

C G T A C G T A
A C G T A C G T

Exploring Satellite DNA & Tandem Repeats with ModDotPlot

Alex Sweeten

Biodiversity Genomics Academy
October 7th, 2024



@alexsweeten



@alexsweeten.
bsky.social



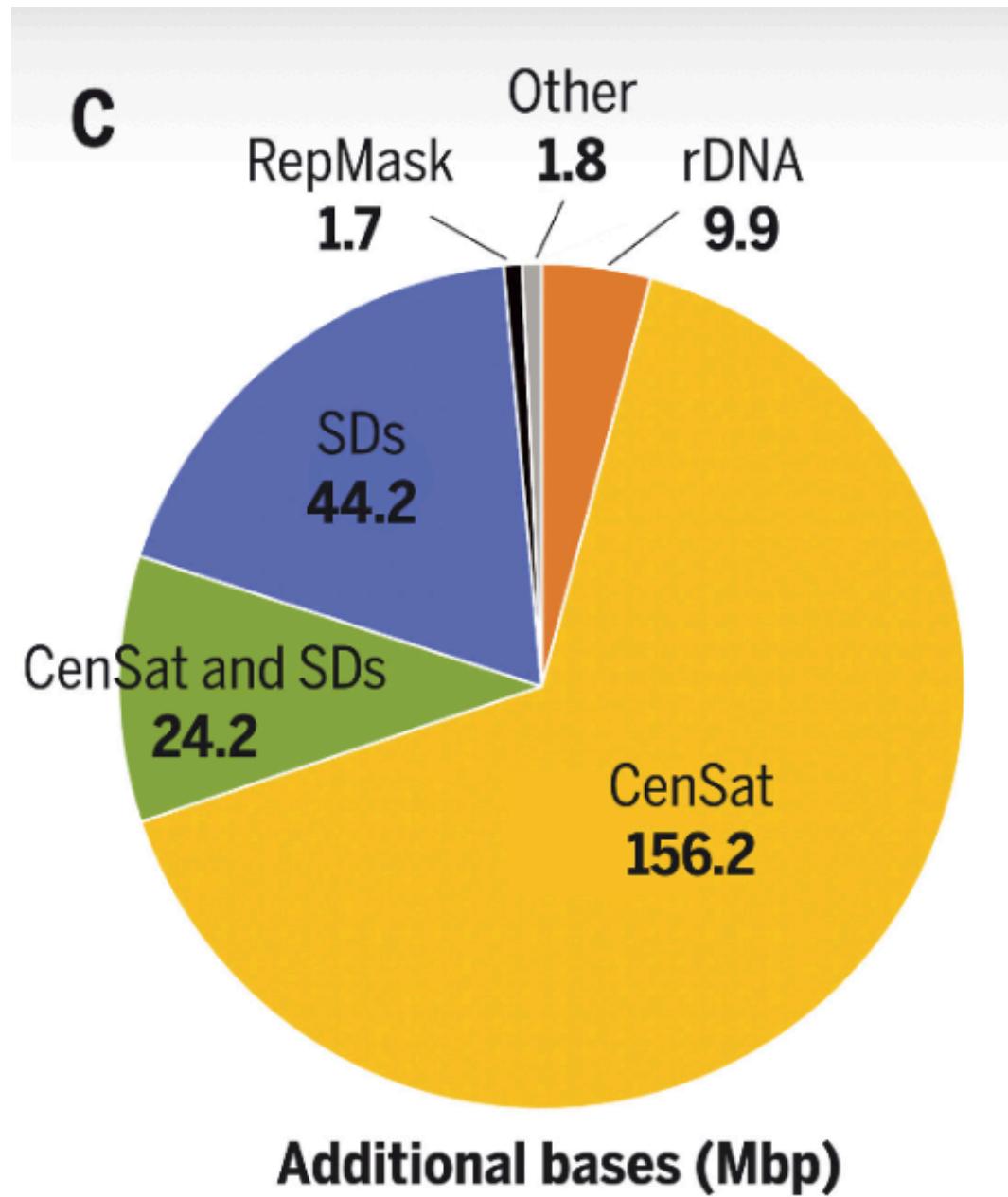
—
The **Forefront**
of **Genomics**
—



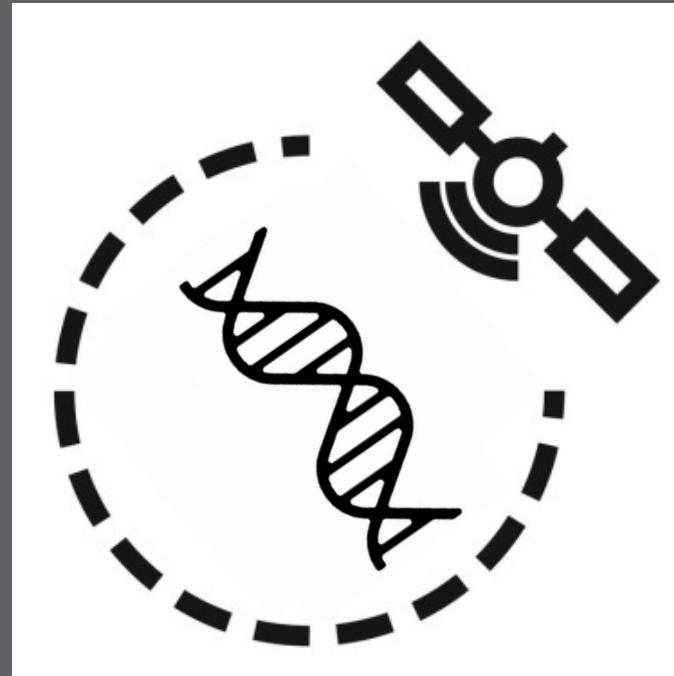
JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

A new era

- The first T2T human genome introduced over 200 Mbp of new sequence.
- Overwhelming majority of this are **satellite & segmentally duplicated sequences**.

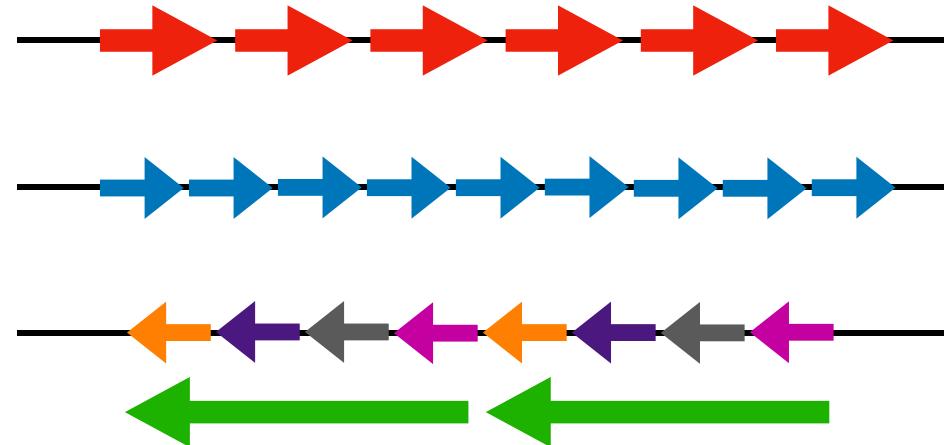
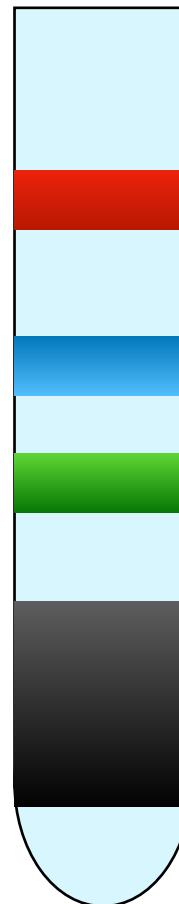


What is Satellite DNA?



Satellite DNA

- Consists of long arrays of **tandemly repeating** sequences.
- Primarily found in **centromeres** and **acrocentric short arms**.
- Compared to bulk DNA, satellites differ in sequence complexity, GC-content, and tandem repeat size.



Who cares about satellites?

Centromeres are Satellites!

- Without them... you die!
- aSatellites required for kinetochore binding.

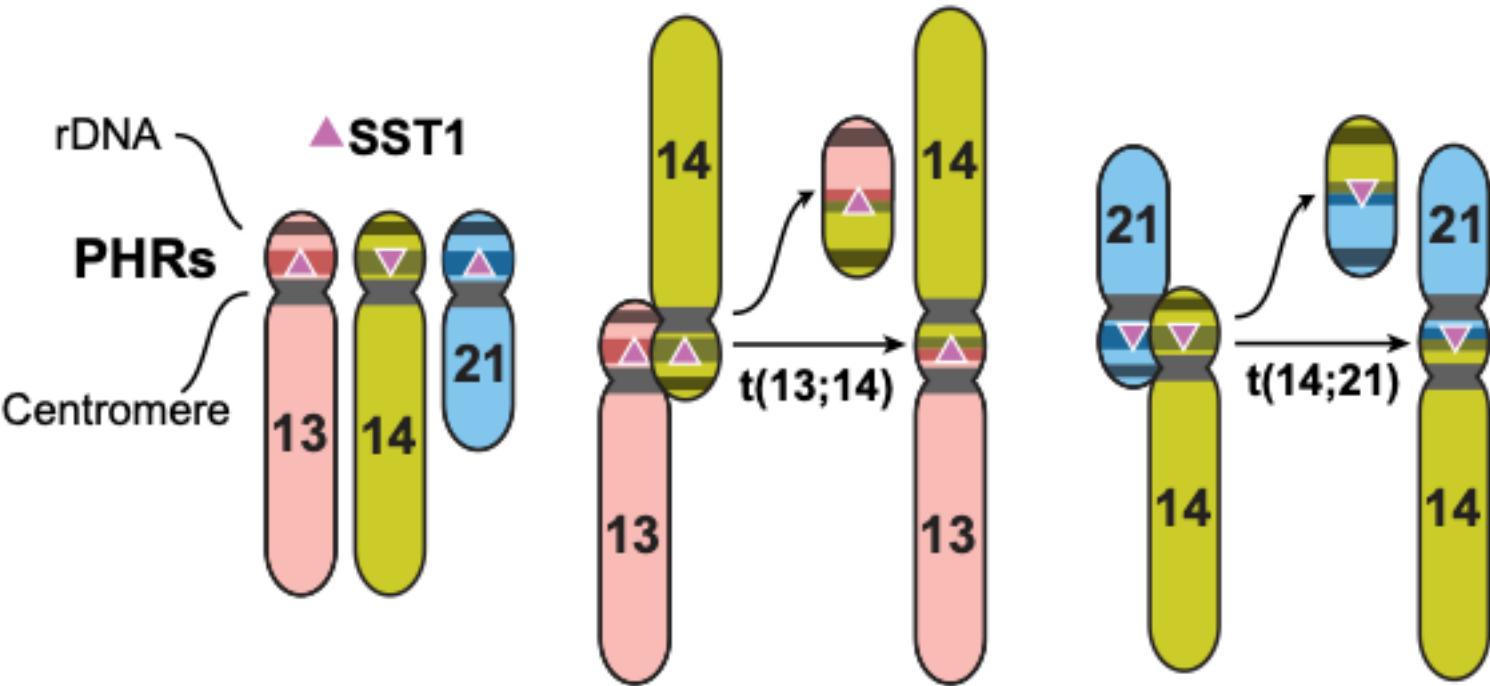


“If I synthesize this assembled genome, do I get the same, living organism back?”



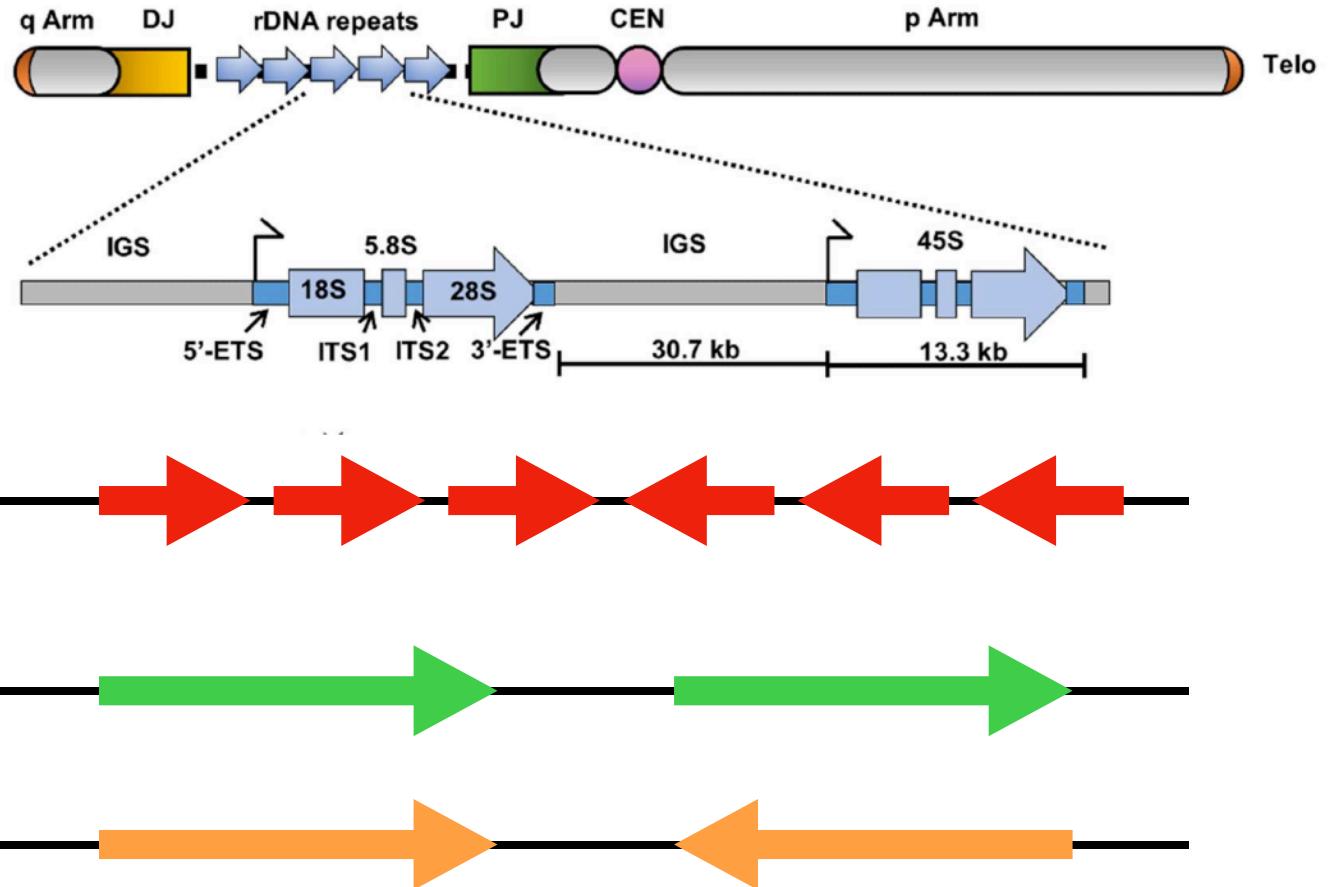
Robertsonian Translocations

- Present in 1/800 humans.
- Acrocentric short arms recombine during meiosis.
- SST1 satellite array is on an inverted duplication in Chr14. This leads to long-arm fusions & loss of short arms.



Other Repeats Worth Caring About

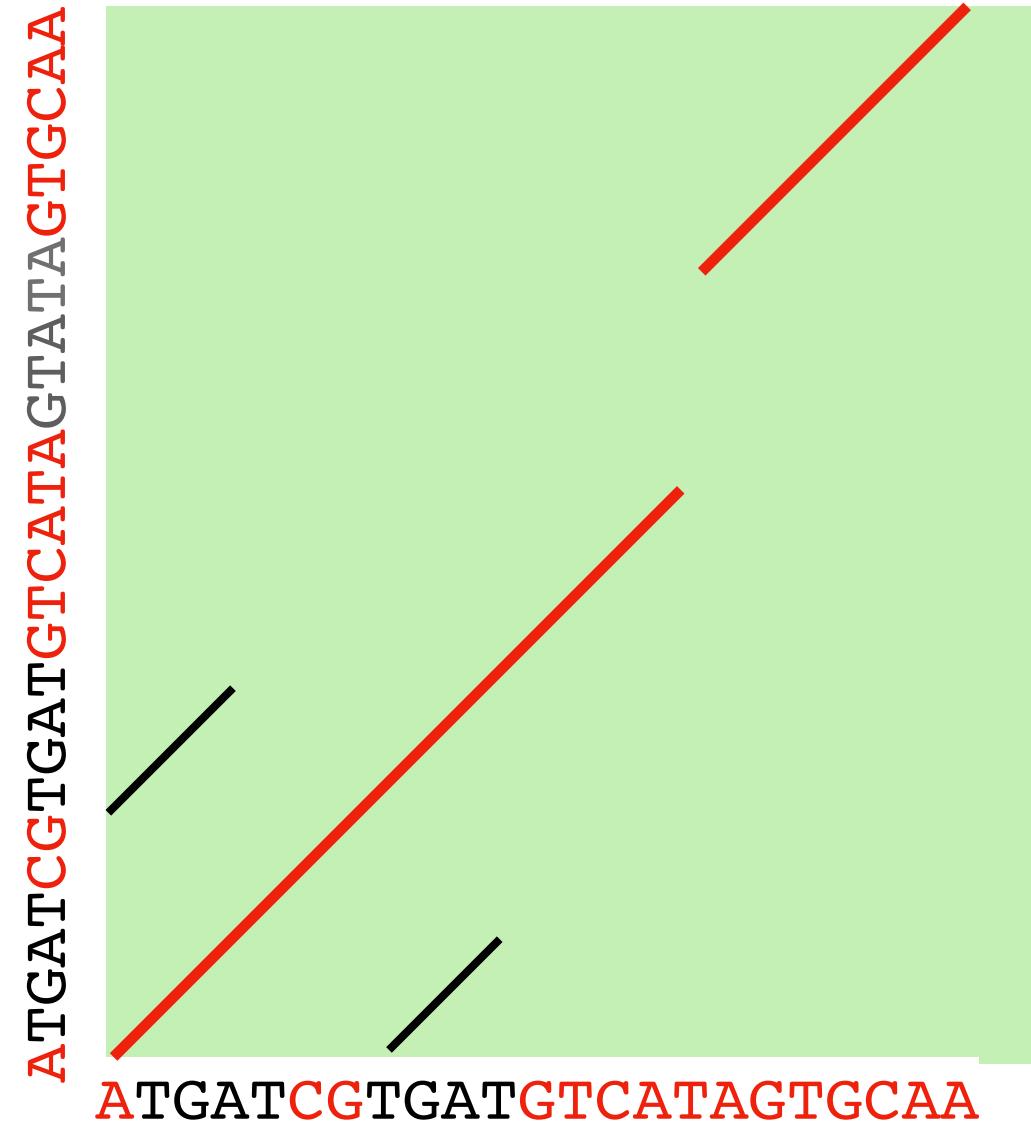
- Tandem gene Arrays:
 - *TSPY* - ChrY
 - *rDNA* - Acrocentrics
- Palindromes
- Large Duplications
- Large Inversions



How can we visualize repeats?

Dotplots

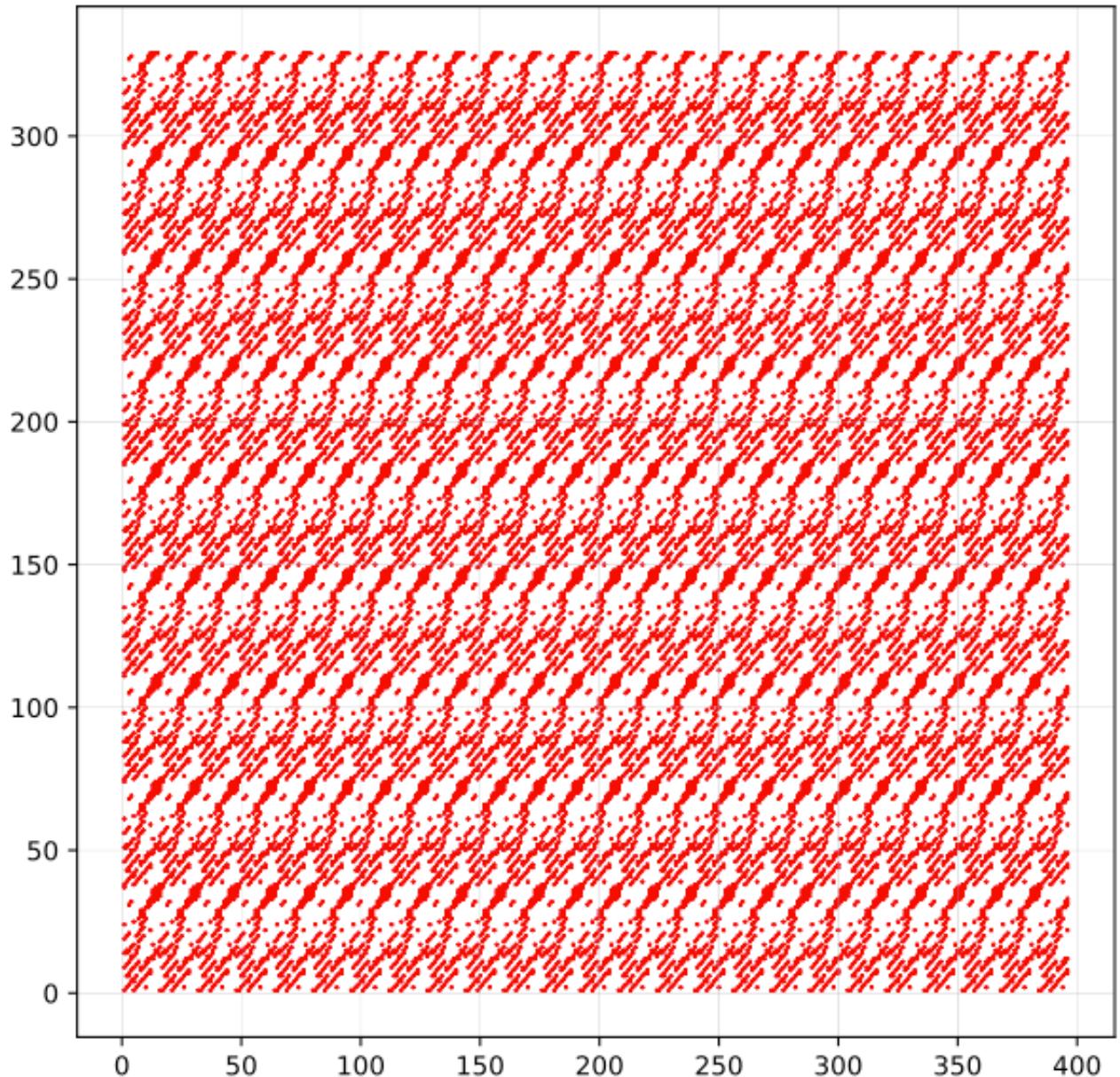
- A matrix consisting of two sequences, drawn with matches of length l (4 in this example).
- Shows self-similarity (self-plot) and orthologous relationships (comparative plot).
- “Vectorized” approach.



(GATTG) n

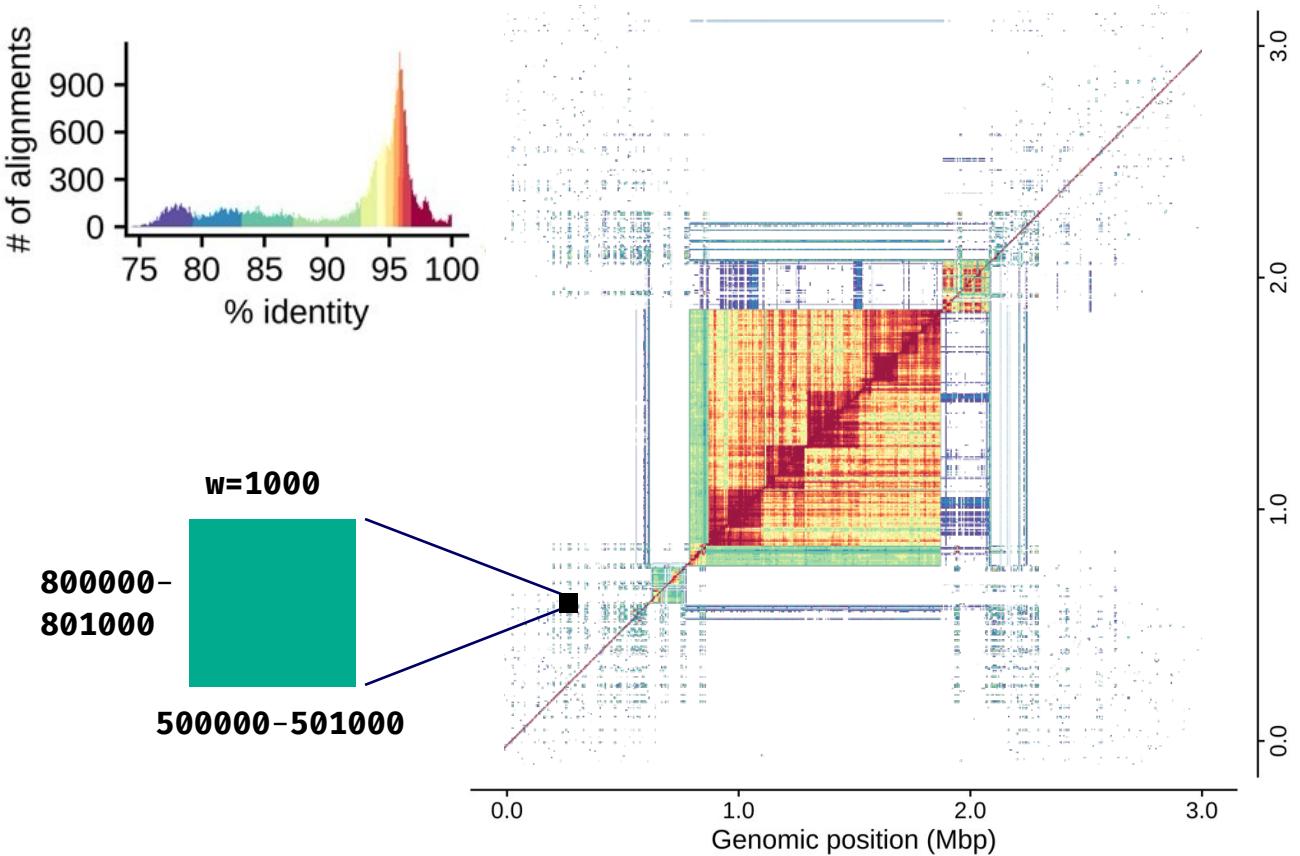
Dotplots

- Visualizing larger repeats quickly becomes overwhelming!!



Dotplots

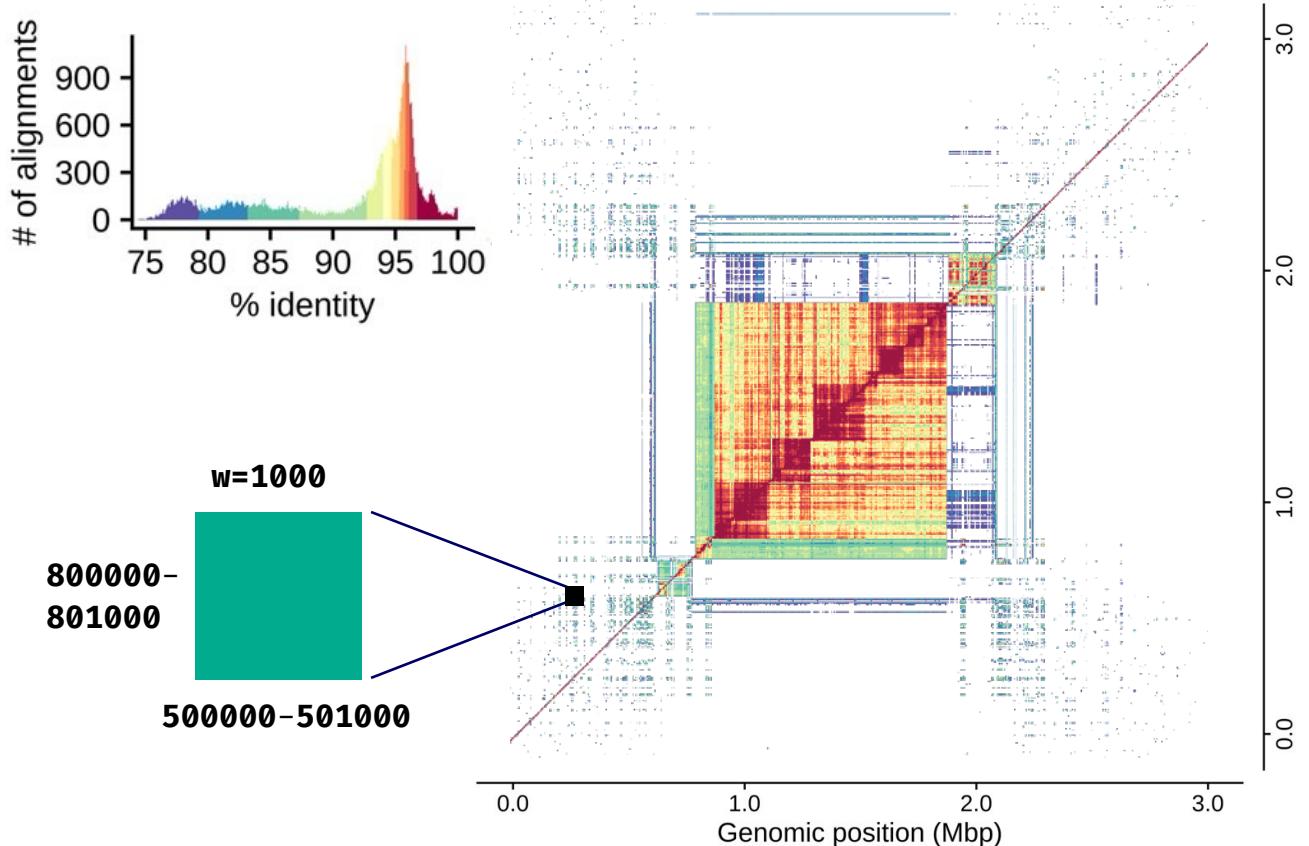
- Convert a vectorized into a rasterized model.
- Each pixel represents a pair of intervals of a fixed size w .
- Can compute sequence identity for each pairwise combination of intervals.



Computing Sequence Identity

$$100 \left(\frac{M}{M+X+I+D} \right)$$

- Use Minimap2 to estimate identity through matches (M), mismatches (X), insertions (I) and deletions (D)
- Scales quadratically with genome size. We can do better!

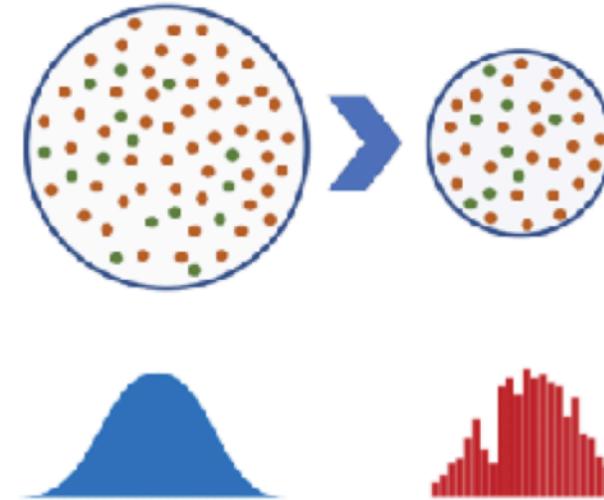




We can save time with sketching!

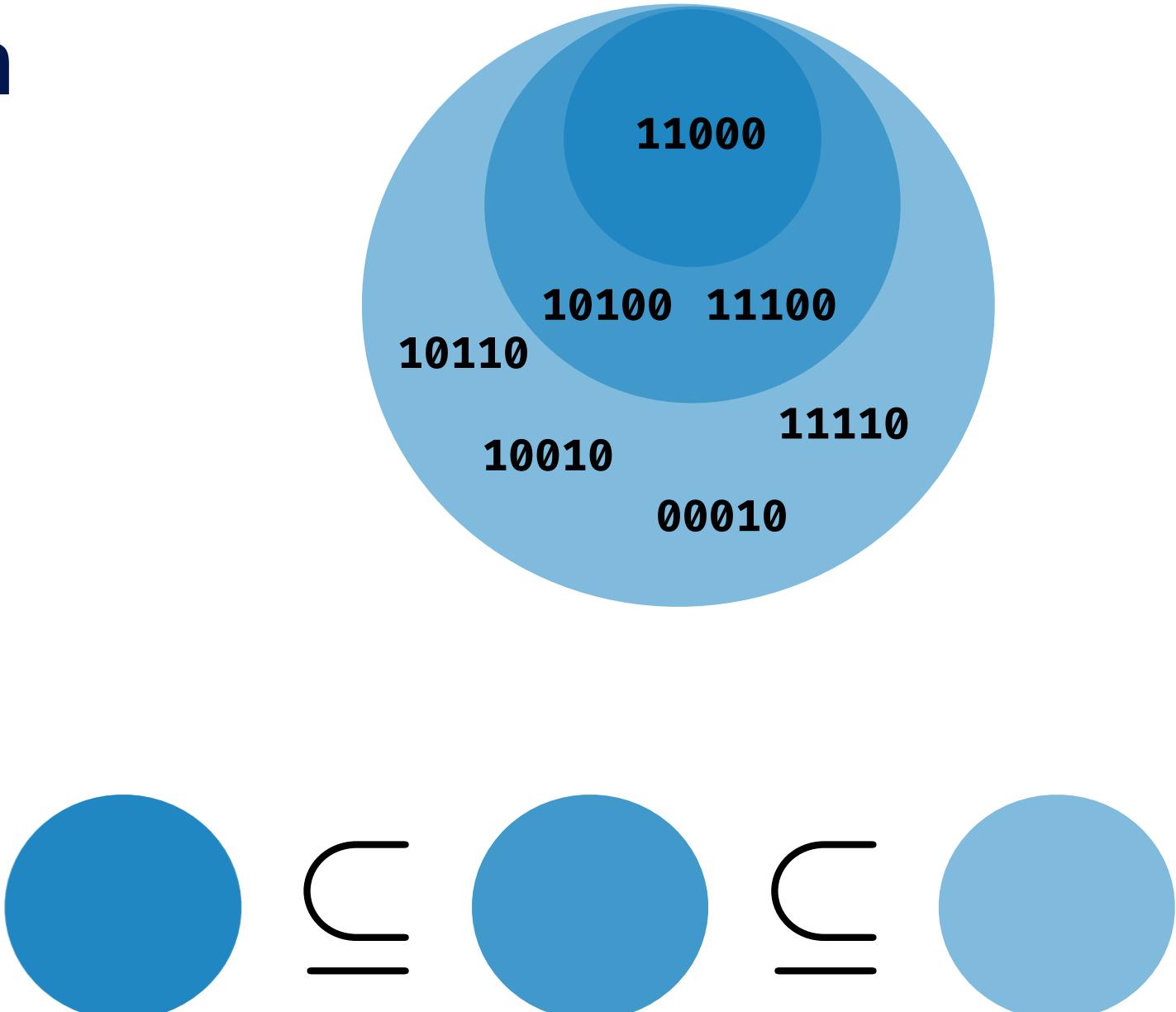
What is sketching?

- “Smart” subsample of the data
- Has a guarantee that certain properties are retained.



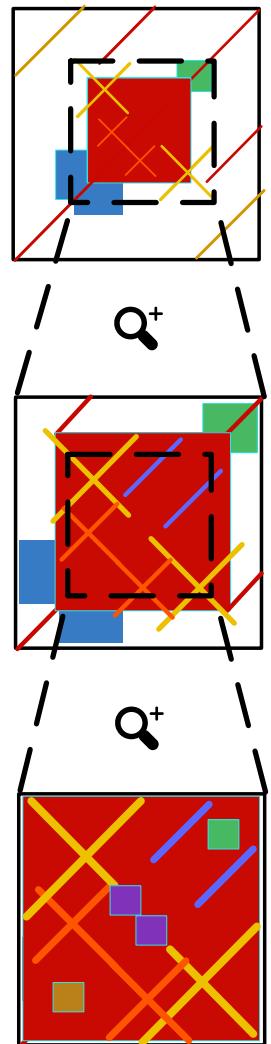
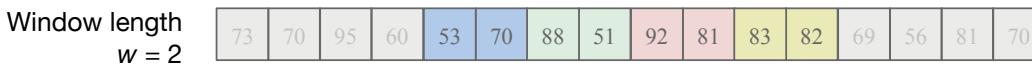
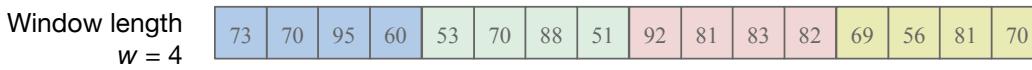
Modimizer Sketch

- Apply a hash function on all k -mers. Check if $k\text{-mer} \bmod s = 0$.
- If s is a power of 2, this creates a guaranteed hierarchy of k -mers
- Efficient: simply look at the last $\log_2(s)$ digits for 0's.



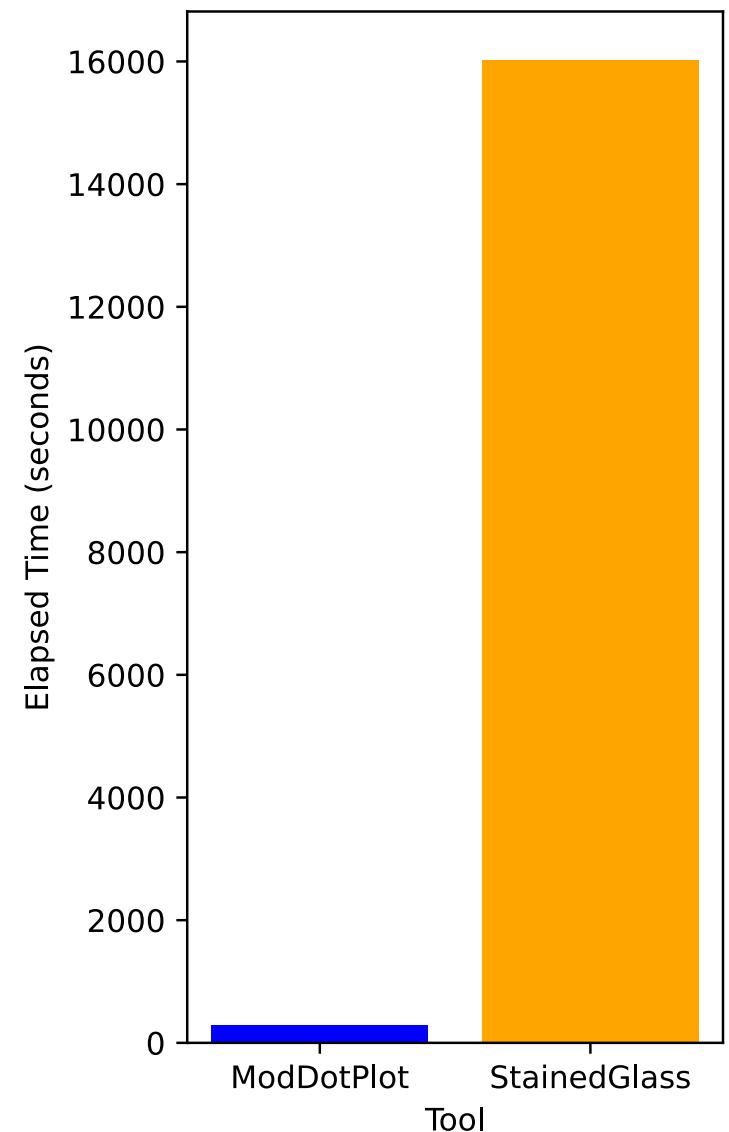
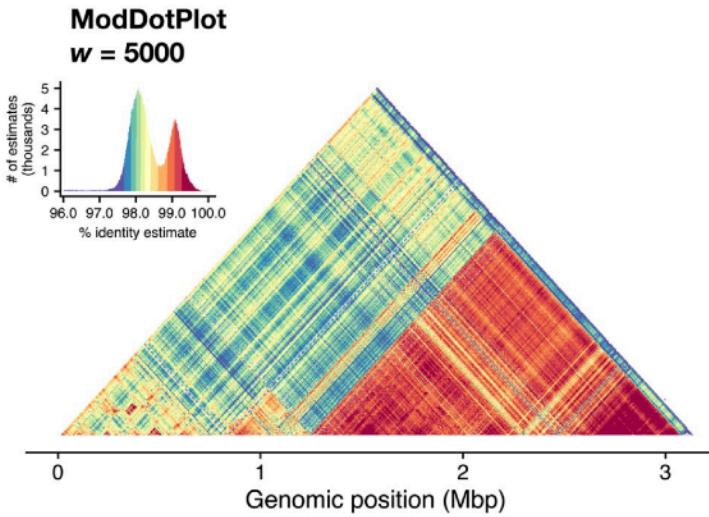
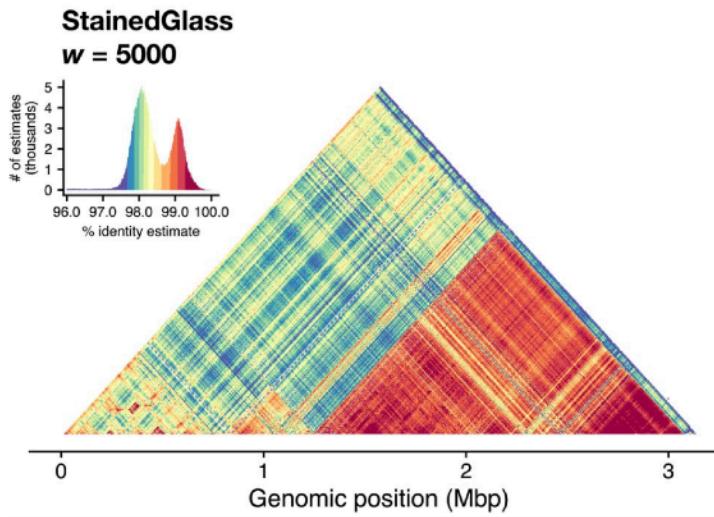
Interactive Plots

- Modimizer hierarchies allow for the construction of dynamically changing dotplots.
- Adapt s to the current field of view.
- In this example, 1 modimizer is used



Speed Comparison

- Computes identity for the *Arabidopsis* whole genome in ~5 minutes, while accurately representing the data!



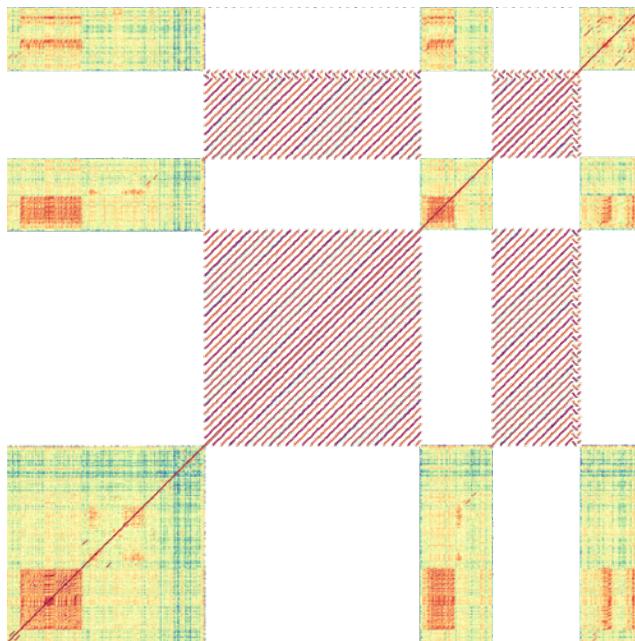
How do I interpret a dotplot?



Interpreting Self-Identity Plots

- ModDotPlot outputs 3 main files:

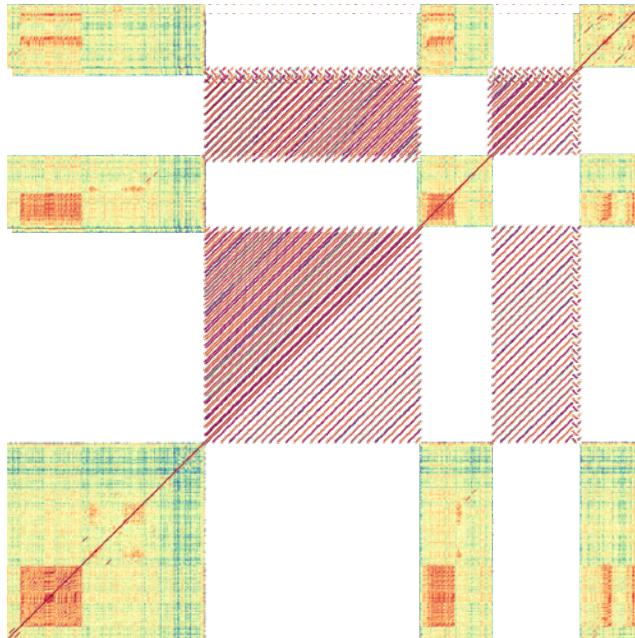
_FULL



Interpreting Self-Identity Plots

- ModDotPlot outputs 3 main files:

_FULL



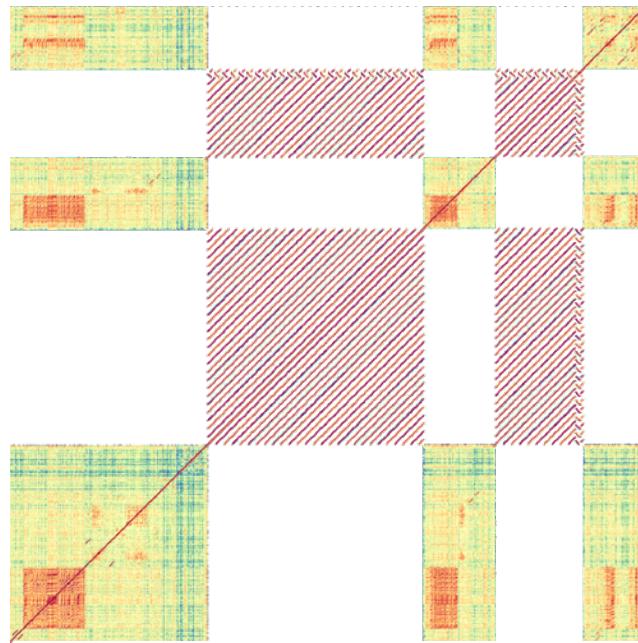
_TRI



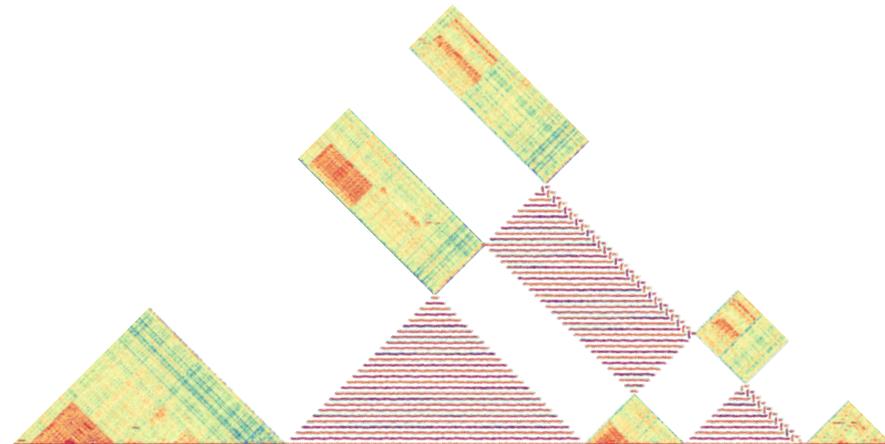
Interpreting Self-Identity Plots

- ModDotPlot outputs 3 main files:

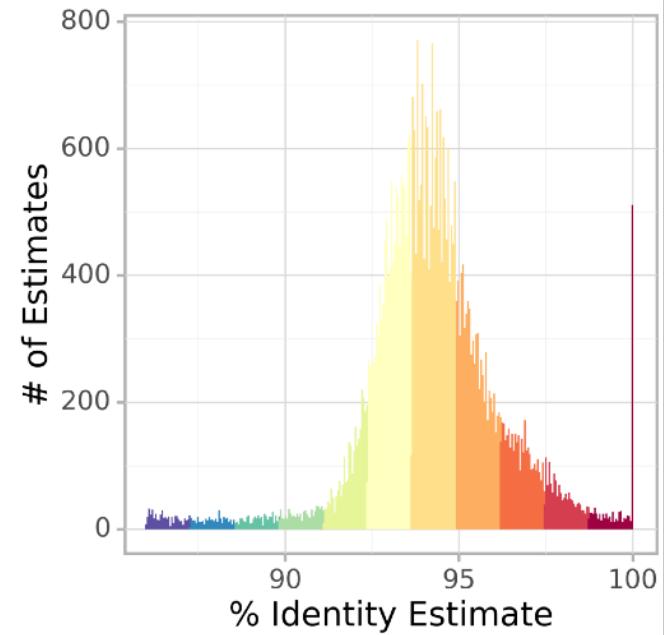
_FULL



_TRI

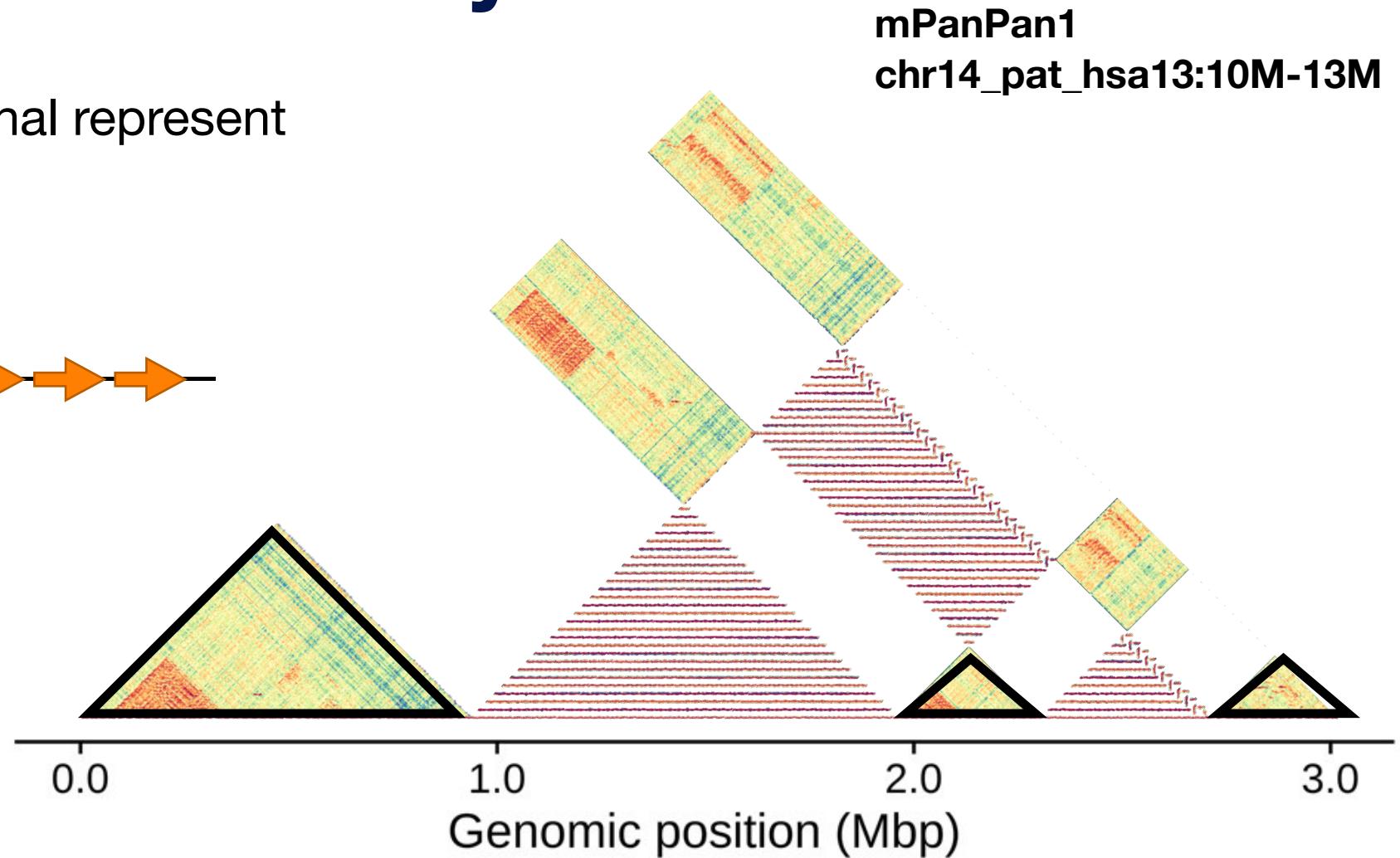
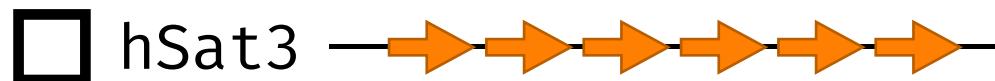


_HIST



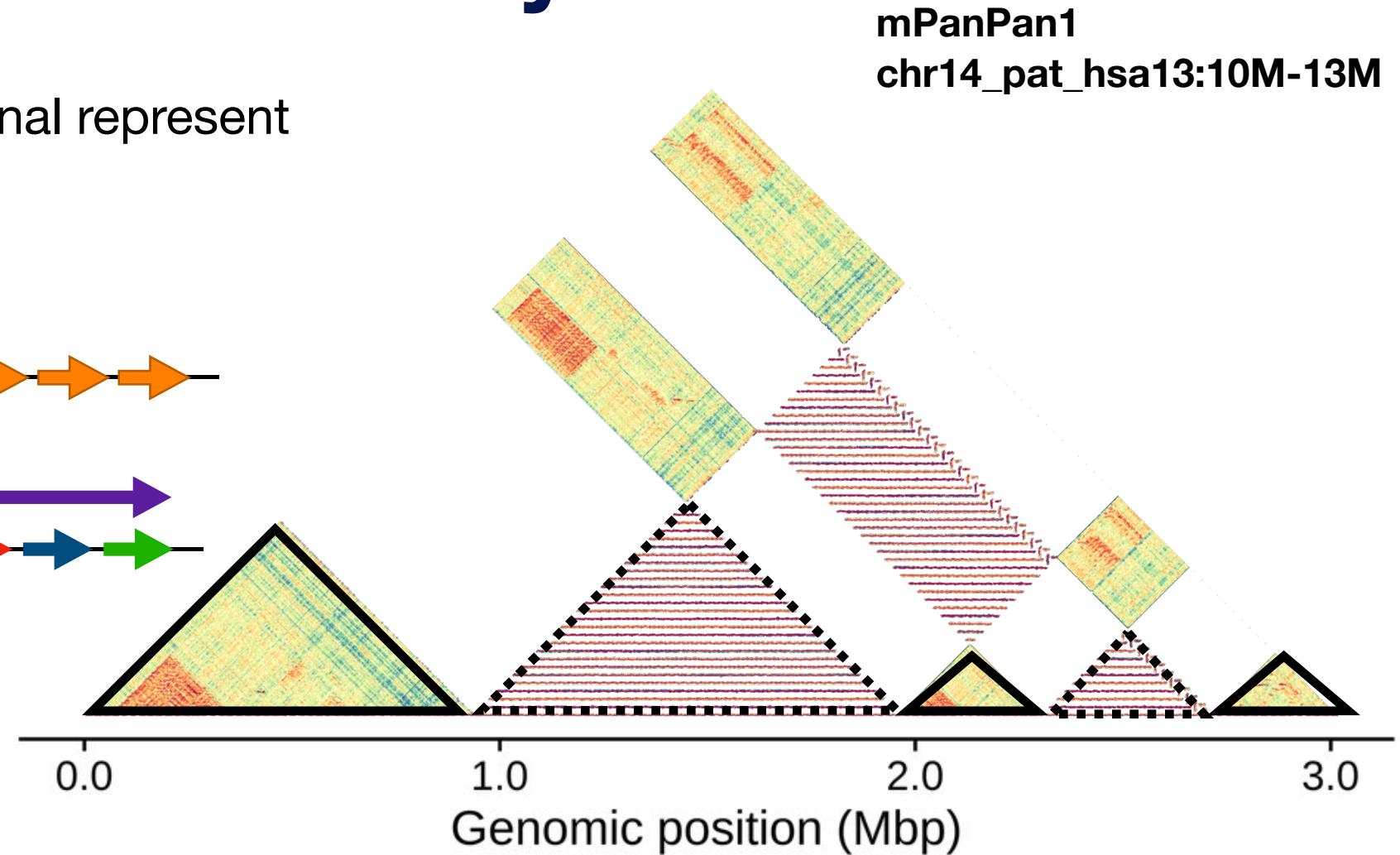
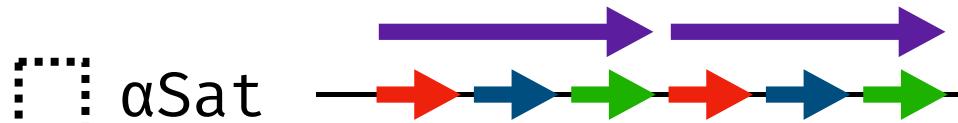
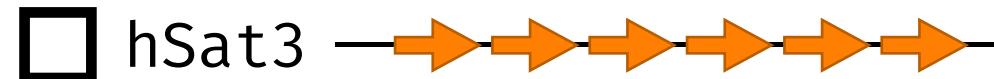
Interpreting Self-Identity Plots

- Triangles on the diagonal represent satellites

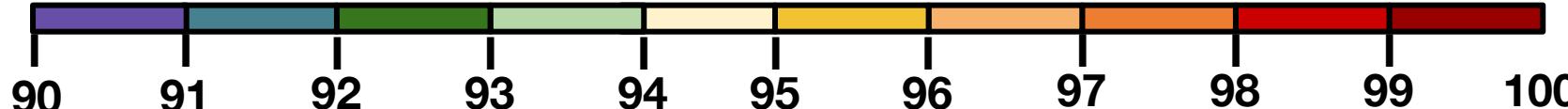


Interpreting Self-Identity Plots

- Triangles on the diagonal represent satellites

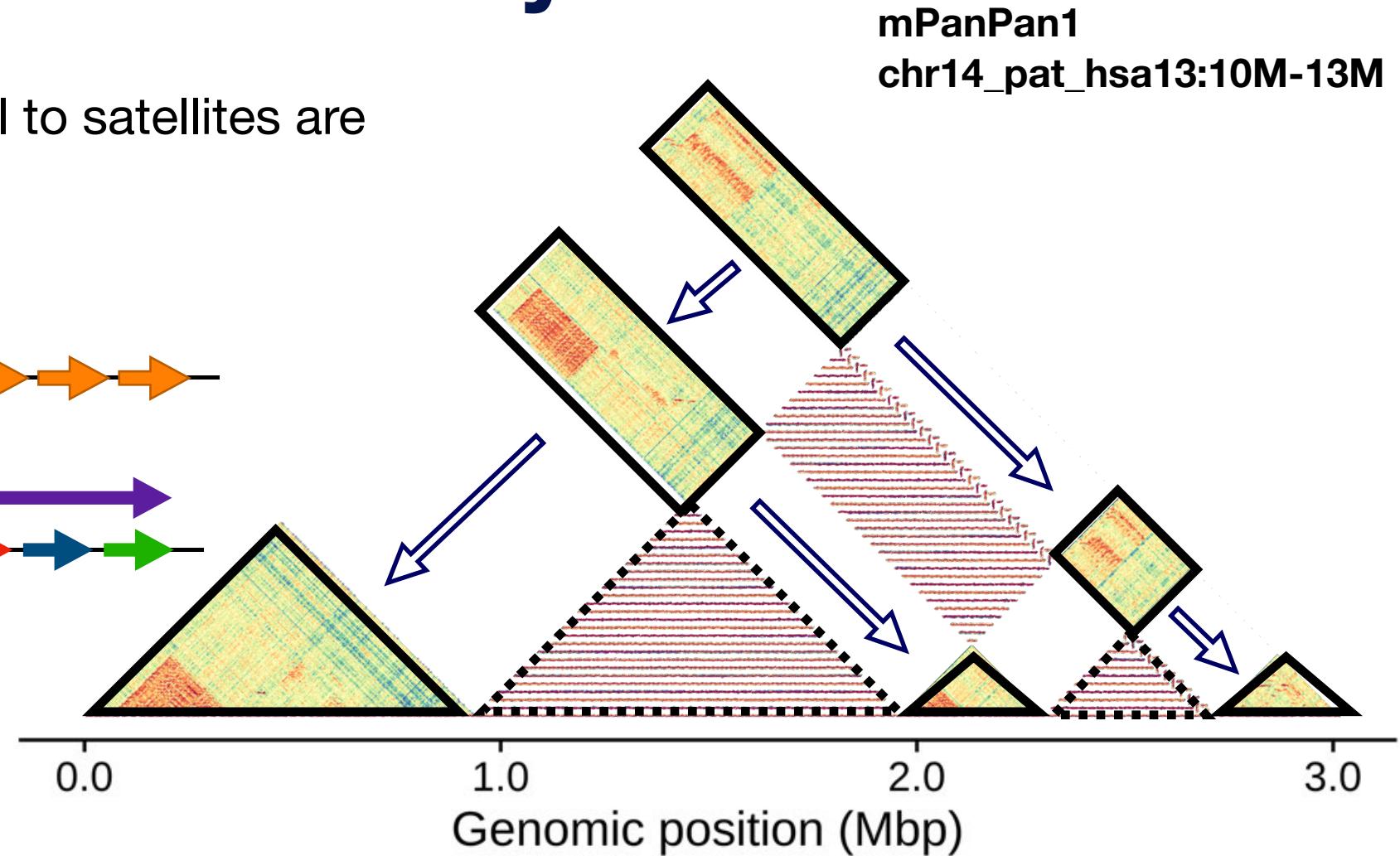
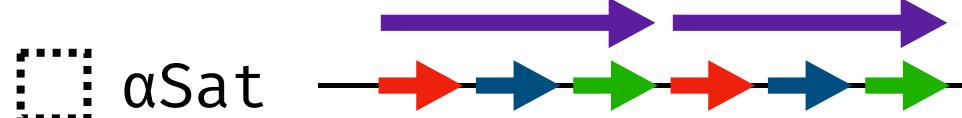
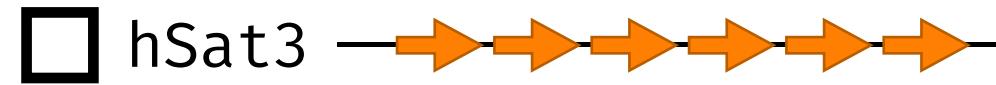


Average
Nucleotide
Identity

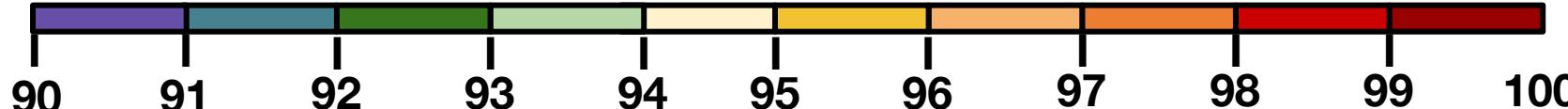


Interpreting Self-Identity Plots

- Rectangles orthogonal to satellites are the same type!

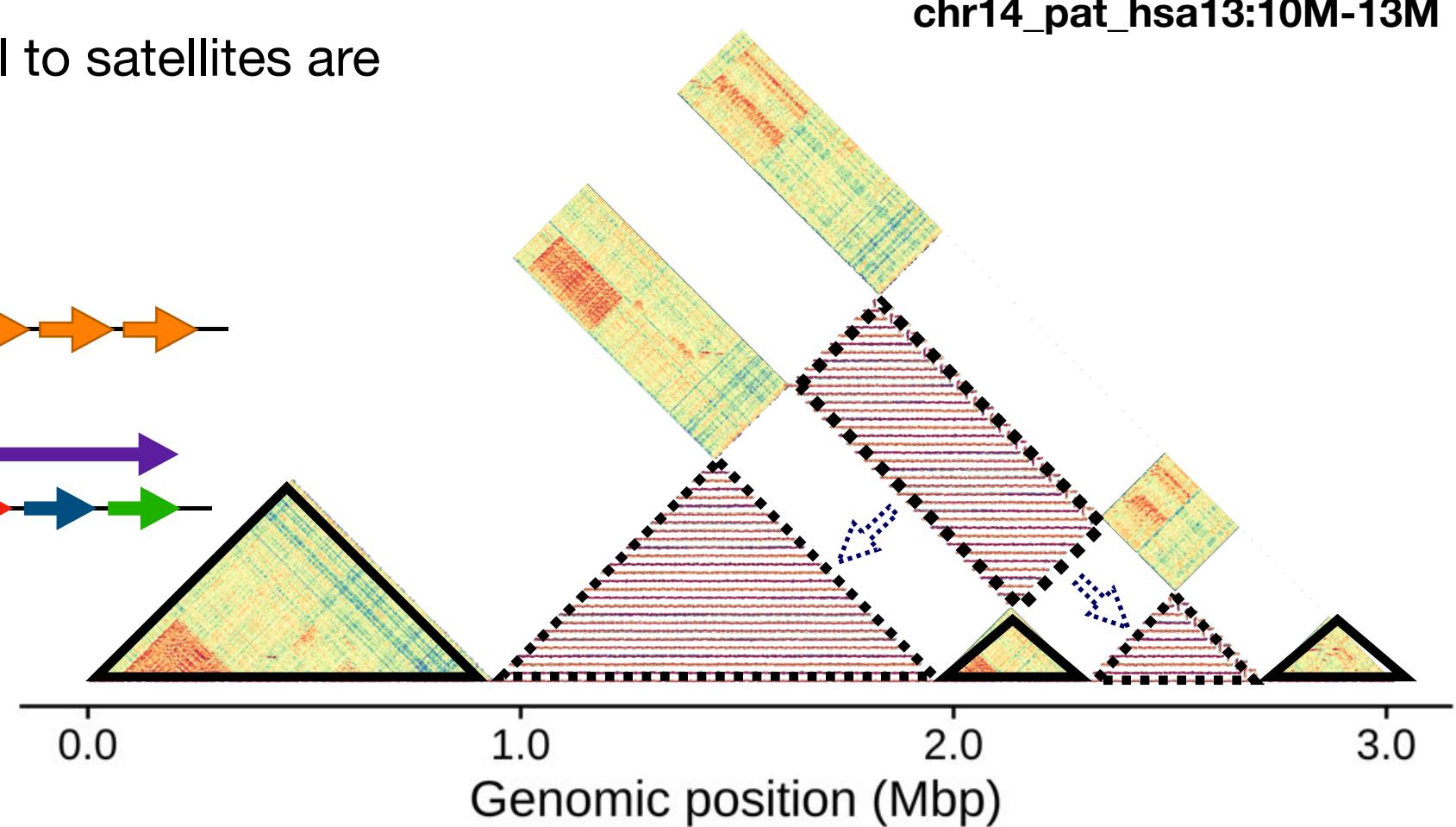
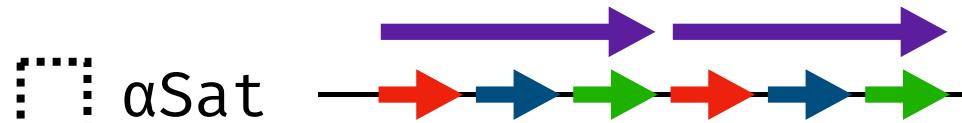
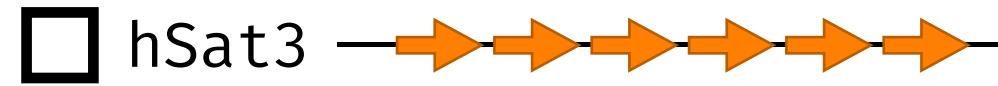


Average
Nucleotide
Identity

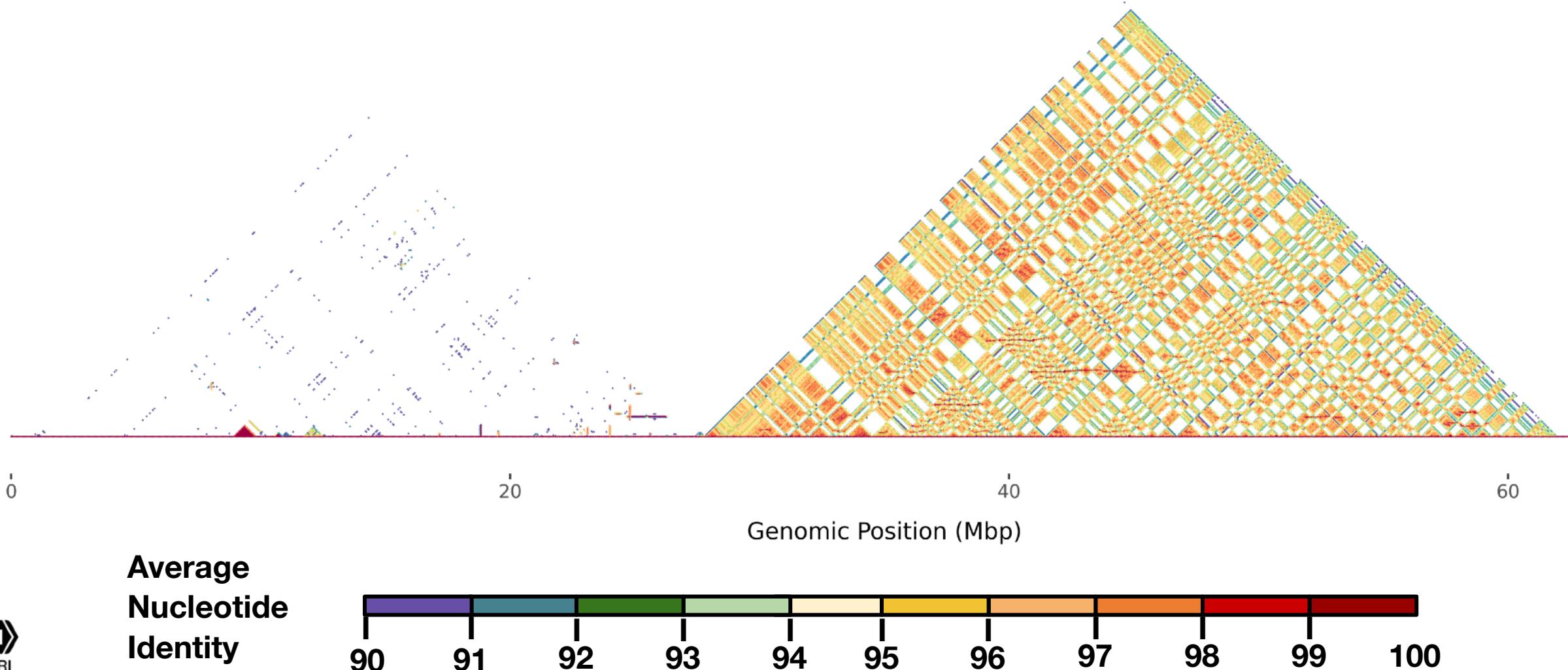


Interpreting Self-Identity Plots

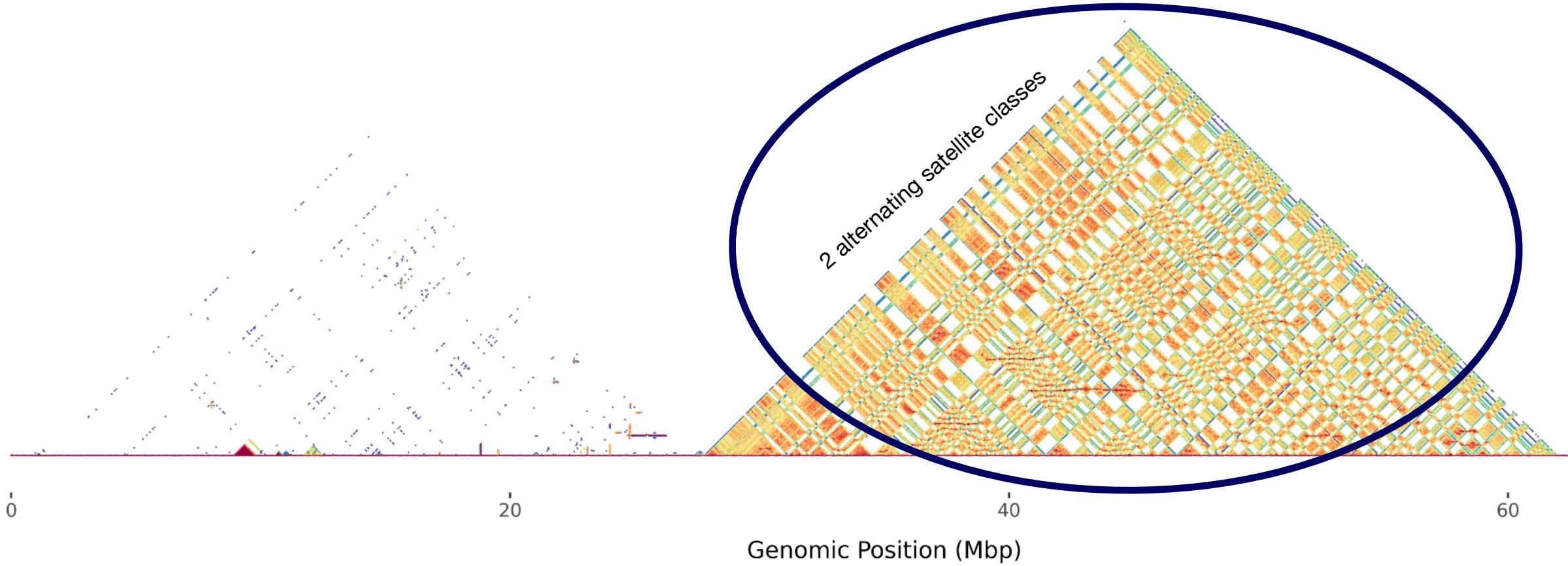
- Rectangles orthogonal to satellites are the same type!



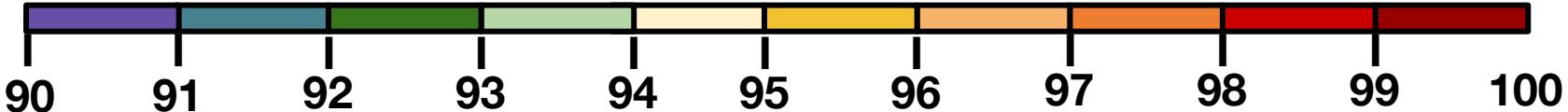
Interpreting Self-Identity Plots



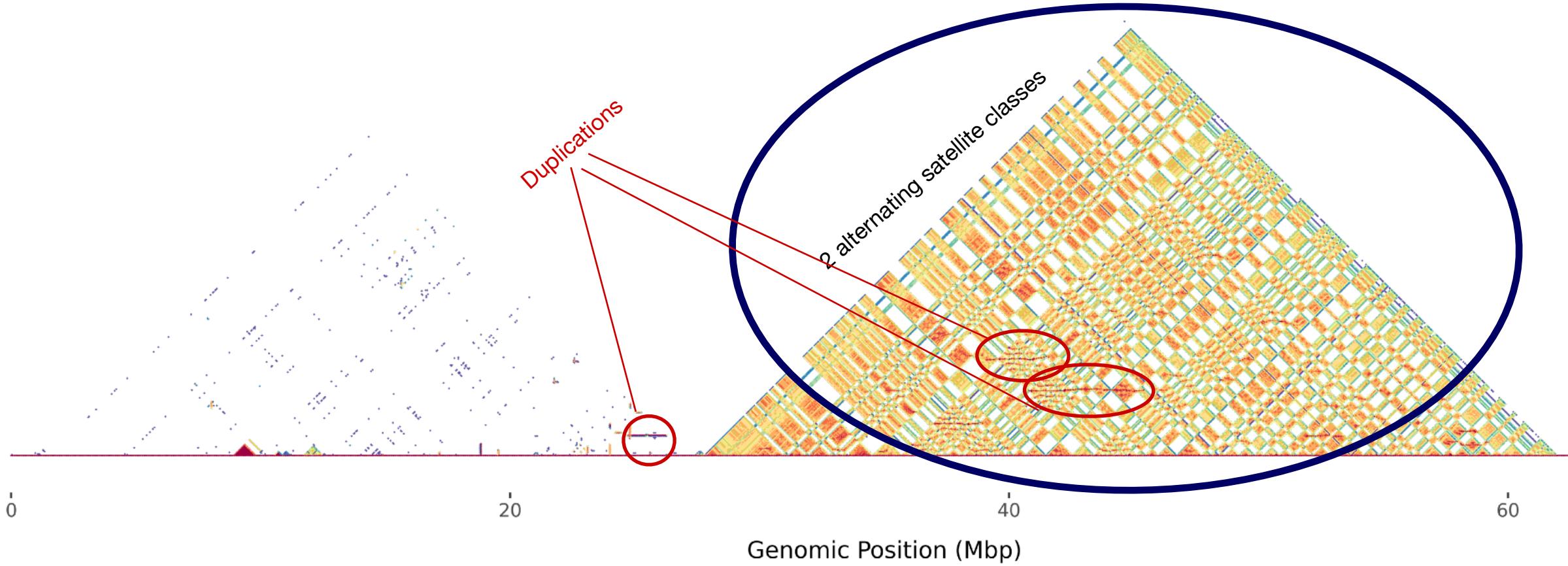
Interpreting Self-Identity Plots



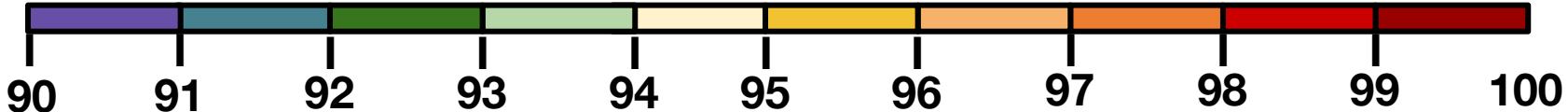
Average
Nucleotide
Identity



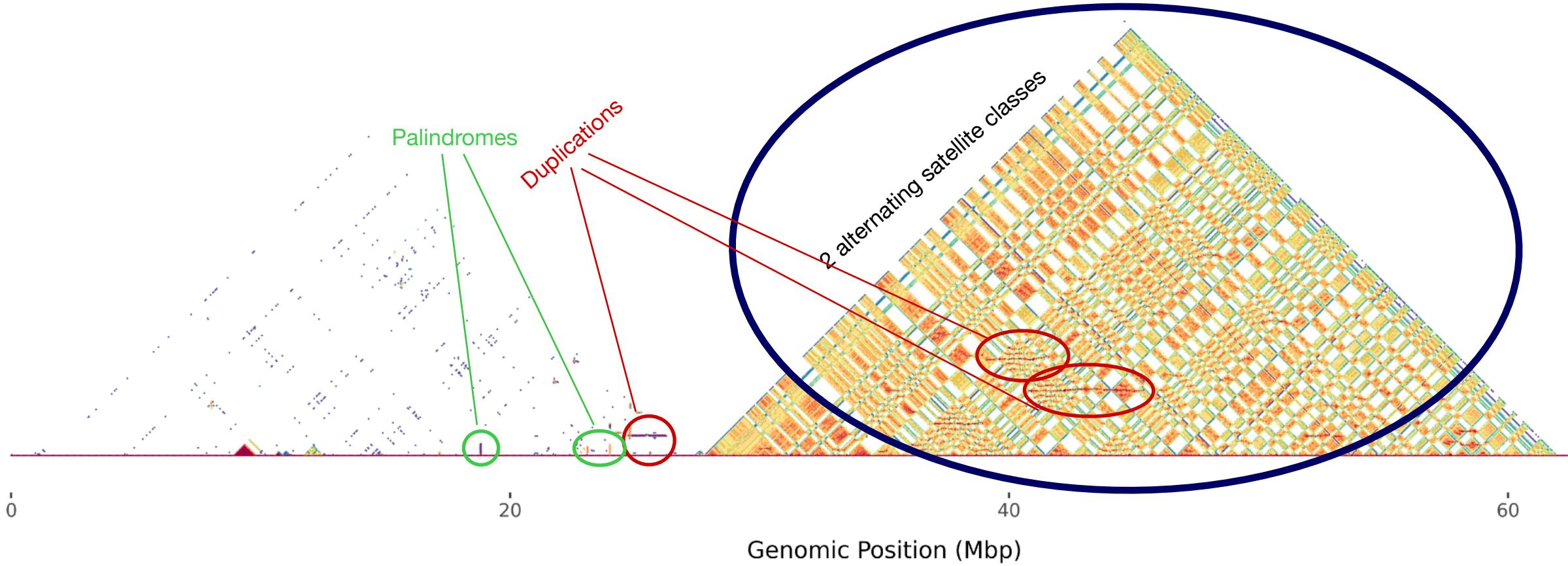
Interpreting Self-Identity Plots



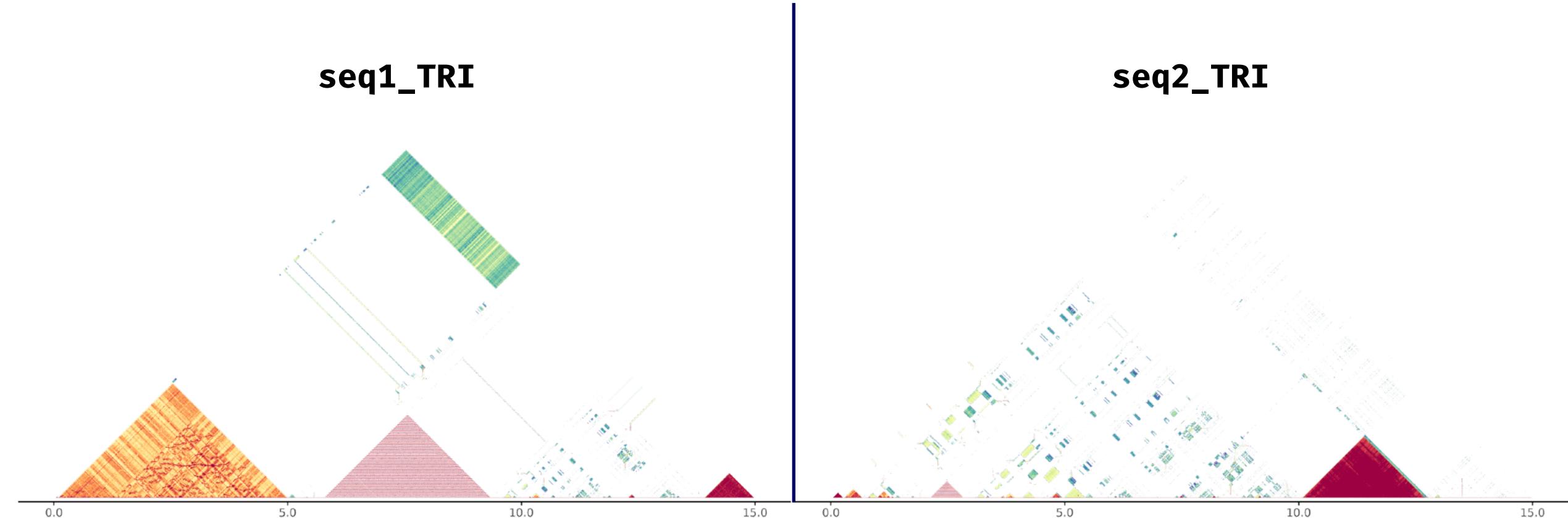
Average
Nucleotide
Identity



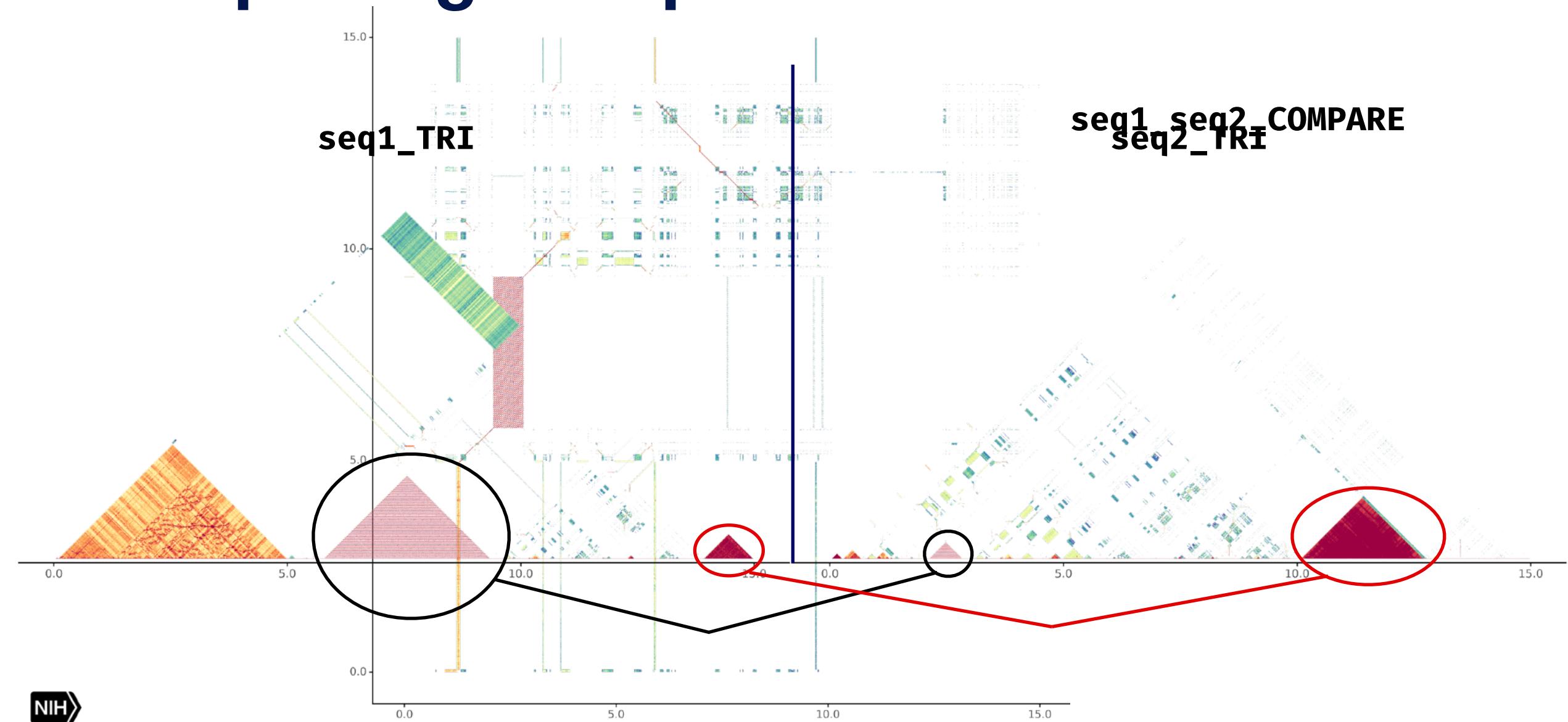
Interpreting Self-Identity Plots



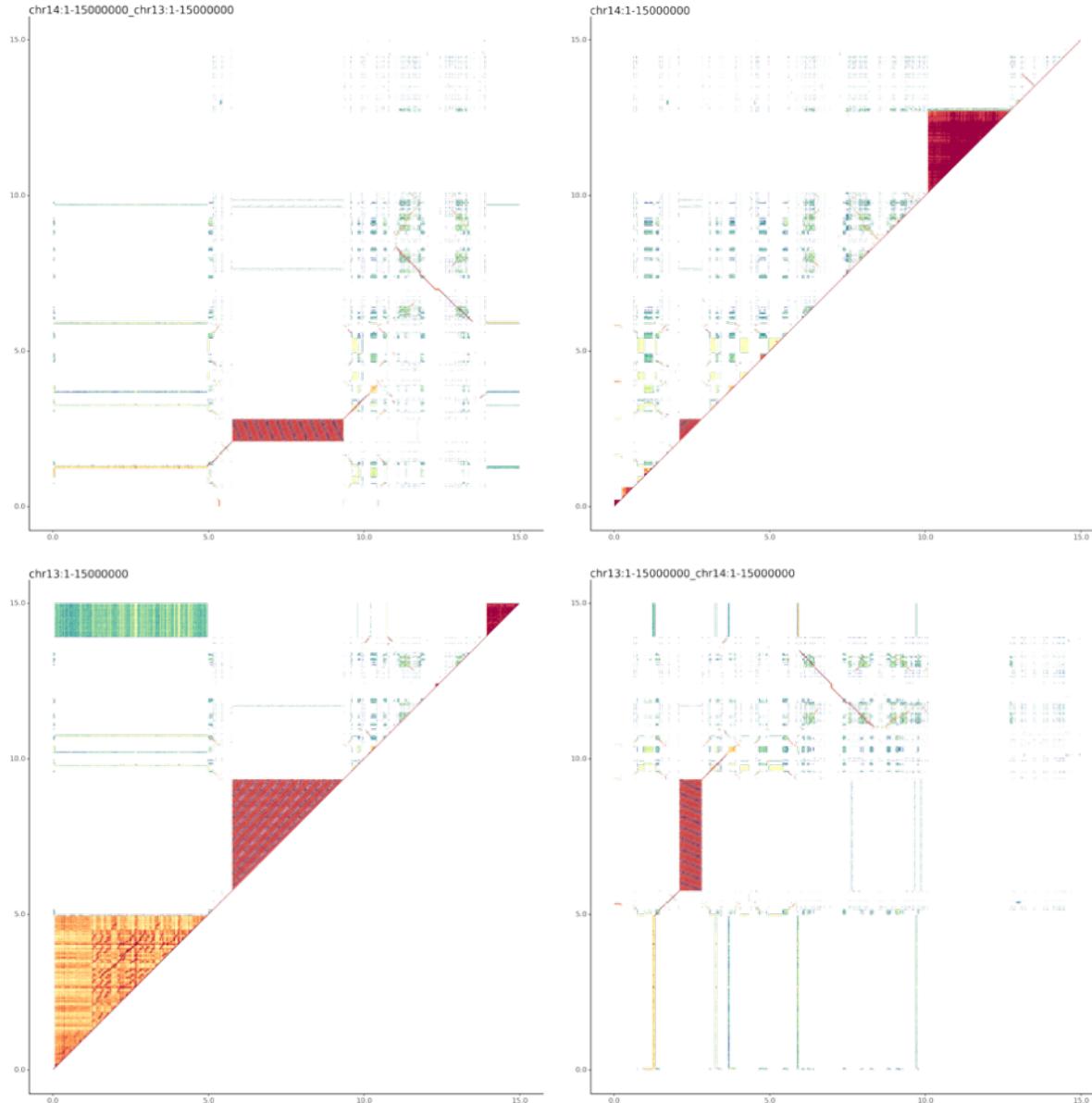
Interpreting Comparative Plots



Interpreting Comparative Plots

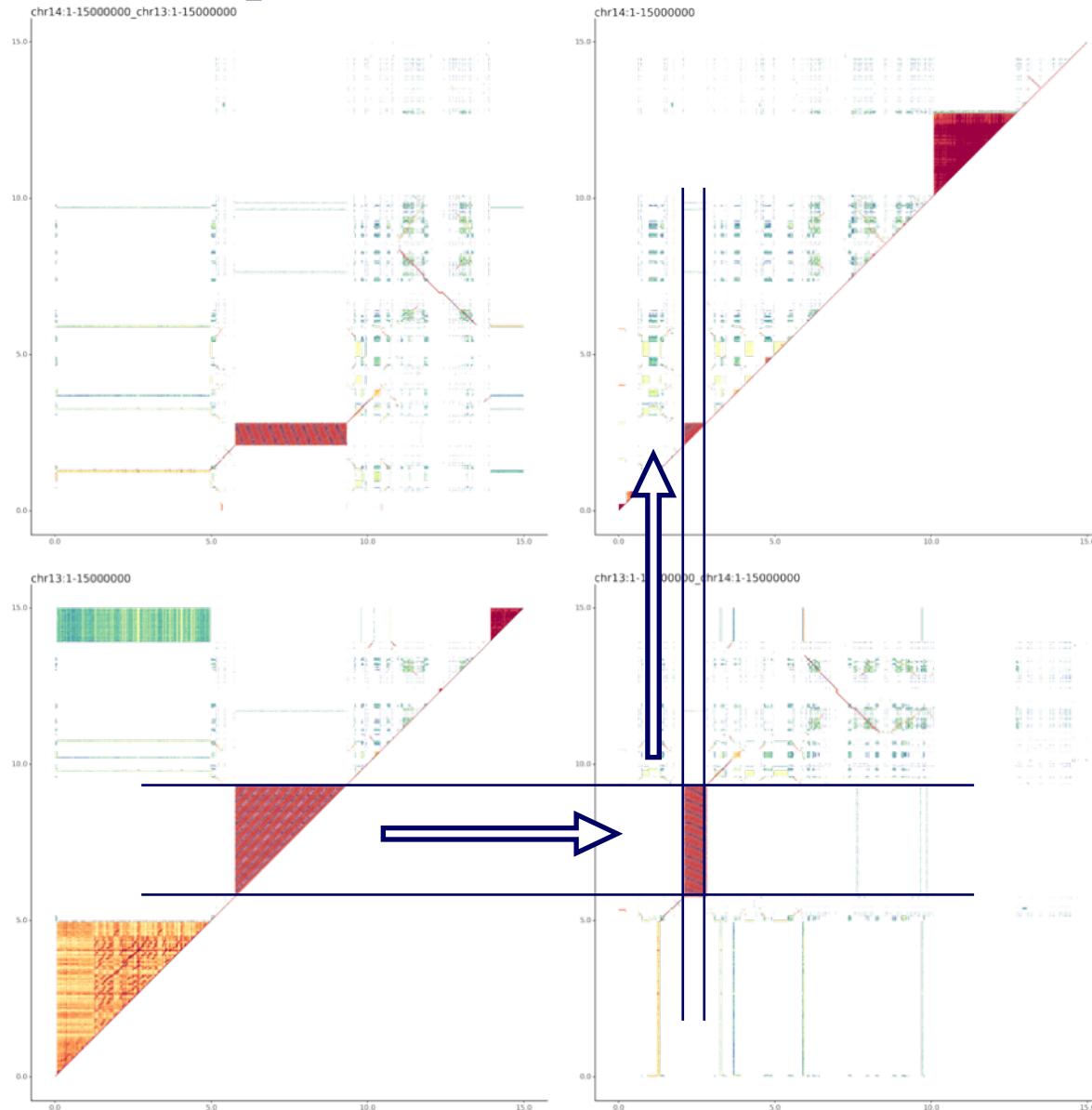
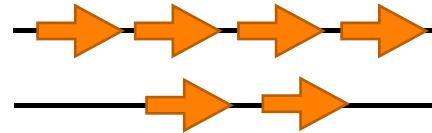


Interpreting Comparative Plots



Interpreting Comparative Plots

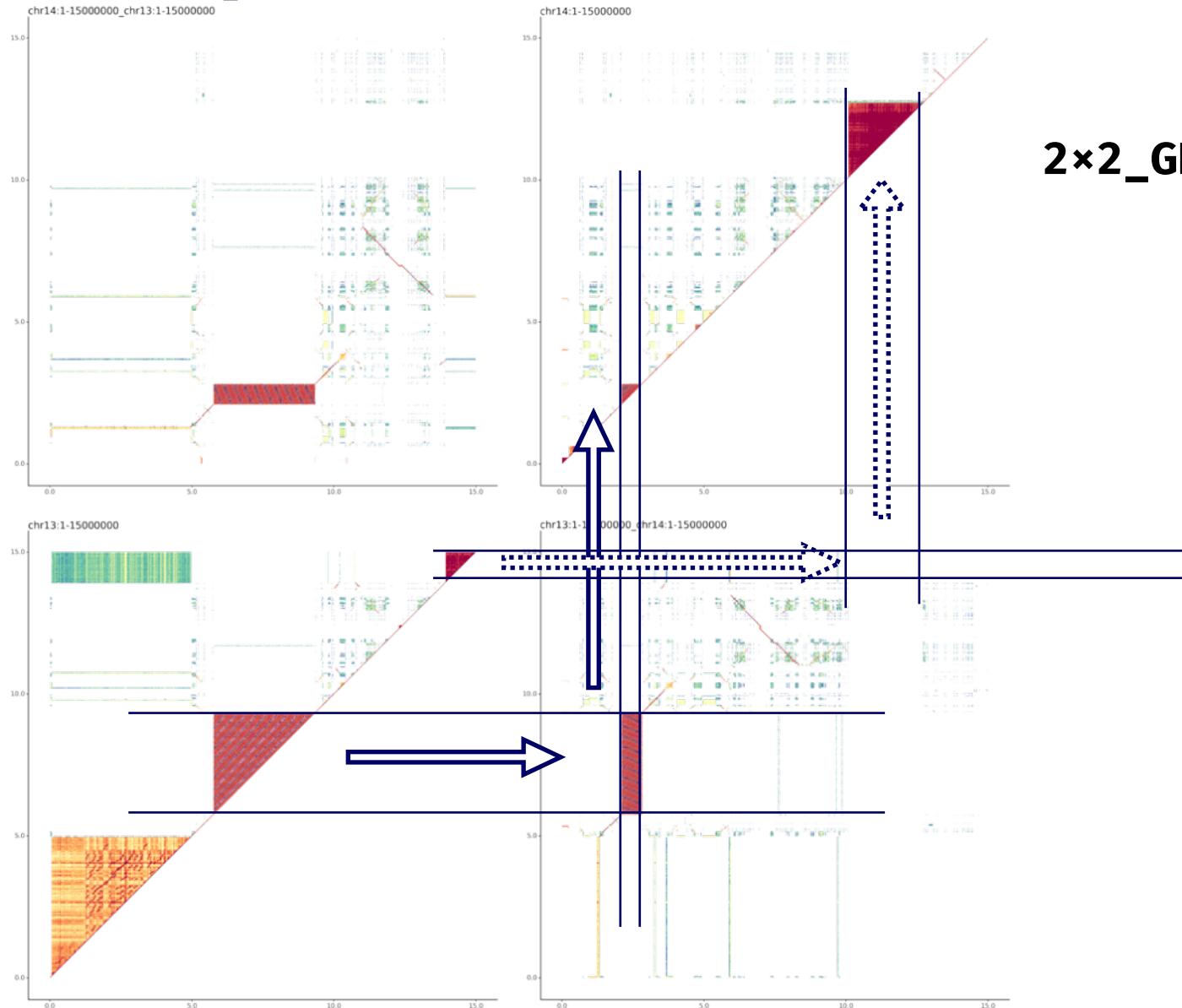
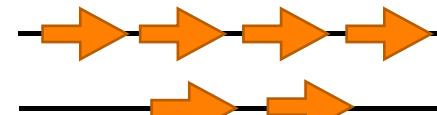
□ rDNA



2x2_GRID

Interpreting Comparative Plots

rDNA



hSat1

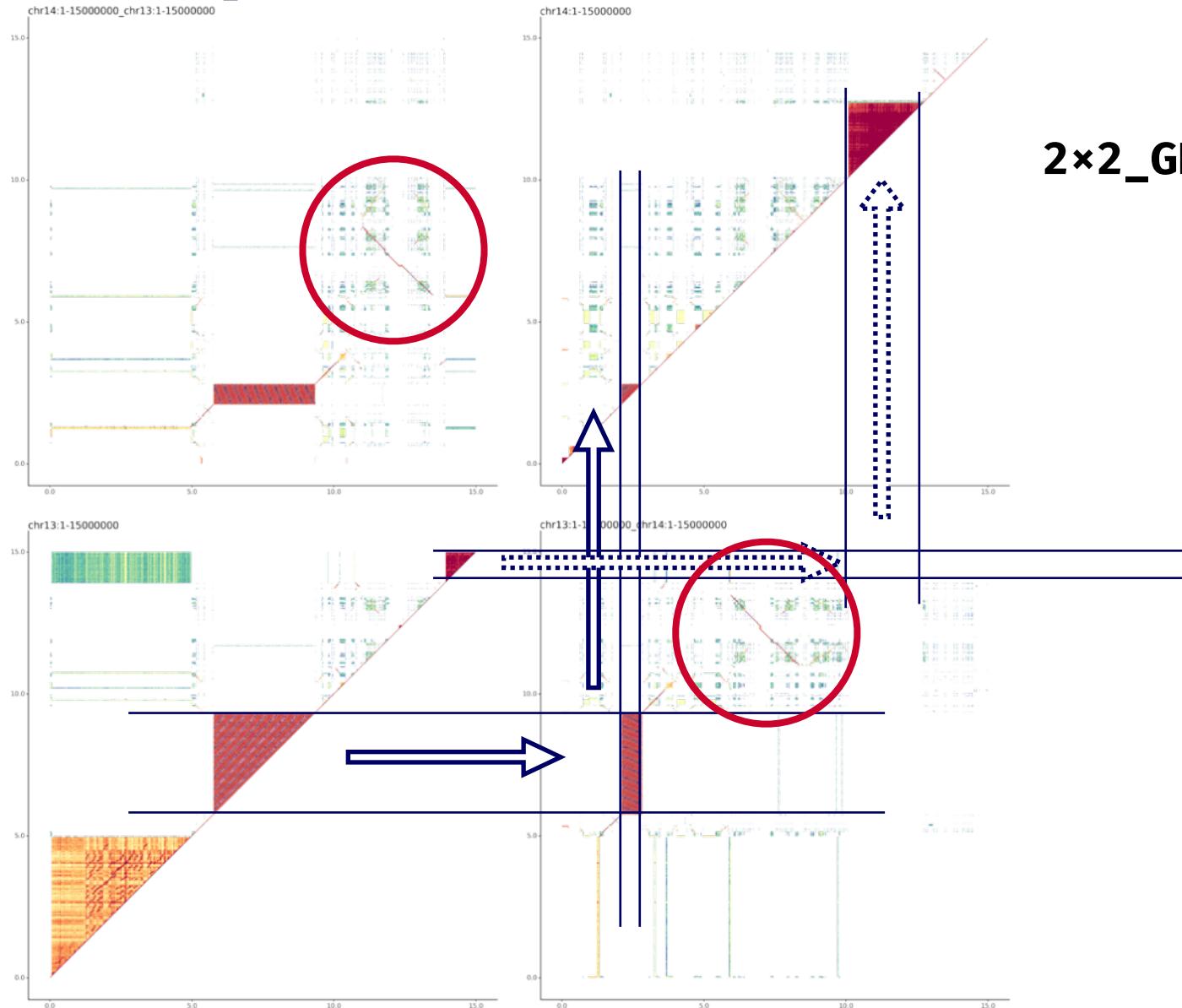
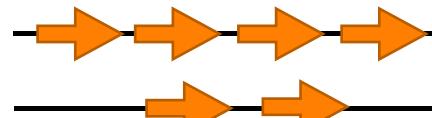


hSat3



Interpreting Comparative Plots

□ rDNA



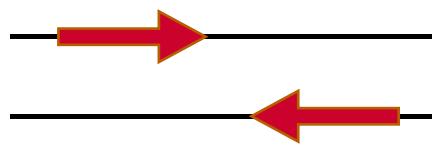
□ hSat1



□ hSat3



○ SST1



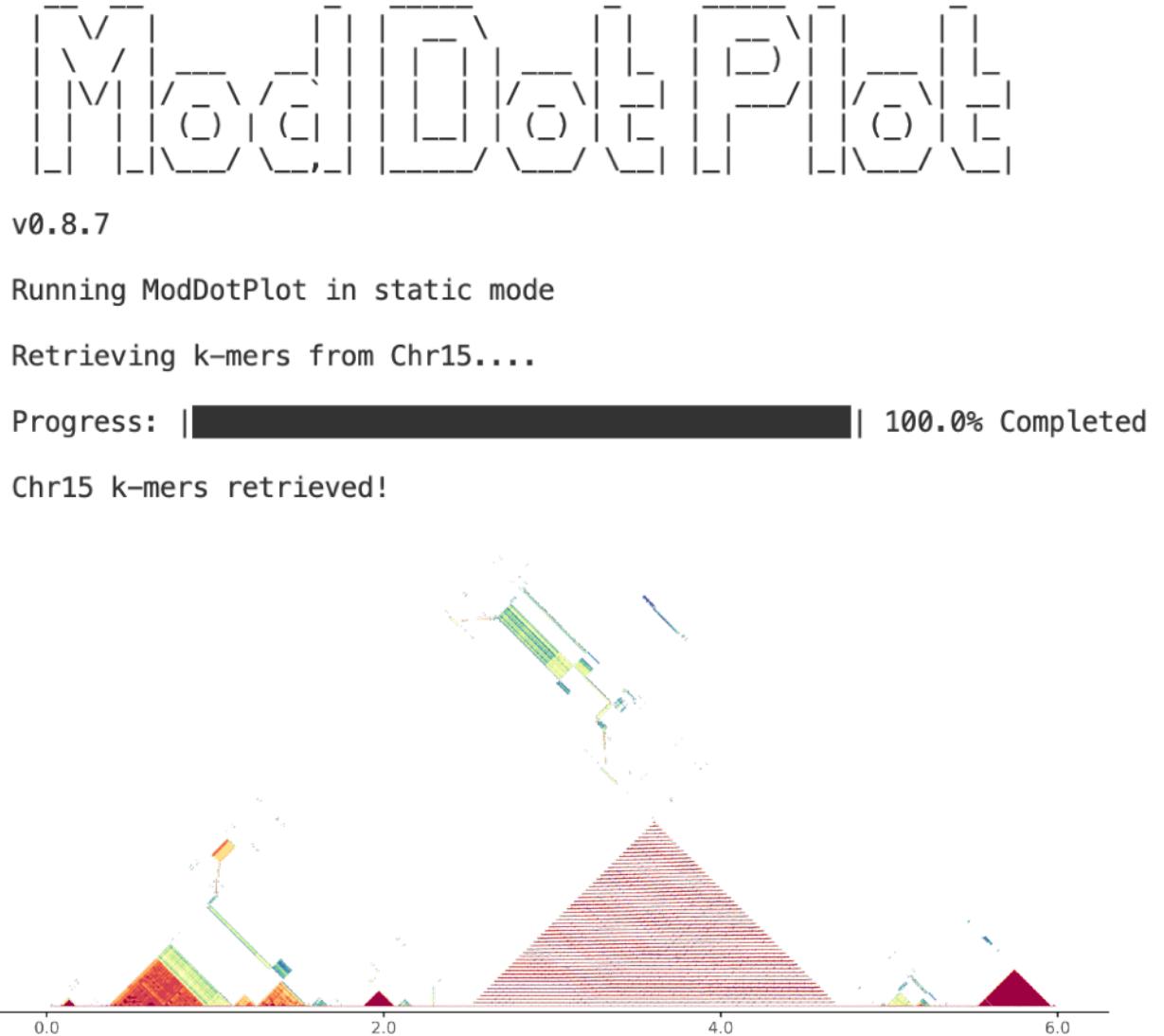
Running ModDotPlot

Running ModDotPlot

- ModDotPlot can be run in one of 2 modes:

- Static Mode:

- Produces static image files
- Good for publication quality figures
- Fast!



Running ModDotPlot

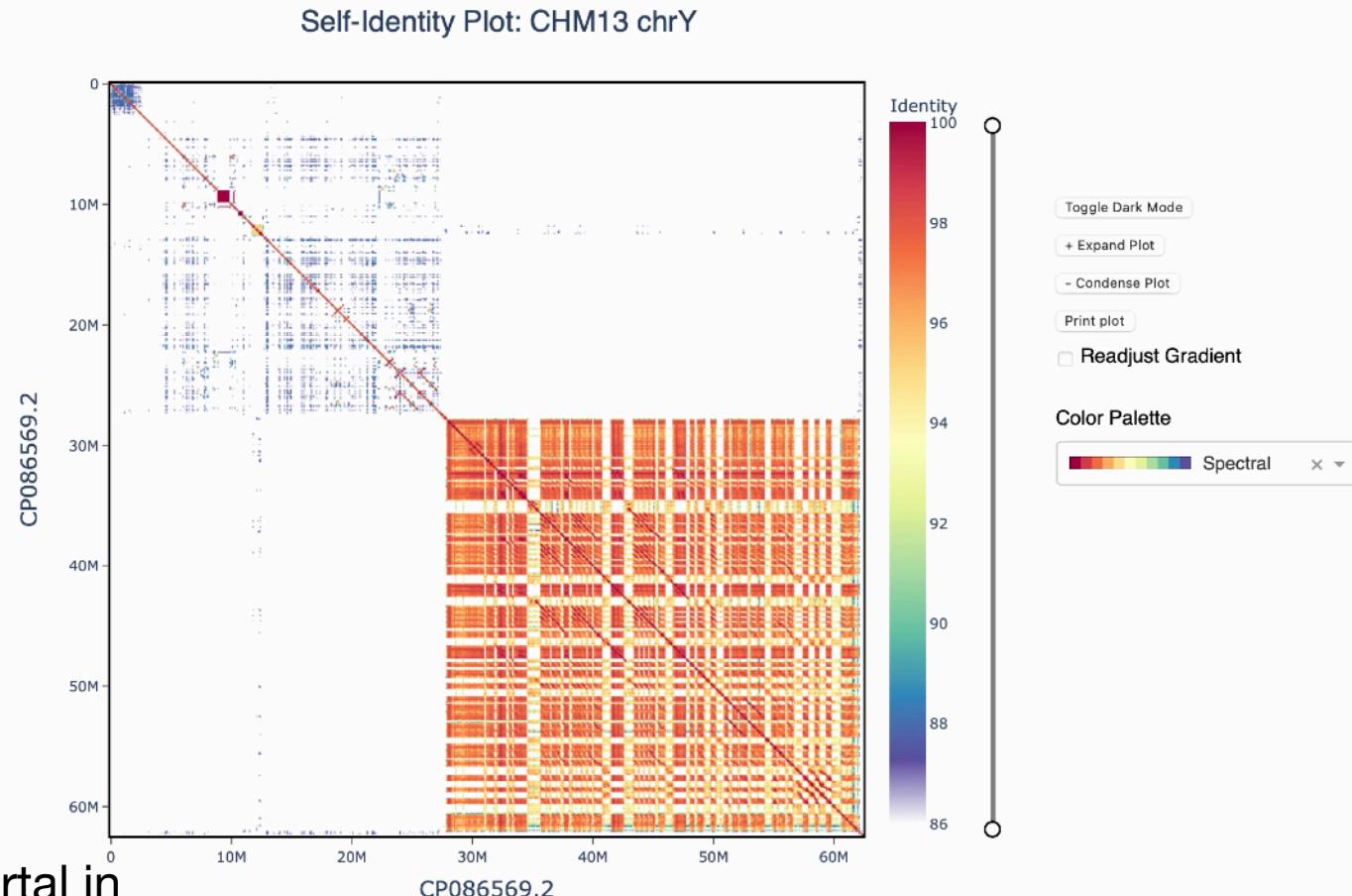
- ModDotPlot can be run in one of 2 modes:

- Static Mode:

- Produces static image files
- Good for publication quality figures
- Fast!

- Interactive Mode:

- Builds a plot hierarchy.
- Launches an interactive viewing portal in your web-browser.
- Slower, but allows for resolution adjustment & exact coordinates.



Running ModDotPlot - Static Mode

- Parameters to be familiar with:

- **-f / --fasta:** Input sequence. *Fasta format
- **-b / --bed:** Input sequence. *Paired end bed format
- **-o / --output-dir:** Output directory. Default: . /

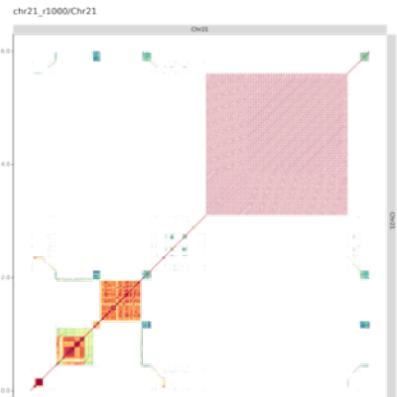
```
>Chr21_segment
AAACCCGTAACCTAACCGTAACCTAACCGTCAACCTAACCTAACCTAACCTAACCTAACCC
TAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAAC
CCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAAC
AACCCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAAC
CTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAAC
```

#query_name	query_start	query_end	reference_name	reference_start	reference_end
Chr21	1	6000	Chr21	1	6000
Chr21	1	6000	Chr21	6001	12000
Chr21	1	6000	Chr21	372001	378000
Chr21	1	6000	Chr21	378001	384000
Chr21	1	6000	Chr21	2196001	2202000
Chr21	1	6000	Chr21	2202001	2208000
Chr21	1	6000	Chr21	2592001	2598000
Chr21	1	6000	Chr21	2598001	2604000
Chr21	1	6000	Chr21	3066001	3072000
Chr21	1	6000	Chr21	3072001	3078000
Chr21	1	6000	Chr21	5766001	5766000

fasta



bed



plots

Running ModDotPlot - Static Mode

- Parameters to be familiar with:
 - **n: Sequence size.** *Not a parameter you can control through ModDotPlot!
 - **-r/--resolution:** Number of intervals (plot resolution). Default: 1000.
 - **-w/--window:** Interval length. *Note this is inversely proportional to resolution!

Running ModDotPlot - Static Mode

- Parameters to be familiar with:
 - **n: Sequence size.** *Not a parameter you can control through ModDotPlot!
 - **-r/--resolution:** Number of intervals (plot resolution). Default: 1000.
 - **-w/--window:** Interval length. *Note this is inversely proportional to resolution!

Runtime quadratically scales with *r* (not *n*!)

Running ModDotPlot - Static Mode

- Parameters to be familiar with:

- **n: Sequence size.** *Not a parameter you can control through ModDotPlot!
- **-r/--resolution: Number of intervals (plot resolution).** Default: 1000.
- **-w/--window: Interval length.** *Note this is inversely proportional to resolution!

```
Computing self identity matrix for Chr21...
Sequence length n: 6000000
Window size w: 6000
Modimizer sketch size: 1500
Plot Resolution r: 1000
Progress: |██████████| 100.0% Completed
```

```
Saved bed file to chr21_r1000/Chr21.bed
```

```
real    0m36.850s
user    0m36.594s
sys     0m0.784s
(venv) gitpod /workspace $
```

Running ModDotPlot - Static Mode

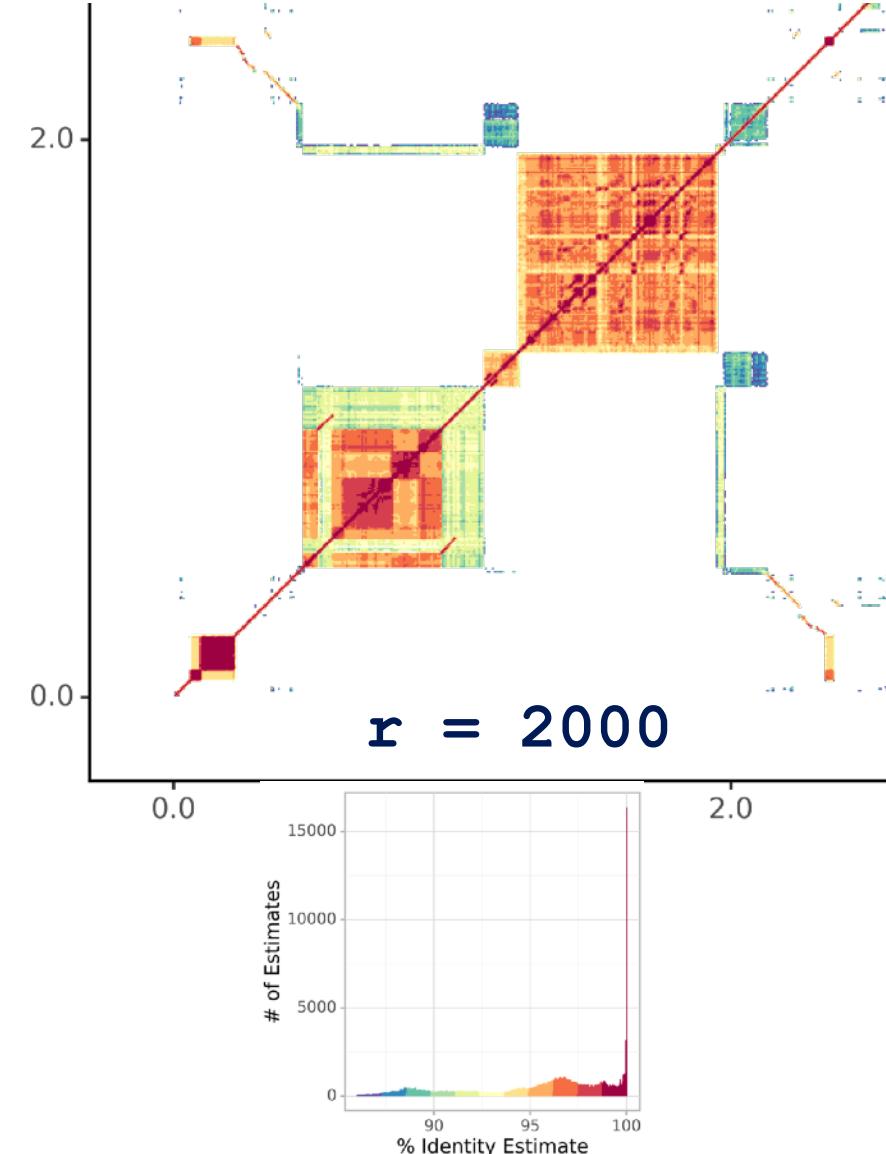
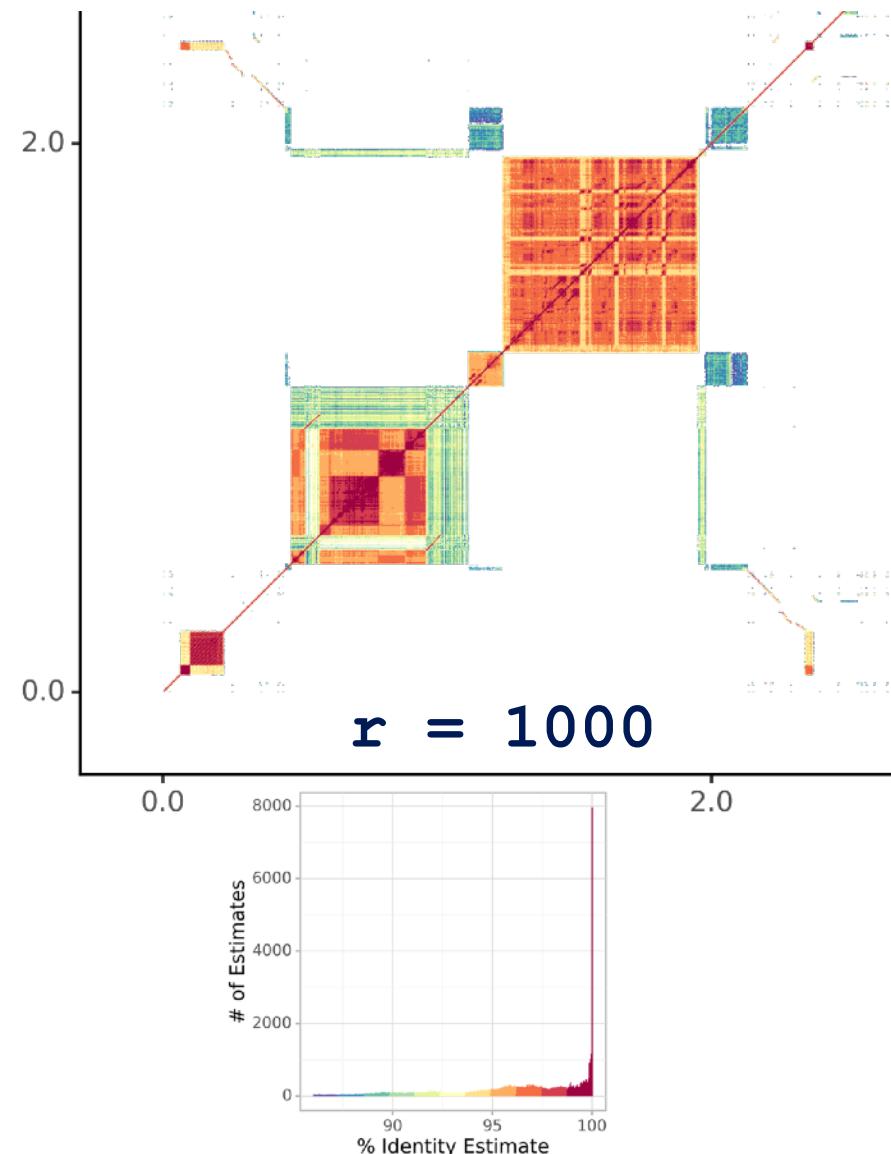
- Parameters to be familiar with:

- **n: Sequence size.** *Not a parameter you can control through ModDotPlot!
- **-r/--resolution: Number of intervals (plot resolution).** Default: 1000.
- **-w/--window: Interval length.** *Note this is inversely proportional to resolution!

```
Computing self identity matrix for Chr21...
Sequence length n: 6000000
Window size w: 6000
Modimizer sketch size: 1500
Plot Resolution r: 1000
Progress: |██████████| 100.0% Completed
Saved bed file to chr21_r1000/Chr21.bed
real    0m36.850s
user    0m36.594s
sys     0m0.784s
(venv) gitpod /workspace $
```

```
Computing self identity matrix for Chr21...
Sequence length n: 6000000
Window size w: 3000
Modimizer sketch size: 1500
Plot Resolution r: 2000
Progress: |██████████| 100.0% Completed
Saved bed file to chr21_r2000/Chr21.bed
real    1m44.242s
user    1m42.177s
sys     0m2.465s
```

Running ModDotPlot - Static Mode



Running ModDotPlot - Static Mode

- Other parameters:

- **-k / --kmer:** *k*-mer length. Default: 21
- **--compare:** Will produce a comparative plot if 2+ fasta/bed files input
- **--grid:** Produce an $n \times n$ grid, used with compare. *Limited to 6 sequences!

```
Computing pairwise identity matrix for Chr15 and chr1:14M-18M...
```

```
Sequence length Chr15: 6000000
```

```
Sequence length chr1:14M-18M: 4000000
```

```
Window size w: 6000
```

```
Modimizer sketch size: 1500
```

```
Plot Resolution r: 1000
```

```
Progress: |██████████| 100.0% Completed
```

Running ModDotPlot - Interactive Mode

- Most parameters stay the same! Different ones:
 - `--port`: Port number to show Dash on. Default: 8050
 - `--save`: Save matrices to folder (-o)
 - `--load`: Load a save matrix hierarchy.

Acknowledgements



Phillippy Lab - NHGRI



AP SK AR NH



BP AS DA BW SS



Schatz Lab - JHU



BGA24

Thanks Damon!

