## Supplemental Material

---

**Algorithm 1** Partition Sequence S into sets of modimizers

---

**Input:** List of hashed $k$-mers $S_k[x_1, x_2, ...x_n]$, window size $w$, sparsity $s$
**Output:** List of modimizer sets $MOD_s = [A_1, A_2, ..., A_r]$

1: **function** BASELAYER($S_k, w, s$)
2:      $n \leftarrow$ length of $S_k$
3:      $r \leftarrow \lceil \frac{n}{w} \rceil$              ▷ Set number of windows (resolution) based on $n$ and $w$
4:      $m \leftarrow \frac{w}{s}$              ▷ Set the expected modimizer density
5:      $MOD_s \leftarrow$ list of size $r$
6:      **for** $i \leftarrow 0$ to $r$ **do**              ▷ Retrieve all modimizers within each interval
7:          $MOD_s[i] \leftarrow \{\}$
8:          $start \leftarrow wi$              ▷ Set non-overlapping interval boundaries
9:          $end \leftarrow \min((start + w), n - 1)$
10:         $MOD_s[i] \leftarrow$ GETMODIMIZERS($S_k[start : end], s, d$)      ▷ Populate list with sets of modimizers
11:      **end for**
12:      **return** $MOD_s$
13: **end function**

14: **function** GETMODIMIZERS($S_k, s, m$)
15:      $A \leftarrow \forall x \in \{S_k[A]\} : x \equiv 0 \bmod s$          ▷ Gather the set of unique modimizers per interval
16:      $\hat{s} \leftarrow \frac{s}{2}$
17:      **while** $|A| < \frac{m}{2}$ and $\hat{s} > 1$ **do**      ▷ Resample at higher density if number of modimizers is under threshold
18:          $A \leftarrow \forall x \in \{S_k[A]\} : x \equiv 0 \bmod \hat{s}$
19:          $\hat{s} \leftarrow \frac{\hat{s}}{2}$
20:      **end while**
21:      **return** $A$
22: **end function**

---

---

**Algorithm 2** Partition Sequence S into a modimizer hierarchy H

---

**Input:** List of hashed $k$-mers $S_k[x_1, x_2, ...x_n]$, minimum window size $\hat{w}$, sparsity $\hat{s}$, resolution $r$
**Output:** List $H = [MOD_{\hat{s}}, MOD_{2\hat{s}}, ..., MOD_{2^{l-1}\hat{s}}]$

1: **function** BUILDHIERARCHY($S_k, \hat{w}, \hat{s}, r$)
2:      $n \leftarrow$ length of $S_k$
3:      $l \leftarrow \lfloor \log_2 \left( \frac{n}{\hat{w}r} \right) \rfloor$      ▷ Initialize number of layers based on min. window size and resolution
4:      $H \leftarrow$ list of size $l$
5:      $H[0] \leftarrow$ BASELAYER($S_k, \hat{w}, \hat{s}$)              ▷ Compute bottom layer
6:      **for** $i \leftarrow 1$ to $l$ **do**          ▷ Iteratively compute subsequent layers from previous layer
7:          $\hat{r} \leftarrow 2^{l-1-i}r$          ▷ Halve the resolution when building subsequent layer
8:          $H[i] \leftarrow$ ADDLAYER($H[i-1], 2^i \hat{s}, 2^i \hat{w}, n, \hat{r}$)
9:      **end for**
10:     **return** $H$
11: **end function**

12: **function** ADDLAYER($MOD_{\hat{s}}[A_1, A_2, ..., A_r], s, w, n, \hat{r}$)
13:      $MOD_s \leftarrow$ list of size $\hat{r}$              ▷ Initialize current layer
14:      $m \leftarrow \frac{w}{s}$              ▷ Update expected sketch size for current layer
15:      **for** $i \leftarrow 0$ to $\hat{r}$ **do**      ▷ Retrieve modimizers from matching intervals within the previous layer
16:          $MOD_s[i] =$ GETMODIMIZERS($(MOD_{\hat{s}}[A_{2i}] \cup MOD_{\hat{s}}[A_{2i+1}]), s, m$)
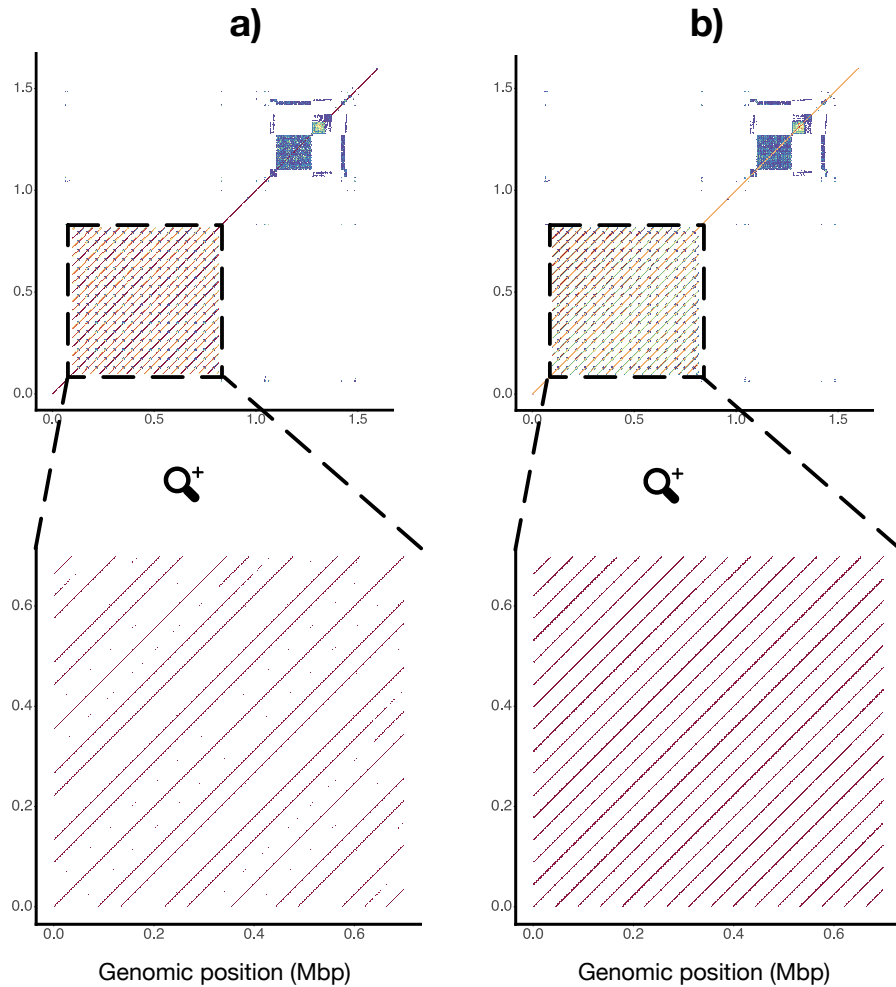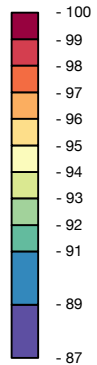17:      **end for**
18:      **return** $MOD_s$
19: **end function**

---

## References

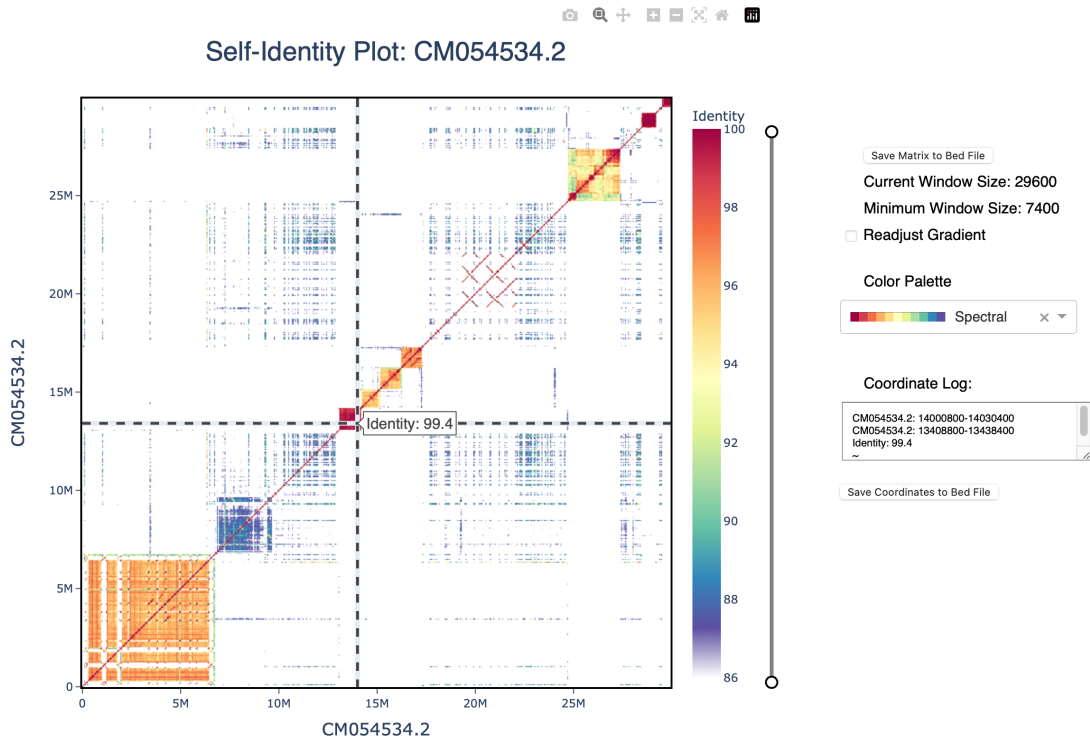**chr14:**
**2,000,000-3,600,000**

% Identity Estimate

**a)**

**b)**

**chr14:**
**2,100,000-2,800,000**

% Identity Estimate

Genomic position (Mbp)
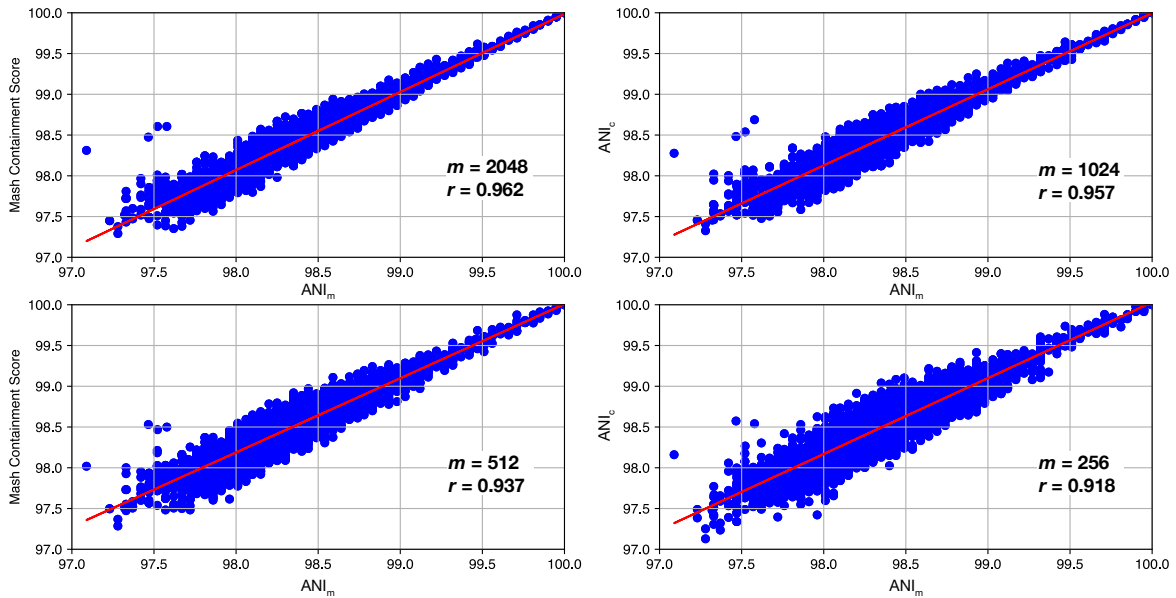
**Supplementary Figure 1.** Screenshots of ModDotPlot run on a human acrocentric short arm (CHM13 chr14:2,000,000-3,600,000), highlighting the rDNA array with and without ModDotPlot's interval extension activated. **a)** With no interval extension, the full 16 copy rDNA array is visible; however, when zoomed in (chr14:2,100,000-2,800,000) and filtered for >98% sequence identity, some rDNA copies disappear from the plot due to registration artifacts. **b)** With intervals extended (i.e. when computing similarity for the cell M(A,B), interval B is extended by $w/2$ in both directions), all rDNA copies appear at all zoom levels.

**Supplementary Figure 2.** Screenshot of ModDotPlot's interactive mode, showcasing the entire Y chromosome of a gibbon (mSym-Syn1, Makova *et. al* 2024)). Despite spanning almost 30 Mbp, ModDotPlot was able to create 3 matrices in under 2 minutes, with around 2.5GB of memory. Screenshot was taken using ModDotPlot version v0.8.0 (git commit `ed190c7`).



**Supplementary Figure 3.** Scatterplot of $ANI_m$ against the $ANI_c$ distances under various sketch sizes. The 10-mer non-canonical HOR region from Figure 4 was excluded from analysis.