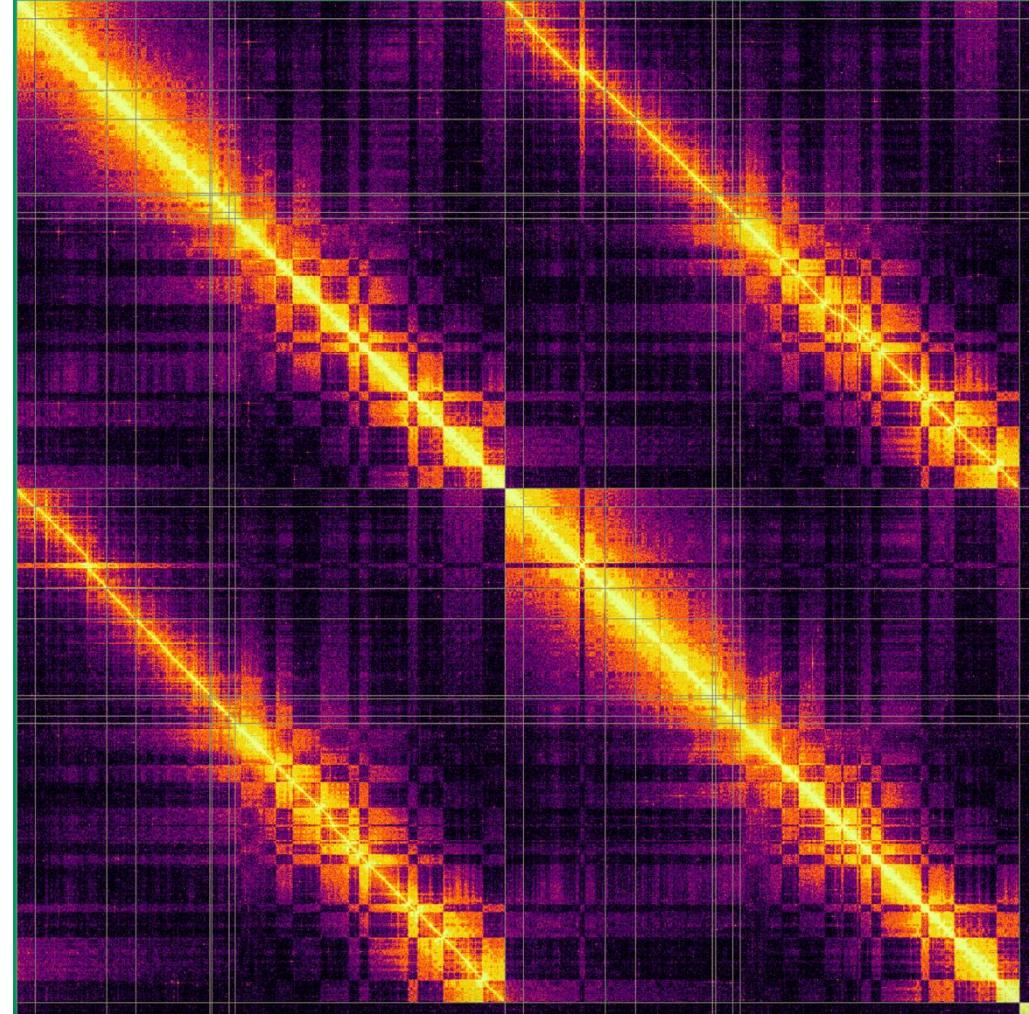
Haplotype bad phasing



Gap track on



Swap bits between the haplotypes

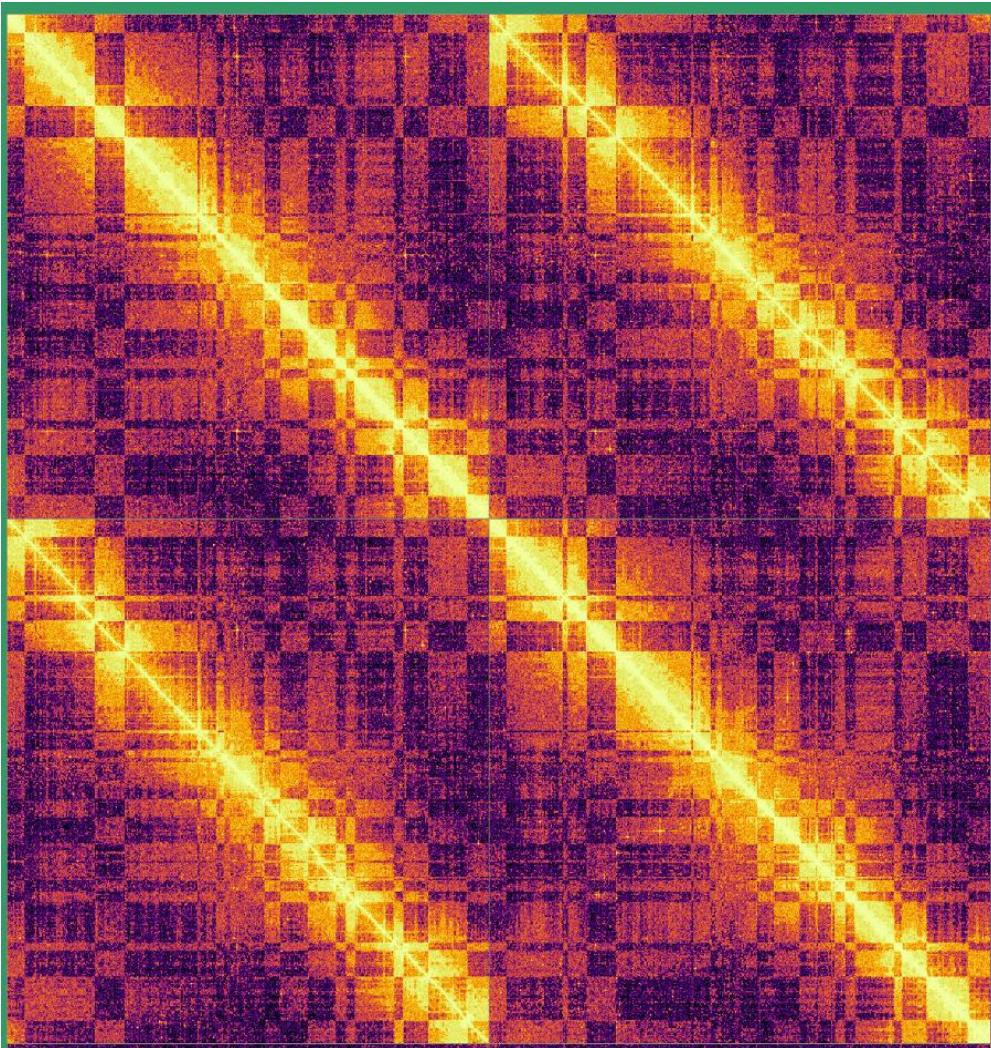


Manual phasing in progress

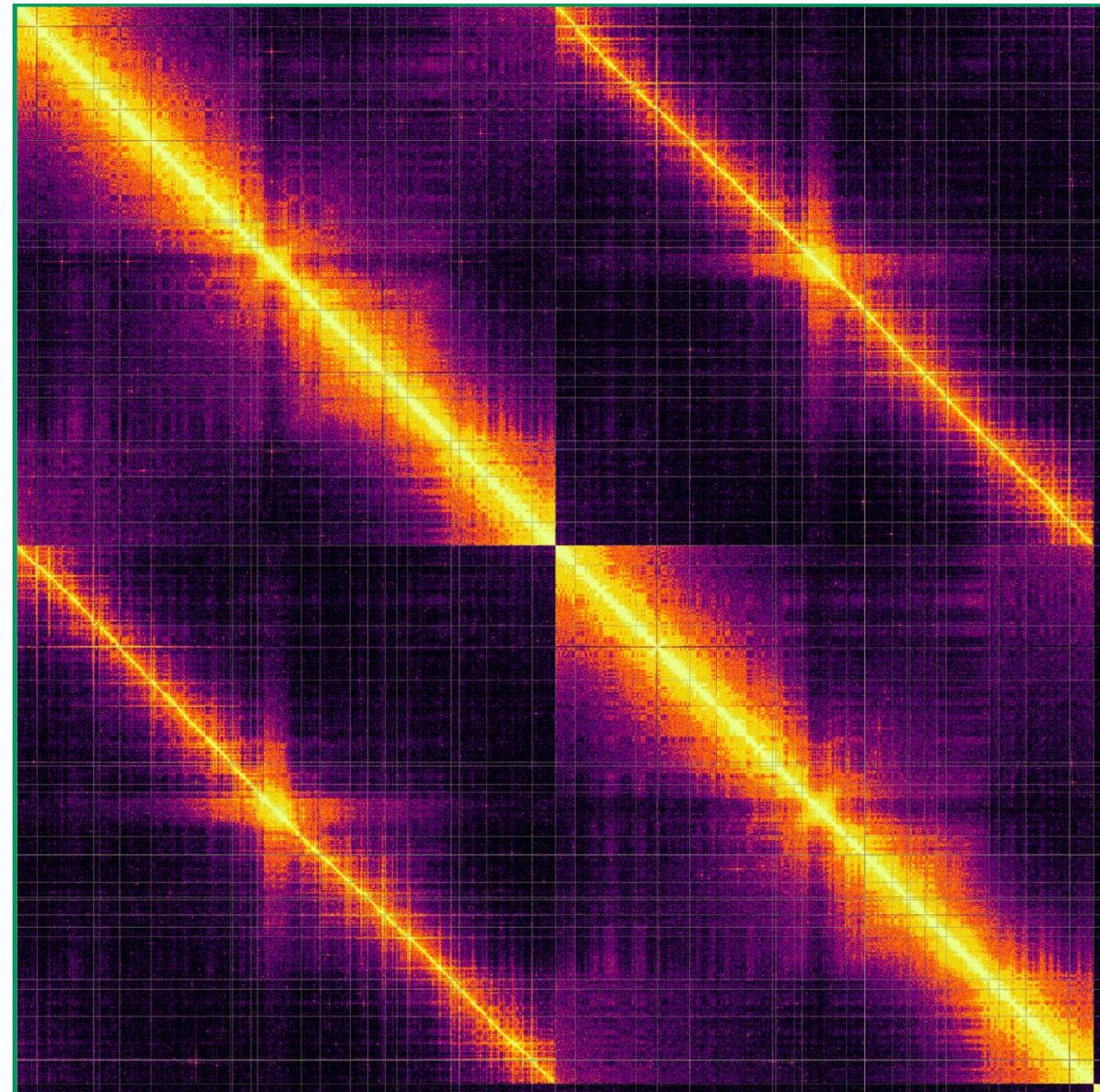
Pixel Cut in Pretext may be helpful



Haplotype bad phasing



After manual phasing





All haplotypes assembly curation

Standard Pipeline Assembly

By Dominic Absolon

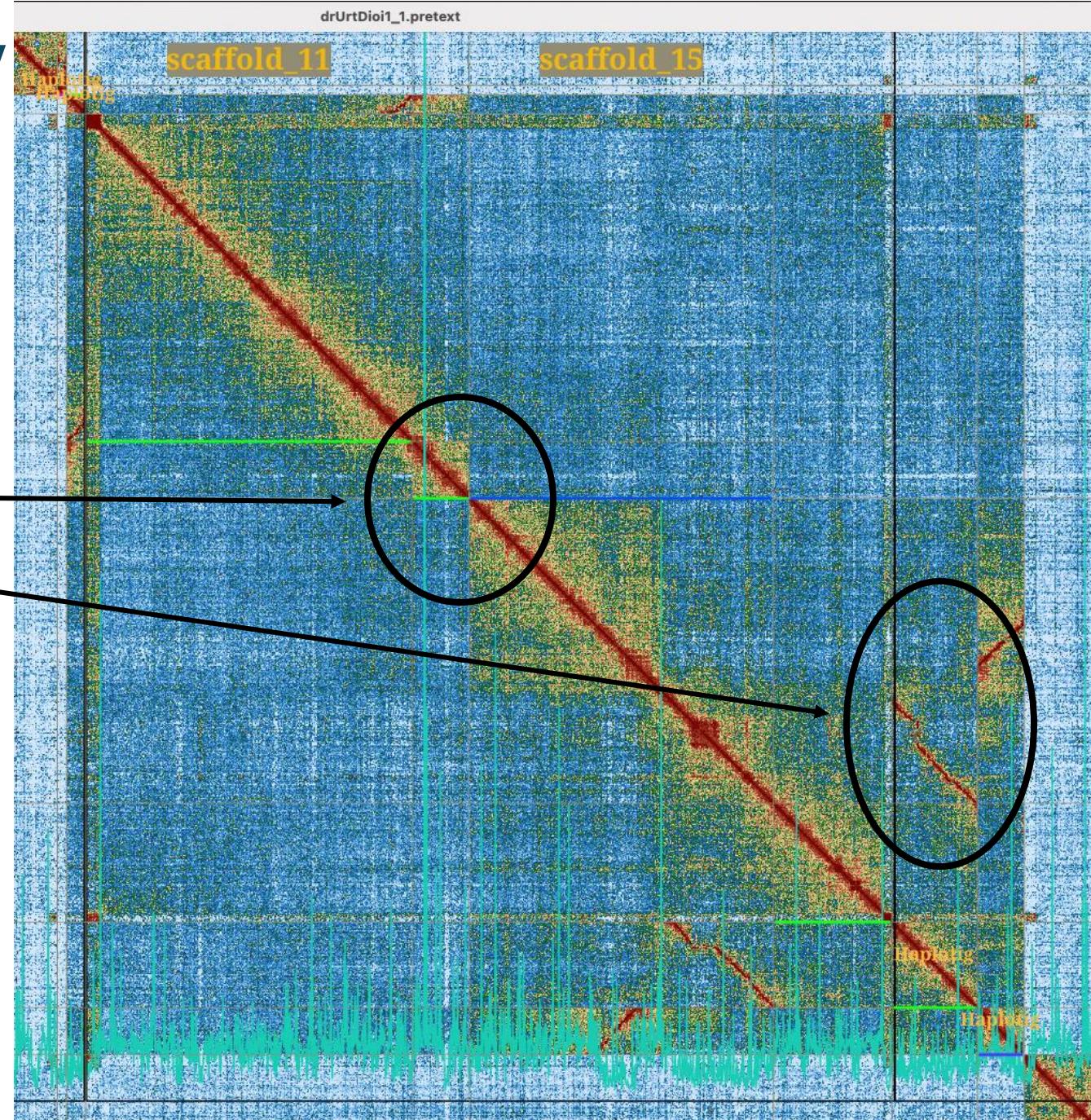
drUrtDioi1 – tetraploid

Initial “primary” assembly had issues:

- Missing sequence
- Over-represented sequences

Primary assembly:

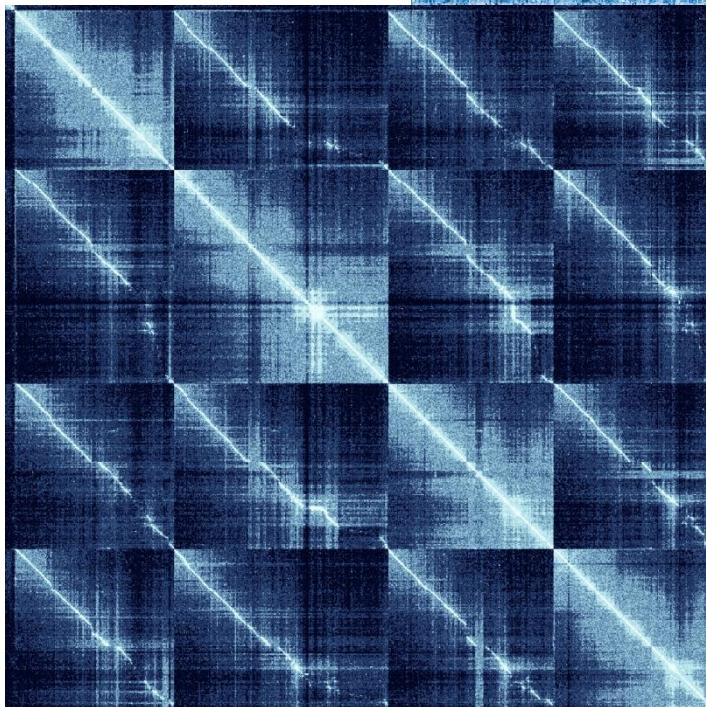
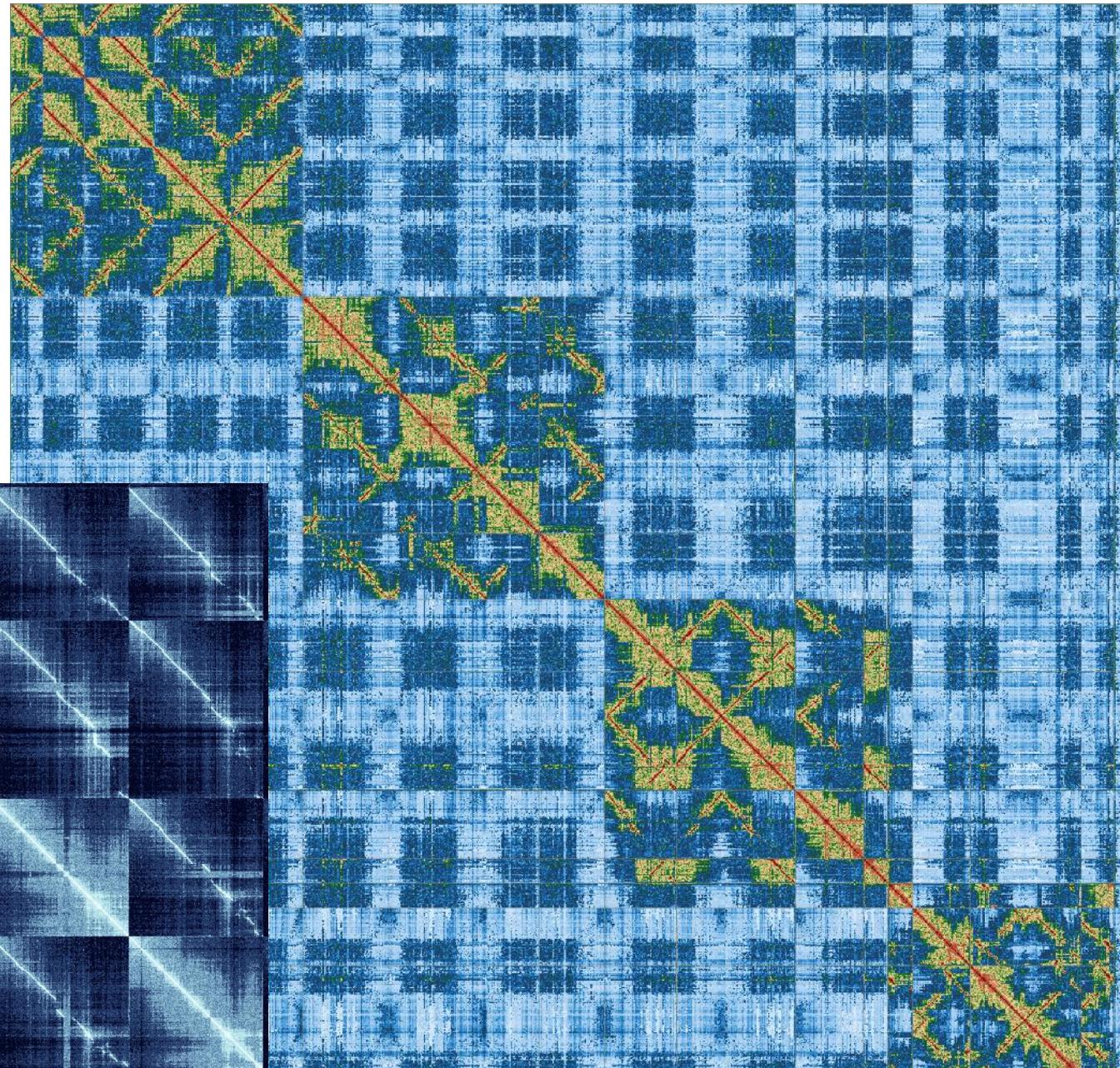
hifiasm (w/ purging) + purge_dups + hicmapping +
yahs



All haplotype map

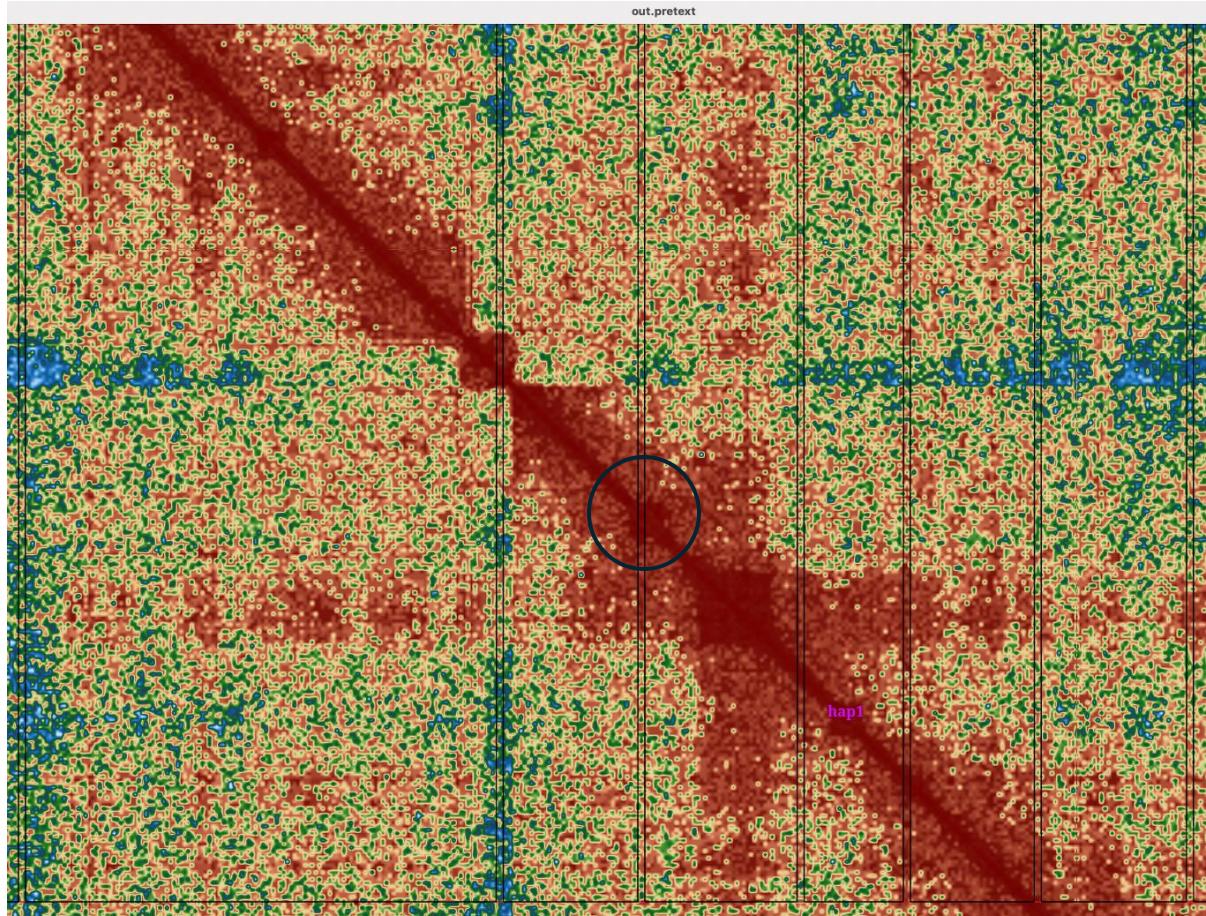
By Dominic Absolon

All haplotype assembly



All haplotype X Single haplotype maps

Curating – all 4 haps – hi res

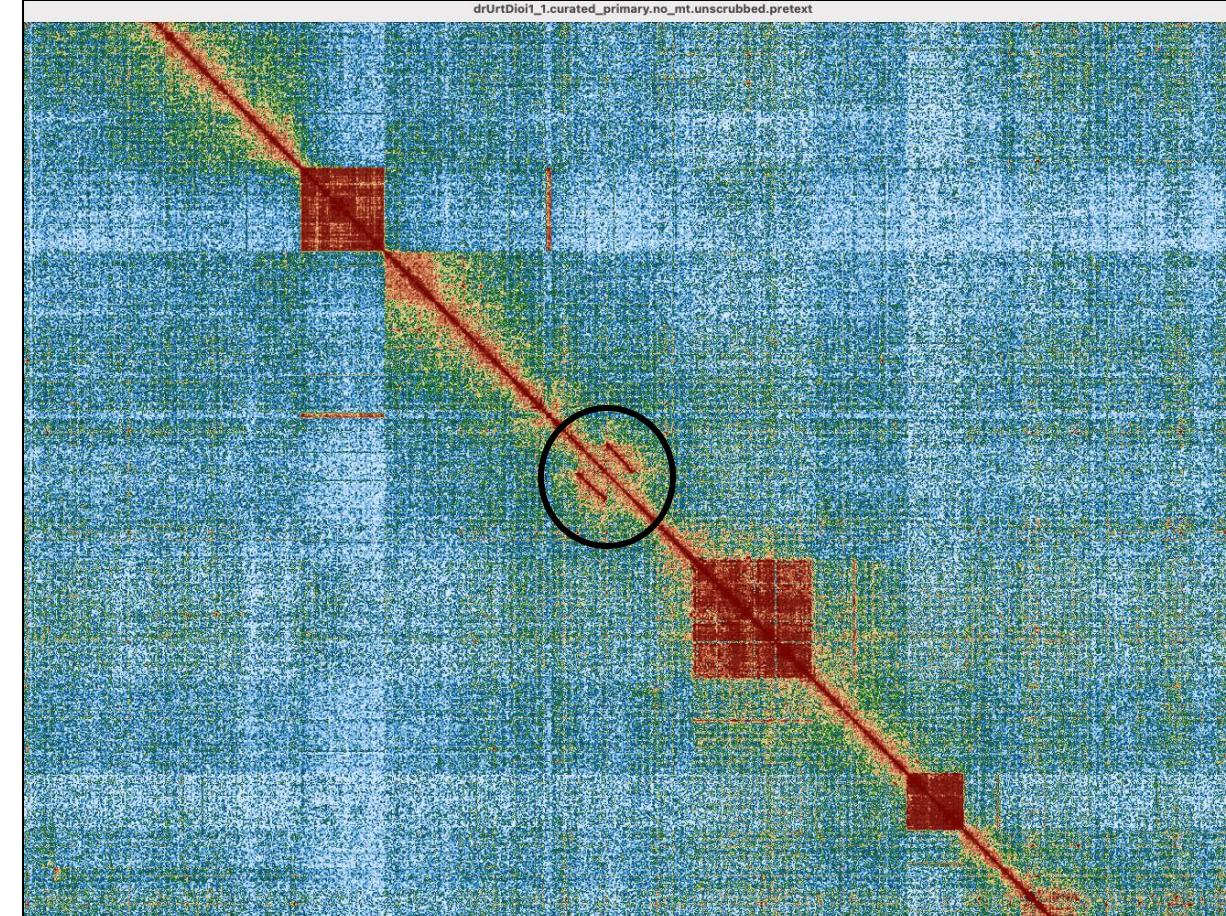


By Dominic Absolon



Resolution split among all haps

Curated – 1 single hap - hi res



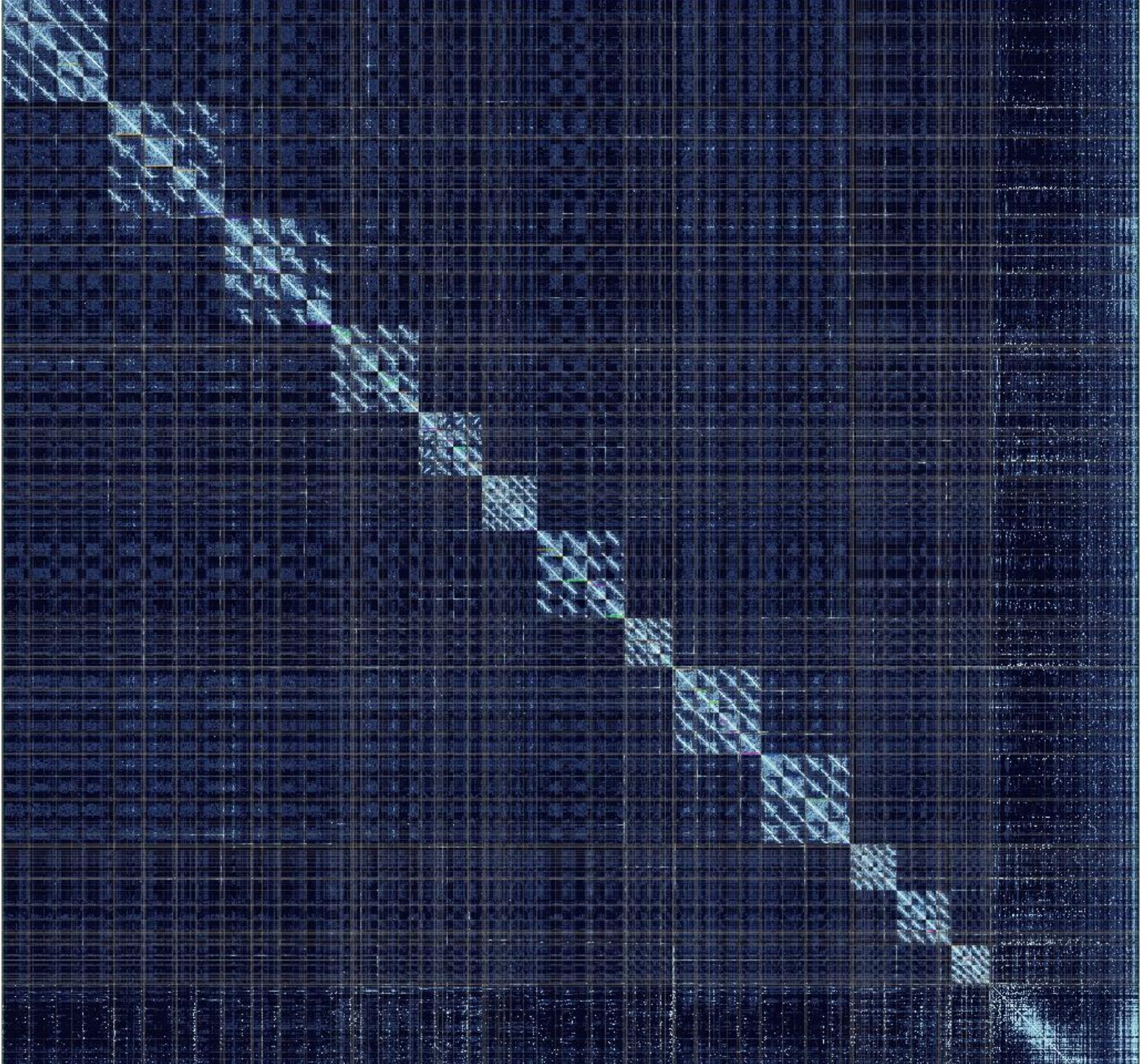
1 hap mapped against all HiC dataset

A working approach:

By Dominic Absolon

HAP1 file:
Chromosome-level curated HAP1

HAP2 file:
Scaffold-level
HAP2, HAP3 and HAP4





Polyplloid genomes

Polyplloid genomes



- Main issues:

1. Polymorphism among haplotypes
2. Translocations between regions of different chromosomes
3. Curated fasta mapping to the whole HiC dataset in the curated map

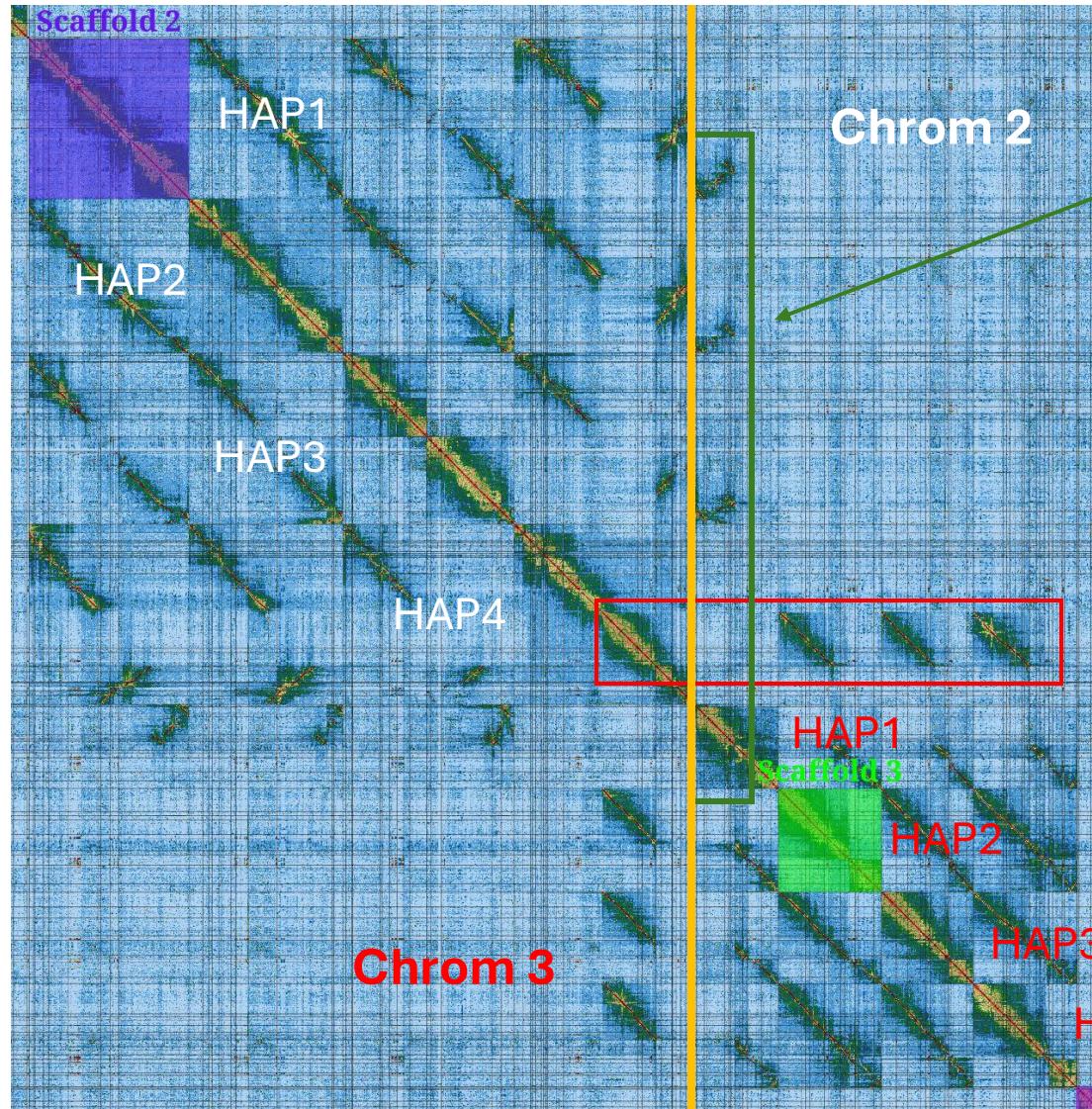


hap1 fasta only

Differences (look like errors) in the curated map

Polyplloid genomes

Translocations between regions of different chromosomes



Chrom 3 affinity with one region of chrom 2 in 3 haps
(except HAP4)

Chrom 2 affinity with one region of chrom 3 in 3 haps
(except HAP1)

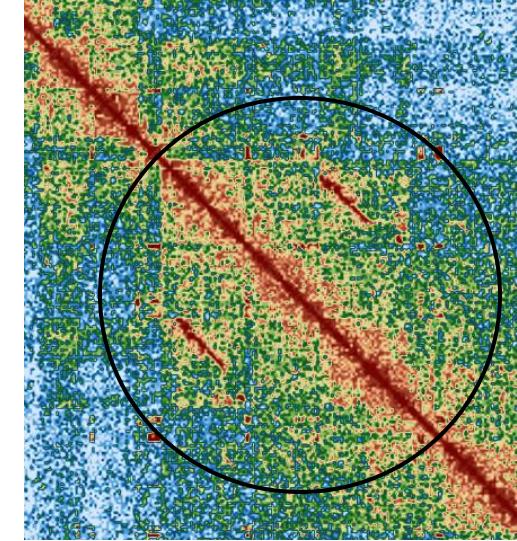
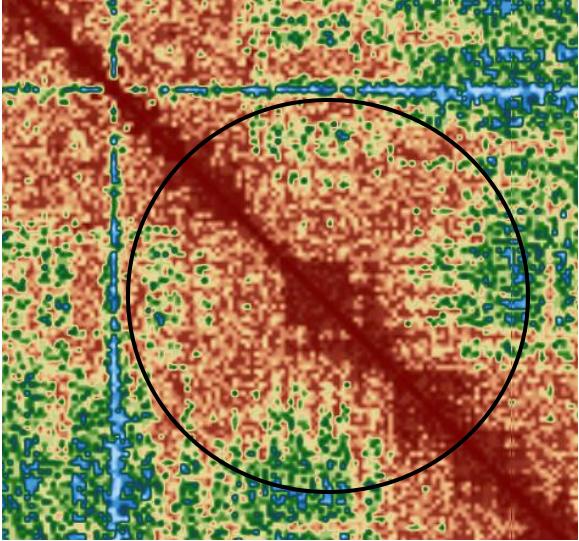
Only 1 hap will be in the final curated fasta
Translocations will not be shown

wgTheLage1

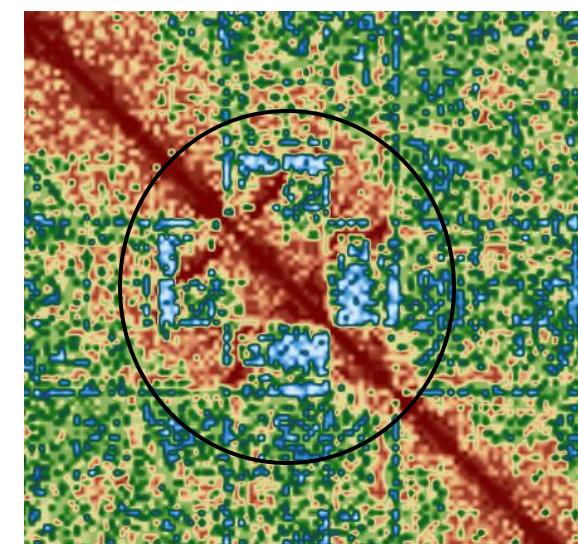
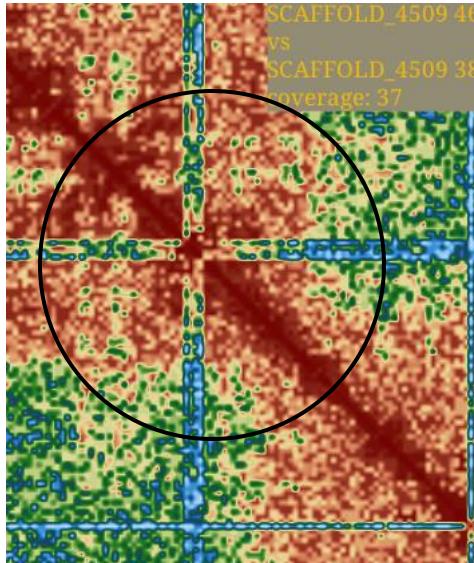
Polyplloid genomes

All set of HiC data mapping to one hap fasta only. Differences (look like errors) on the curated map

All haplotypes map

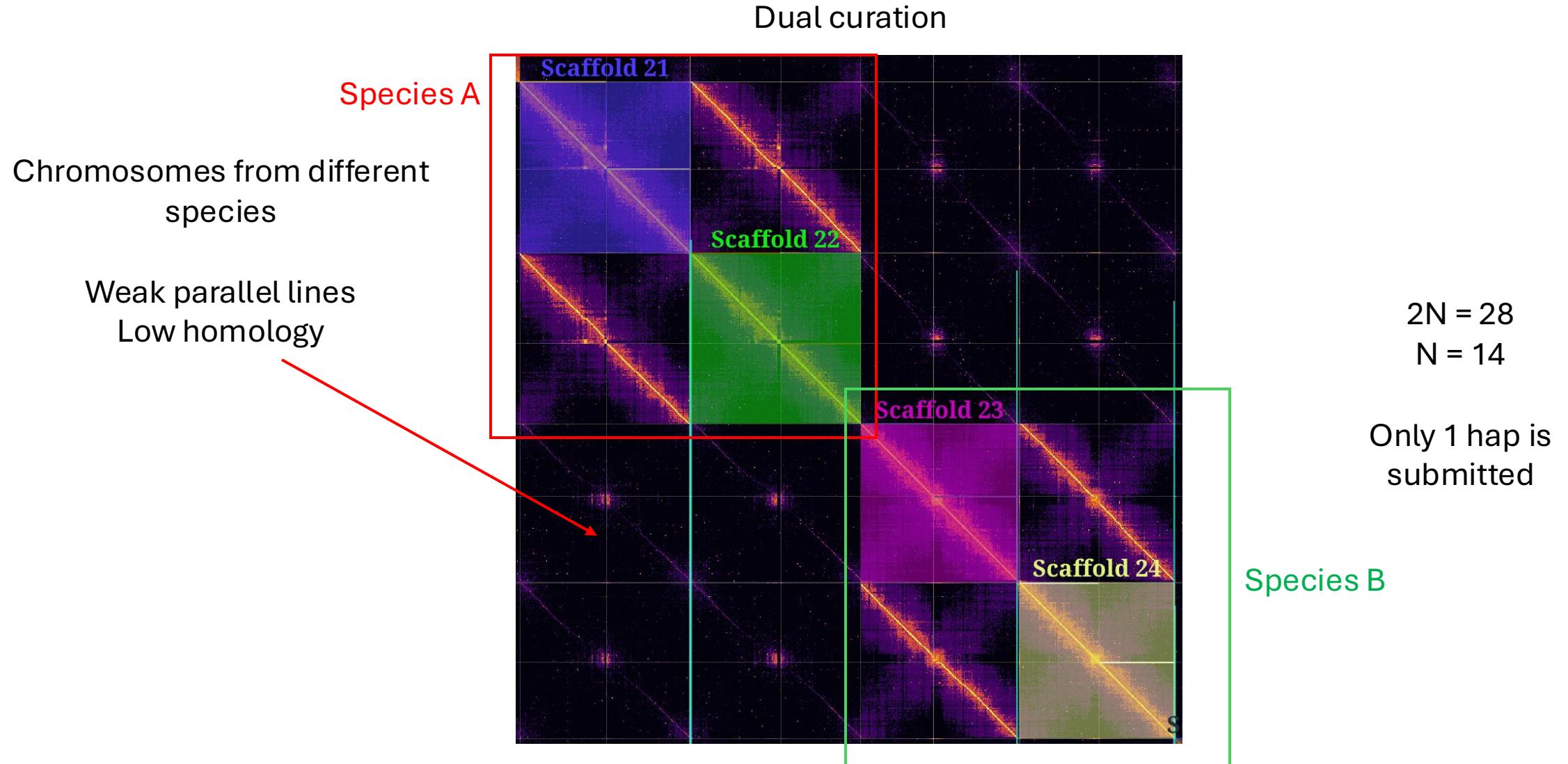


One hap only
curated map



Allopolyploid genomes

Two or more complete sets of chromosomes from different species

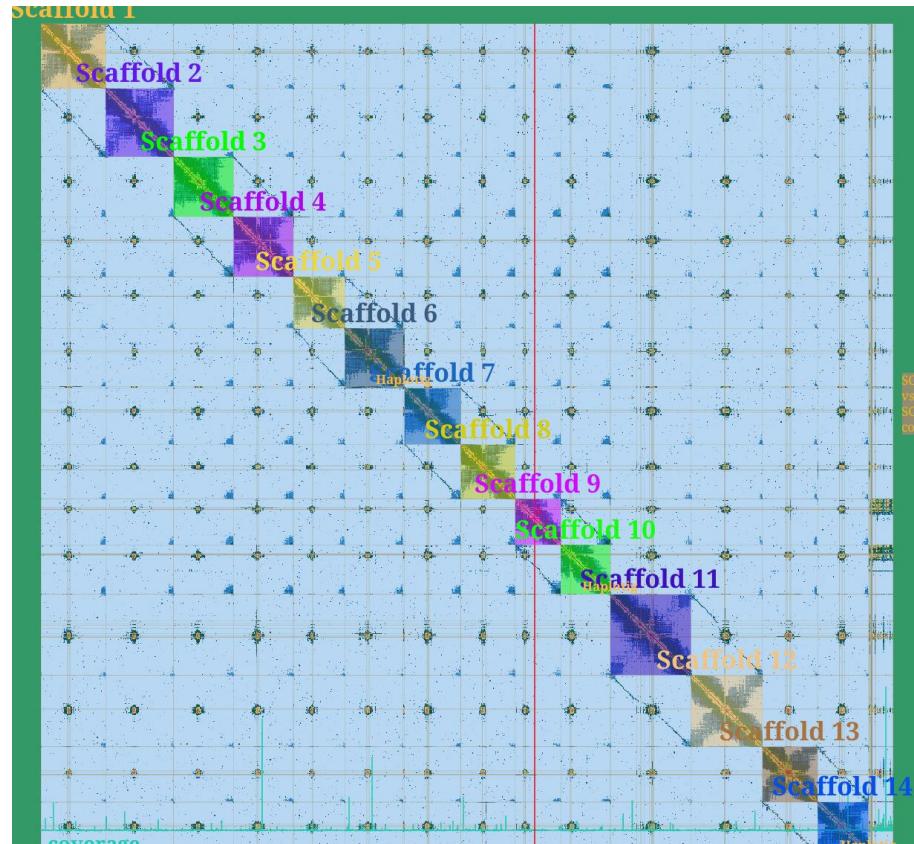


Allopolyploid genomes

Allopolyploid genomes

Single hap curation

N = 14



Curated haploid map



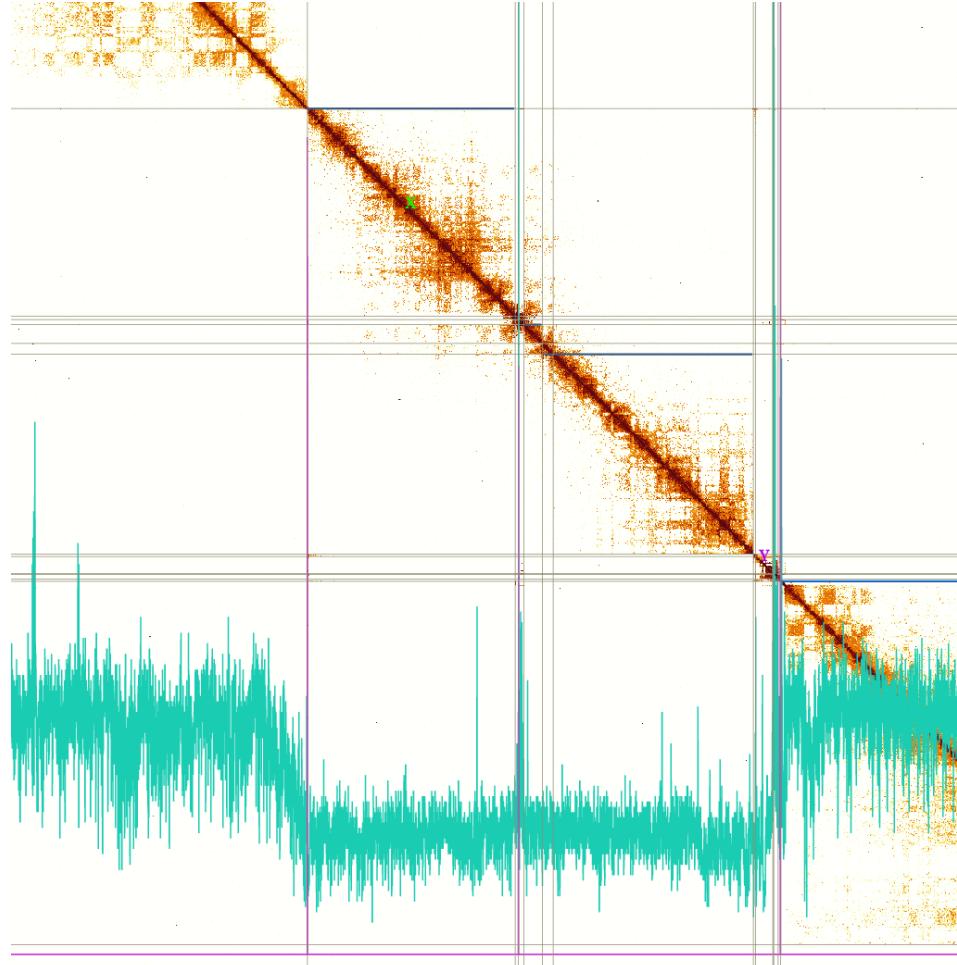


Sex chromosomes

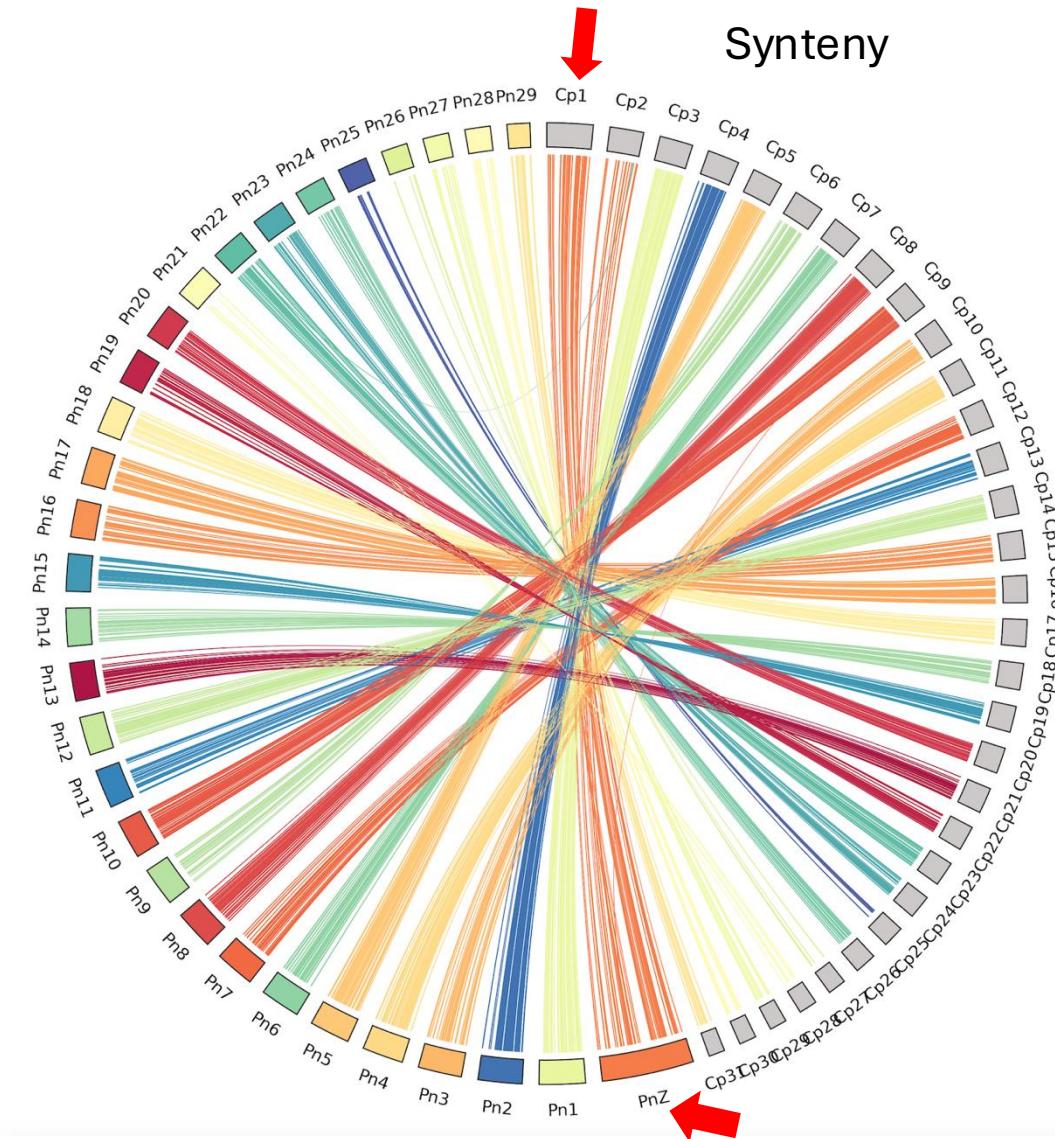
How do we usually identify/assign sex chroms?



PacBio read coverage

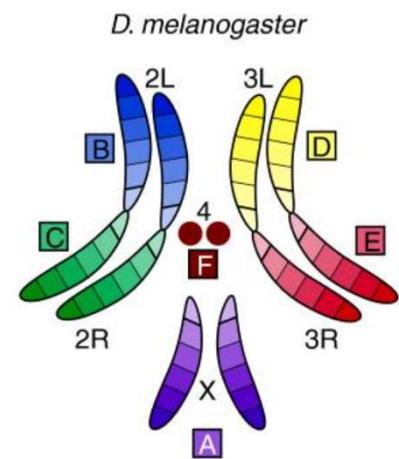
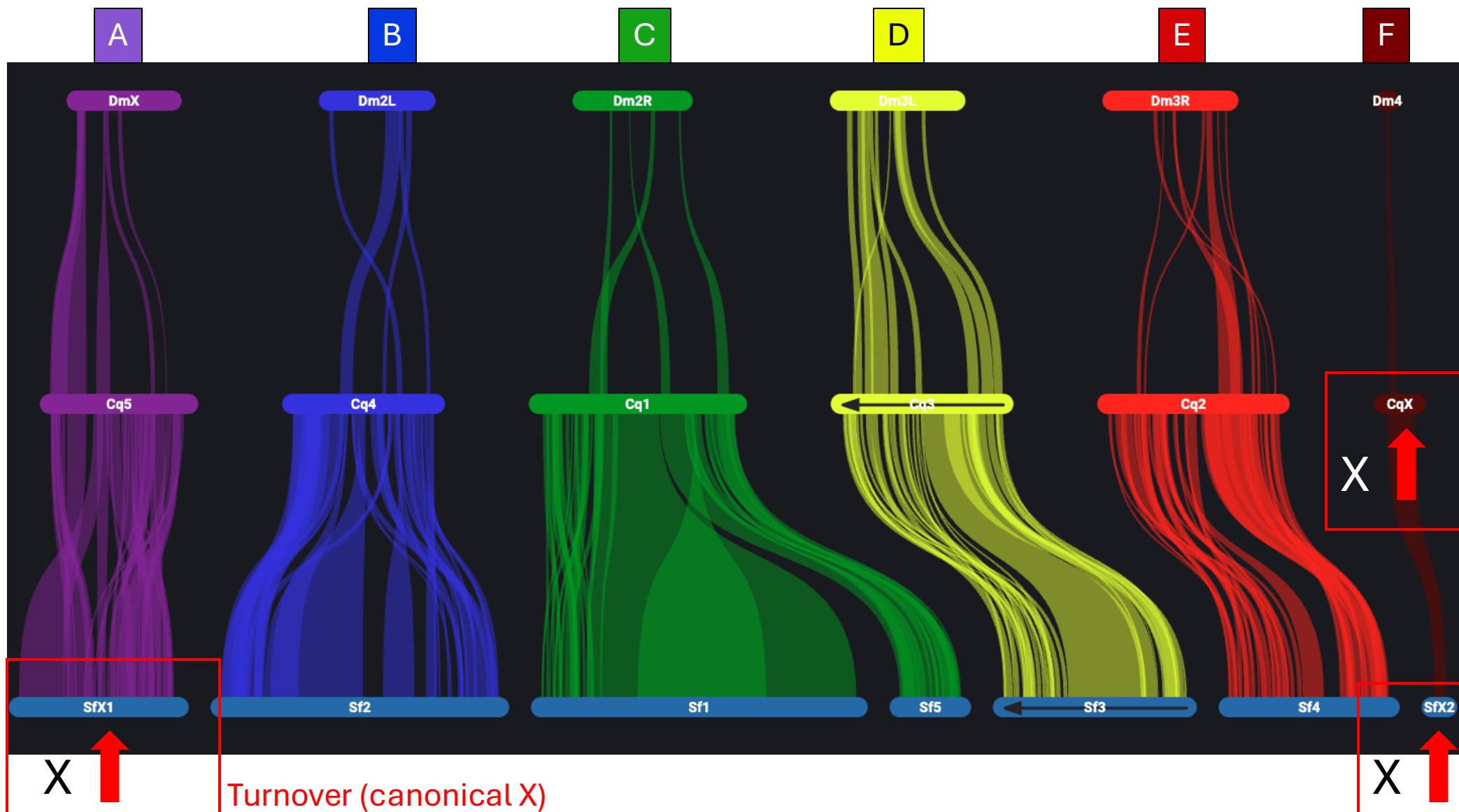


Synteny



Painting Muller elements onto Conopidae with *D. melanogaster*

Turnover is frequently observed in Diptera



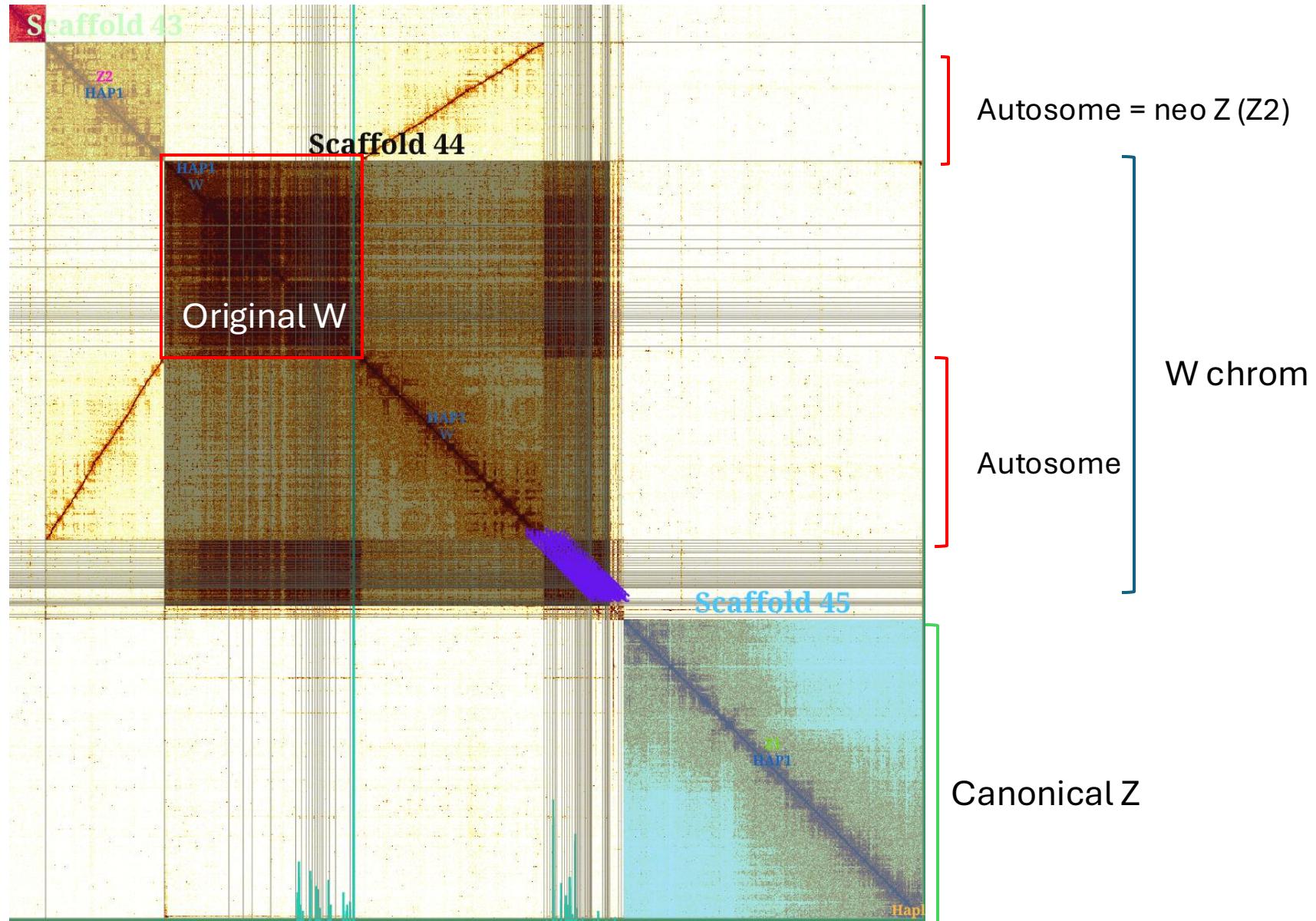
IdConQuad1

Turnover
(canonical X)

idSicFerr1
(version 2)

Turnover
(neo X) X2

Autosome + sex chrom fusion = neo sex chroms



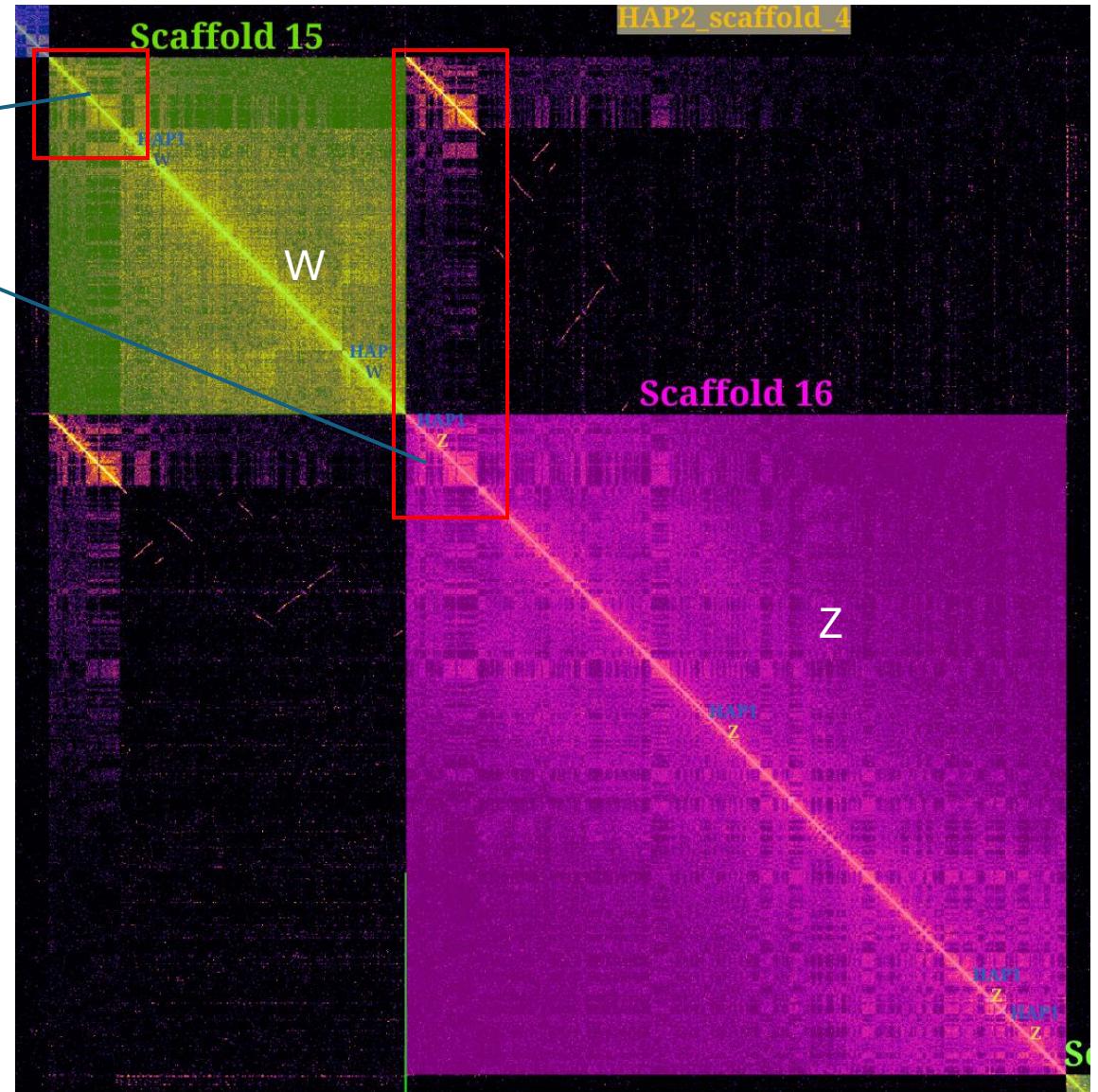
Homology regions between sex chromosomes = PAR



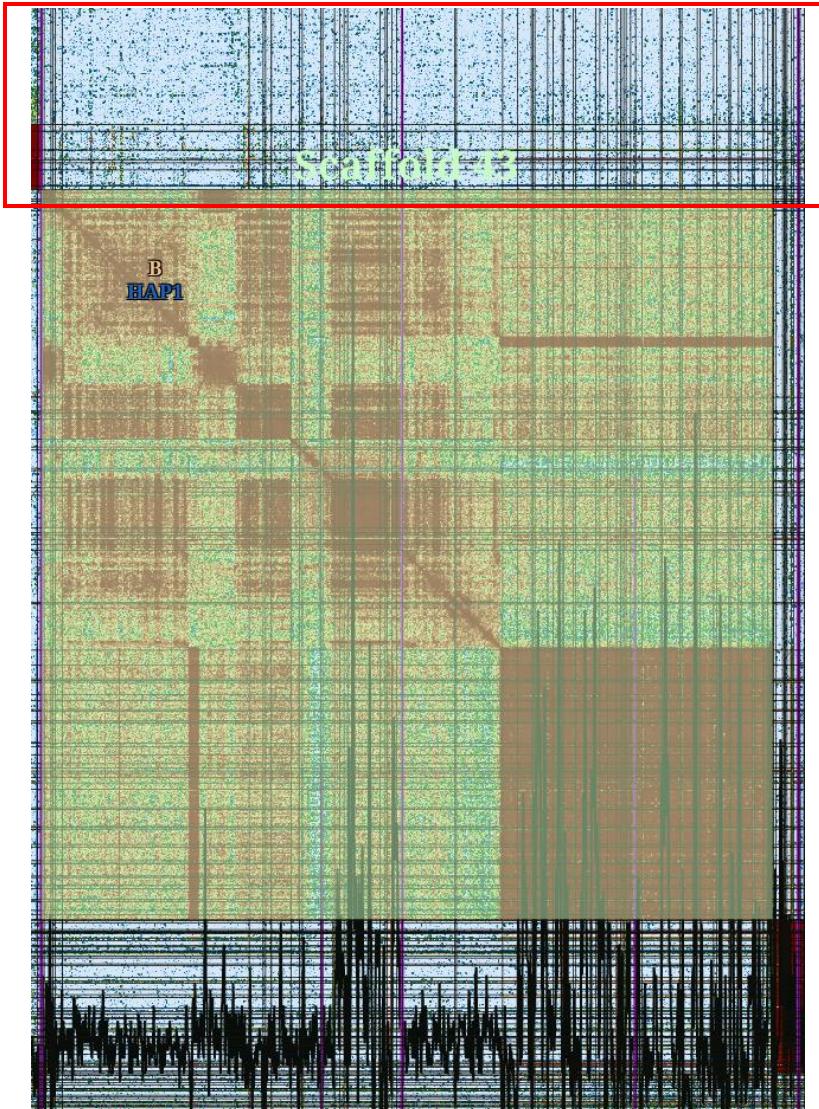
They should be assembled at
the same region/orientation
in both sex chromosomes

Used to self orientate sex
chroms

PARs (Pseudo-Autosomal-Regions) are highly similar sequences usually found at one end of a sex-chromosome pair. The rest of the sex chromosomes are typically highly diverged.

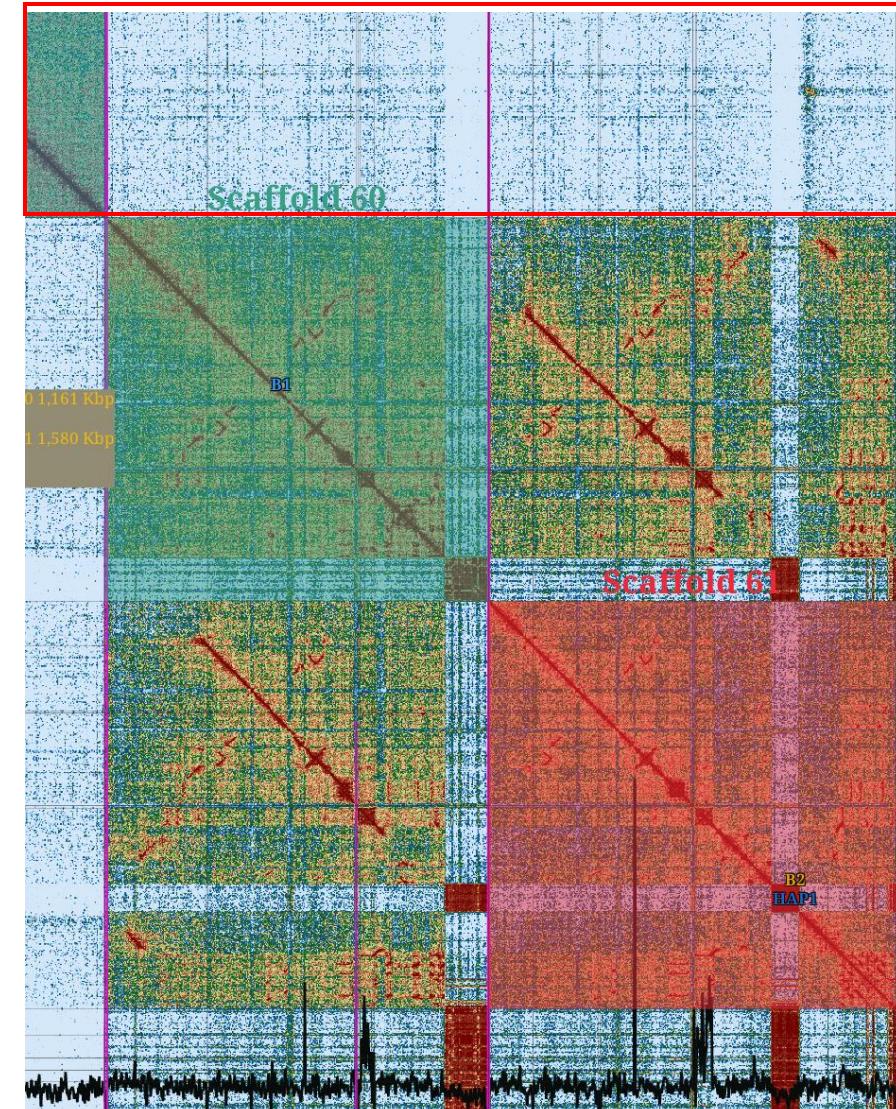


B chromosomes



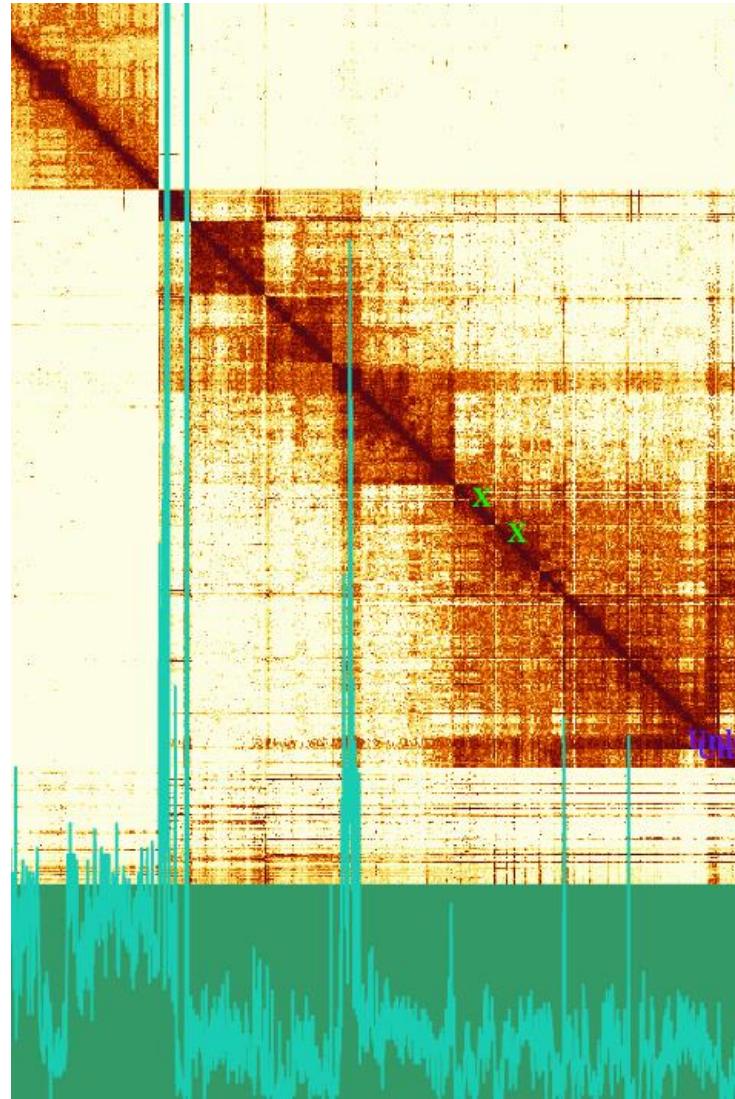
Duplicated BUSCO
hits only

HiC background with
the autosomes is
denser



HiC data may be mandatory to identify sex chromosomes

PacBio from male sample
HiC from unknown sex sample
(potential male)

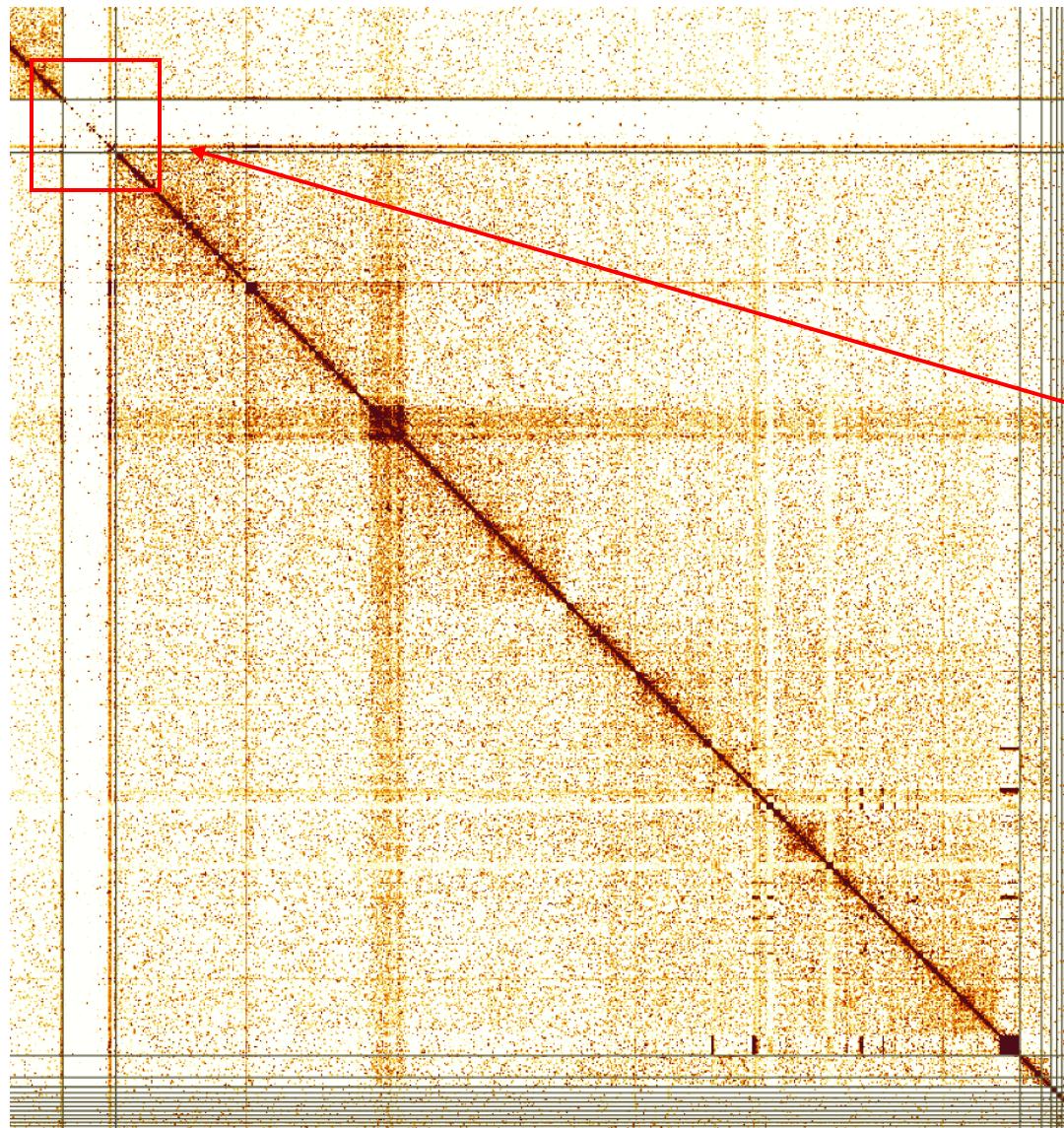


X is half coverage

Where is the Y chromosome?

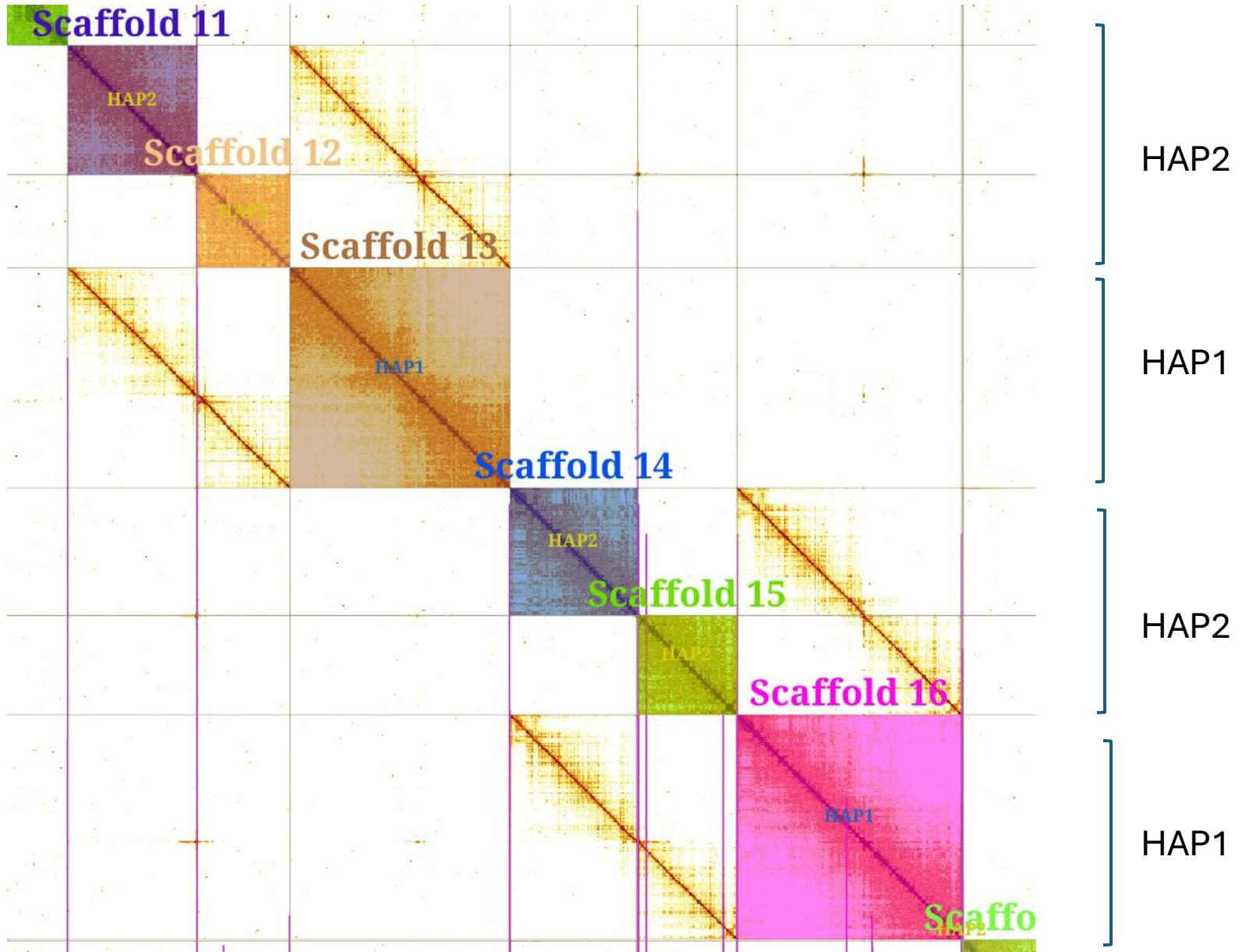
HiC data may be mandatory to identify sex chromosomes

HiC from a female sample
No HiC signal for Y



Here is the Y
chromosome!

Heterozygous fissions/fusions



iEreGorg1

Hands-on

<https://github.com/epaule/Physalia-Manual-Genome-Curation/blob/main/Session4.md>