

1 Enrichment of non-B-form DNA at *D. melanogaster* centromeres

2 Venkata S. P. Patchigolla¹ and Barbara G. Mellone^{1,2}

3 1. Department of Molecular and Cell Biology, University of Connecticut, Storrs CT, 06269, USA

4 2. Institute for Systems Genomics, University of Connecticut, Storrs CT, 06269, USA

6 Abstract

7 Centromeres are essential chromosomal regions that mediate the accurate inheritance of genetic
8 information during eukaryotic cell division. Despite their conserved function, centromeres do not contain
9 conserved DNA sequences and are instead epigenetically marked by the presence of the centromere-specific
10 histone H3 variant CENP-A (centromeric protein A). The functional contribution of centromeric DNA sequences
11 to centromere identity remains elusive. Previous work found that dyad symmetries with a propensity to adopt
12 non-canonical secondary DNA structures are enriched at the centromeres of several species. These findings lead
13 to the proposal that such non-canonical DNA secondary structures may contribute to centromere specification.
14 Here, we analyze the predicted secondary structures of the recently identified centromere DNA sequences from
15 *Drosophila melanogaster*. Although dyad symmetries are only enriched on the Y centromere, we find that other
16 types of non-canonical DNA structures, including DNA melting and G-quadruplexes, are common features of all
17 *D. melanogaster* centromeres. Our work is consistent with previous models suggesting that non-canonical DNA
18 secondary structures may be conserved features of centromeres with possible implications for centromere
19 specification.

21 Introduction

22 Eukaryotes share a common mechanism to faithfully segregate genetic information during each cell cycle by
23 which chromosomes are attached to microtubule fibers and are physically pulled towards opposite poles by the
24 kinetochore. Centromeres are essential chromosomal regions that specify the site for the assembly of the
25 kinetochore and are epigenetically marked by chromatin enriched in the histone H3 variant centromeric protein
26 A (CENP-A). CENP-A has been shown to be sufficient for kinetochore assembly and *de novo* recruitment of CENP-
27 A in *D. melanogaster* somatic cells (Chen et al., 2014; Mendiburo et al., 2011; Palladino et al., 2020). Despite
28 their conserved and essential function, centromeres are among the most rapidly evolving regions of genomes
29 (Melters et al., 2013). This rapid evolution has been proposed to be a result of intra-genomic conflict whereby
30 centromeres act as selfish genetic elements driving the rapid evolution of centromeric proteins (Henikoff et al.,
31 2001; Malik and Henikoff, 2009). Furthermore, in organisms such as fungi, nematodes, insects, plants, and
32 vertebrates, centromere function is largely independent of the presence of centromeric DNA sequences, relying
33 instead on the presence of CENP-A chromatin (reviewed in (McKinley and Cheeseman, 2016)). Thus, for most
34 species, the functional significance of centromeric DNA sequences in dictating (or at least contributing to)
35 centromere identity remains unclear.

36 In an effort to identify genetic characteristics shared amongst the centromeres of diverse eukaryotes,
37 Kasinathan et al. (Kasinathan and Henikoff, 2017) surveyed centromeric DNA sequences from mouse, chicken, *S.*
38 *pombe* and humans for the presence of <10-bp dyad symmetries (a.k.a. inverted repeats), which are known to
39 adopt unconventional secondary structures such as stem-loops or cruciform extrusions. The authors found that
40 the centromeres of species such as the African Green monkey, chicken, and the fission yeast *S. pombe* were
41 enriched in these motifs. Centromeres enriched in dyad symmetries also showed a predicted propensity to form
42 non-canonical secondary DNA structure under stress, such as that resulting from DNA supercoiling caused by

43 transcription or replication. Non-canonical DNA structures are known as non-B-form DNA and collectively
44 represent any deviation from double stranded B-DNA (the right-handed helix with 10-nt per turn). High
45 likelihood of predicted cruciforms correlated with enrichment in dyad symmetries and other structures, such as
46 melt DNA, were also predicted for some species. Interestingly, centromeres devoid of dyad symmetries, such as
47 those of humans, contain binding sites for CENP-B, a protein that binds specifically to CENP-B box DNA motifs
48 found within α -satellite (Verdaasdonk and Bloom, 2011). CENP-B binding results in the bending of DNA (Tanaka
49 et al., 2001), which in itself represents another non-canonical DNA structure. Based on these analyses, the
50 authors proposed that non-canonical secondary structures may have been selected for during centromere
51 evolution, with a possible role as a structural cue for centromere specification (Kasinathan and Henikoff, 2017).
52 Various non-B structures such as hairpins (Jonstrup et al., 2008), R-loops (Kabeche et al., 2018) and i-motifs
53 (Garavis et al., 2015a; Garavis et al., 2015b) have been observed *in vitro* and *in vivo*, consistent with this model.
54 How widespread centromeric non-B-DNA structures across species may be remains unknown.

55 The centromeres of *D. melanogaster* were not identified and characterized until recently through a
56 combination of long-read sequencing, chromatin immunoprecipitation, and OligoPaints Fluorescence In-Situ
57 Hybridization (FISH). Chang et al. identified five contigs that make up at least part of the centromeres (Chang et
58 al., 2019) (**Fig. 1A**). The contigs for centromeres X, 3 and 4 contain an island of complex DNA enriched in
59 retroelements flanked by simple satellite repeats. For centromere 2, only a short contig was identified, which
60 contains a small island with a single truncated retroelement flanked by simple satellites. Lastly, the contig for the
61 Y centromere consists of a large island and no satellite DNA. FISH on mitotic chromosomes and extended
62 chromatin fibers show that for centromeres X, 2 and 4, the CENP-A domain spans a region larger than the contig
63 itself, which, based on cytological analyses, can be inferred to be made up of unassembled simple satellites
64 (Chang et al., 2019).

65 Here, we use several prediction algorithms to survey the presence of non-B-DNA-form at the centromeres of
66 *D. melanogaster*. Although we show that inverted repeats and cruciform extrusions are not a predominant
67 feature at *D. melanogaster* centromeres, we find evidence for the enrichment of other predicted non-canonical
68 secondary structures such as melted DNA and G-quadruplexes.

69

70 **Results and discussion**

71 **Dyad Symmetries are not common features of *D. melanogaster* centromeres**

72 To determine if *D. melanogaster* centromeres are enriched in <10-bp DNA dyad symmetries as previously
73 reported for the centromeres of other species (Kasinathan and Henikoff, 2017), we used the program
74 Palindrome from the EMBOSS suite. We used five contigs (one for each of the X, 2, 3, 4 and Y chromosomes)
75 that are highly enriched in CENP-A chromatin immunoprecipitations and were confirmed to be associated with
76 CENP-A using OligoPaint FISH on extended chromatin fibers as the bona fide *D. melanogaster* centromeres
77 (Chang et al., 2019) (**Fig. 1A**). For our controls, we used several composition and length-matched random
78 genomic sequences for each of the centromere contigs (see Methods). We plotted the EMBOSS palindrome
79 output by calculating the dyad density, obtained by adding the number of base pairs that are part of a dyad
80 divided by the sequence length, and found that only the Y centromere displays dyad symmetry densities higher
81 than control average (**Fig. 1B-G**). These analyses suggest that dyad symmetries are not major features of *D.*
82 *melanogaster* centromeres and thus are unlikely to play a role in centromere specification. A lack of dyad
83 symmetries was previously reported for human, great apes and *M. musculus* centromeres (Kasinathan and
84 Henikoff, 2017).

85 **Enrichment of predicted non-B-form DNA structures at centromeric contigs using SIST**

86 The EMBOSS palindrome algorithm identifies dyad symmetries based on sequence analysis. However, this
87 algorithm does not take into account the predicted thermodynamics of DNA and thus does not provide
88 information on the secondary structures it is likely to adopt. Superhelical transitions occur in DNA when negative
89 supercoiling drives susceptible regions to acquire forms alternative to native B-DNA that are energetically
90 favorable. To determine if centromeres are susceptible to adopt non-B-form DNA, we used a computational
91 algorithm that models stress-induced structural transitions (SIST) for multiple non-canonical DNA secondary
92 structures: Z-DNA, DNA melting (*i.e.* strand separation), and cruciform extrusions (Zhabinskaya et al., 2015). SIST
93 was previously used by Kasinathan et al. to show higher probability to adopt non-B-form DNA for centromeres
94 enriched in dyad symmetries (Kasinathan and Henikoff, 2017).

95 We ran segments of DNA in 5,000-bp blocks every 2,500-bp and took the maximum values for the
96 overlapping regions whenever different. DNA transitions depend on temperature; since *D. melanogaster* is an
97 ectotherm species, we ran SIST at five different temperatures at which *D. melanogaster* may be found (18°C,
98 22°C, 25°C, 30°C and 35°C) and determined enrichment probabilities for centromeres compared to their
99 respective control regions. The probability of Z-DNA formation, which has not been previously analyzed for
100 centromeres, is lower than controls for each of the centromeres irrespectively of the temperature and thus is
101 unlikely to be associated with centromeres (**Fig. 2A**). As for cruciforms, only the centromere of the Y
102 chromosome shows higher probability than controls at all temperatures (**Fig. 2B**). These findings are consistent
103 with the observation that the Y is the only centromere showing an enrichment of inverted repeats (**Fig. 1F**),
104 which are thought to adopt cruciform extrusions (Hamer and Thomas, 1974; Leach, 1994). Our findings in
105 *Drosophila* are consistent with previous analyses on the centromeres of fission yeast, African green monkey and
106 on human neocentromeres, where the probability of DNA melting was found to be higher than that of controls
107 (Kasinathan and Henikoff, 2017).

108 (Kasinathan and Henikoff, 2017). Interestingly, at 25°C and 30°C, all of the centromeres have higher
109 probability than controls for DNA melting (melt). Centromere 2 and 4 display higher melting probability than
110 controls also at 35°C. The Y displays higher DNA melting probability than controls at all temperatures greater
111 than 22°C. At 18°C, none of the centromeres displays higher probability of DNA melting (**Fig. 2C**). When we
112 plotted the overall probability of forming all three types of non-B DNA, we noticed that it increases with higher
113 temperatures (**Fig. 2D**); this is likely due, at least in part, to the contribution of DNA melting to this probability.
114 Cell and organism growth are regulated by temperature and the temperatures at which different organisms
115 thrive are vastly different across eukaryotic species. If the ability of centromeres to adopt non-B DNA
116 conformations needed for proper centromere function during cell division is also affected by temperature, this
117 could be a factor under selection during evolution, contributing to the diversity of centromeric DNA sequences
118 observed across lineages.

119 DNA melting is accurately predicted at actively transcribed regions that display strand separation *in vivo*
120 (Zhabinskaya et al., 2015). As centromeres from across species have been shown to display low transcriptional
121 activity (reviewed in (Mellone and Fachinetti, 2021)), the enrichment for this particular non-canonical DNA
122 structure is especially interesting. DNA melting may facilitate transcription, which in turn could facilitate histone
123 turnover or the formation of secondary DNA/RNA structures at centromeres, contributing to centromere
124 specification (Kasinathan and Henikoff, 2017; Talbert and Henikoff, 2020).

125

126 **Enrichment of non-B-form DNA in centromeric contigs using GQuad**

127 Previous work proposed that non-B-form DNA may be an evolutionary conserved signature required for
128 centromere specification. Yet, aside from the Y centromere, which is enriched in inverted repeats and has higher
129 probability of forming cruciforms than controls (**Fig. 1F** and **2B**), all other *D. melanogaster* centromeres show
130 higher probability than controls only for DNA melting. As SIST only predicts 3 types of non canonical DNA
131 structures, we wanted to expand our analysis to additional non-B-form DNA types. For this purpose, we used
132 Gquad, a package that can predict 7 different non-B DNA structures: a-phased DNA repeats, G-quadruplexes,
133 intramolecular triplexes (H-DNA), slipped DNA, short tandem repeats (STR,) triplex forming oligonucleotides
134 (TFO), and Z-DNA. Gquad provides the positions and probability for specific non-B-form DNA using scores
135 ranging from one asterisk (low likelihood) to three asterisks (high likelihood). In the absence of experimental
136 data identifying non-B-form DNA and of a non-B-form DNA database for *D. melanogaster*, sequences known to
137 form non-B-form DNA are not available as positive controls to determine the accuracy of our predictions. A
138 previous study used inter-pulse duration (IPD) values (*i.e.* the time it takes to add a nucleotide during single-
139 molecule sequencing) from PacBio long-read sequencing data to infer non-B-form DNA (Guiblet et al., 2018).
140 When we plotted the average IPD values of regions predicted to form non-B-DNA (*e.g.* Z-DNA) identified by
141 Gquad with a likeliness of two asterisks or greater in a 300-bp window centered on the sequence predicted to
142 form Z-DNA, we observed IPD values that were twice as high, suggesting that the predictions generated by
143 Gquad are accurate (**Fig. 3A**). Next, we calculated all the likelihoods for each type of non-B-DNA and combined
144 them such that if a particular base pair was predicted to form non-B-form DNA of more than one type, the
145 likeliness of the two were added together. To determine the significance of enrichment we used the two-sample
146 Kolmogorov-Smirnov (KS) test. Through this analysis, we find that all centromeres are significantly enriched for
147 non-B-DNA (**Fig. 3B-F**). Since the values for the 7 types of non-B-DNA are combined in this analysis, we next
148 wanted to determine which types of non-B-DNA are contributing most to the enrichment of non-B form DNA at
149 the centromeres found with Gquad. For this, we analyzed the enrichment of individual type and find that of the
150 7 non-canonical DNA forms, the ones that contribute the most are slipped DNA, STR, and G-quadruplexes (**Fig.**
151 **3G**).

152 Next, we sought to determine which types of repeats are contributing most to the likelihood of adopting
153 non-canonical DNA secondary structures by ranking the average Gquad values for all repeats in the *D.*
154 *melanogaster* genome. We find that simple satellite DNAs contribute the most, as they are consistently ranked
155 higher than other elements (**Table S1**). Short satellites are known to be prone to form non-canonical DNA
156 structures, particularly slipped DNA (Sinden et al., 2007). If centromeres need to be marked by unconventional
157 DNA structures in order to function or be stable, a potential explanation for why satellite DNA is found at many
158 regional centromeres across species could be that it can adopt non-B DNA.

159 To determine the prevalence of non-B-DNA at centromeric contigs compared to the rest of the genome
160 (irrespective of GC content), we ranked all contigs that make up the genome based on the average Gquad
161 likelihood. We find that all centromeric contigs fall within the top 37% of the 180 contigs picked up by Gquad as
162 containing some form on non-B DNA, with centromeres X, 2 and 4 ranking 6th, 15th and 22nd, respectively
163 (**Table S2**). These findings indicate that, although the centromeres may not rank the highest, they are among the
164 most likely sequences in the genome to form non-B-DNA.

166 **G-quadruplexes are common features of *D. melanogaster* centromeres**

167 To confirm our prediction of G-quadruplexes at the centromeres with an additional algorithm, we used
168 G4Hunter, a more recent program that gives a G-quadruplex propensity score as output. Unlike Gquad,

169 G4hunter takes into account G-richness and G-skewness of a given sequence. Furthermore, this algorithm was
170 validated on published sequences known to form G-quadruplexes as well as with biophysical methods (Bedrat et
171 al., 2016). We ran G4Hunter using a stringent threshold value of 1.5 and found that all centromeres, except the
172 3 and X centromeres, are enriched in G-quadruplexes compared to their respective controls (**Fig. 4A-E**). Having
173 observed enrichment of G-quadruplexes with two independent methods, we conclude that G-quadruplexes are
174 likely to be common features of *D. melanogaster* centromeres. G-quadruplexes play a role in transcriptional
175 regulation, translation and replication (Bedrat et al., 2016). One possibility is that the higher prevalence of G-
176 quadruplexes at the centromeres may contribute to centromere transcription homeostasis.

177 Collectively, our computational predictions suggest that *D. melanogaster* centromeres are enriched in non-B
178 DNA secondary structures. Our findings are consistent with the model that non-canonical DNA forms may be
179 evolutionarily conserved features of centromeres with possible functions in centromere specification. Under
180 such paradigm, the only feature under selection at centromeres would be their secondary DNA structure. Since
181 a myriad of primary DNA sequence combinations can accommodate such secondary conformations, such
182 mechanism would enable ample opportunity for adaptation under intra-genomic conflict (Kasinathan and
183 Henikoff, 2017).

184

185 **Acknowledgments**

186 We wish to thank Patrick Grady, Sivakanthan Kasinathan, Asna Amjad, Marzia Cremona, Kateryna Makova,
187 Amanda Larracuenta, and Emmy Karim for discussions and suggestions and the UConn HPC and CBC for
188 computing resources. This work was funded by National Institute of Health grant R35GM131868 to BGM and a
189 Beckman Scholar fellowship to VSPP.

190

191 **Figure Legends**

192 **Figure 1: Dyad symmetries are not common features of *D. melanogaster* centromeres.** (A) Schematic of
193 the DNA organization of *D. melanogaster* centromere contigs. (B-F) Dyad symmetry density plots for *D.*
194 *melanogaster* centromeres. Only the Y contig (Y_Contig26; yellow box) showed a significant enrichment. $P < 0.05$,
195 one-sample t-test. (G) Example of inverted repeats from the Y centromere contig (base pairs 181–390).

196

197 **Figure 2: Enrichment of predicted non-B-form DNA at centromeres contigs using SIST.** Diagram
198 summarizing the SIST output. Results for Z-DNA (A), cruciform (B), and melt DNA (melt) (C) are shown for each of
199 the centromeres at five different temperatures ($^{\circ}\text{C}$). Different colors represent significance as outlined in the
200 legend. (D) Average probability of non-B DNA formation for each centromere contig at different temperatures.

201

202 **Figure 3. Enrichment of predicted non-B-form DNA in centromeric contigs using GQuad.** (A) Plot showing
203 the average IPD value for sequences predicted to form Z-DNA by GQuad with a likelihood of greater than two
204 asterisks (see text for details). Z-DNA is centered around 150-bp. (B-F) Data distribution of likelihoods for each of
205 the centromeres as a combination of all non-B DNA predicted by Gquad. Asterisks represent $p < 0.05$ (KS test). (G)
206 Pie chart showing the relative contributions of different non-B DNA types identified by Gquad.

207

208 **Figure 4: G-quadruplexes are common predicted features of *D. melanogaster* centromeres.** (A-E) Graphs of
209 the average G-quadruplex density for each centromere contig predicted by G4Hunter. Asterisks represent p
210 < 0.05 (One-sample t-test). Note that several control regions were not predicted to form any G-quadruplexes.

211 **Supplemental material**

212

213 **Supplemental data**

214 **Table S1. Table ranking the average Gquad value for all repeats in the *D. melanogaster* genome.** Repeats
215 associated with centromere contigs are highlighted in yellow.

216 **Table S2. Table ranking all contigs that make up the genome based on the average Gquad likelihood.** Only
217 contigs with an assigned likelihood are included (180 out of 190 total contigs in the genome). Centromeric
218 contigs are highlighted in yellow

219

220 **Methods**

221 **Genome data**

222 The genome used in this paper is from Chang and Larracunte 2019 (Chang and Larracunte, 2019). The
223 centromere contigs used for this analysis were Contig79 for centromere X, Contig119 for centromere 4,
224 Y_Contig26 for centromere Y, Contig 3R_5 for centromere 3 and tig00057289 for centromere 2 (Chang et al.,
225 2019).

226

227 **Source code**

228 Code used to perform the analysis in this manuscript is available from GitHub
229 (<https://github.com/venkata14/dmel-nonb>).

230

231 **Generation of controls regions**

232 The controls used for the analysis were 50 random segments of the genome that are both the same size and
233 have a similar GC content within 10% as the respective centromeric contig. A maximum of two controls with a
234 50,000-bp overlap was allowed.

235

236 **Detection of dyad symmetries using EMBOSS palindrome**

237 EMBOSS Palindrome (<https://www.bioinformatics.nl/cgi-bin/emboss/help/palindrome>) was used to detect
238 dyad symmetries with the minimum palindrome being 5, the maximum palindrome being 100, allowing a gap
239 limit of 20 and allowing overlapping dyad symmetries. We analyzed the output by calculating the dyad density,
240 which we defined as the sum of the lengths of all palindromic regions identified by Palindrome divided by the
241 length of the entire contig containing it. that contain that position. For a sequence, the length-normalized dyad
242 density was defined as the sum of the values for each position divided by the sequence length.

243

244 **Prediction of Z-DNA, DNA melting and cruciform transitions using SIST**

245 The probabilities of Z-DNA, Cruciform transitions and DNA melting were predicted using SIST (Zhabinskaya et
246 al 2015) as described in Kasinathan et al. (Kasinathan and Henikoff, 2017). We used default parameters with the
247 algorithm type "A" which uses the trans_compete C++ codes along with five different temperatures: 18°C, 22°C,
248 25°C, 30°C, 35°C for this analysis. For sequences greater than 10kb in length, we slid a 5,000-bp window in
249 2,500-bp steps and analyzed these sub-sequences using SIST. The SIST predictions were then reassembled by
250 taking the maximum SIST value for any given base pair.

251 To determine the the average probability of non-B-DNA formation for each temperature for all centromeres,
252 we added the average value of Z-DNA, cruciform, and melt formation at each temperature.

253

254 **Prediction of non-B-DNA using Gquad**

255 Gquad (v2.2-1; <https://cran.r-project.org/web/packages/gquad/gquad.pdf>) consists of multiple R packages
256 that predict individual forms of non-B-DNA. We ran R packages on the heterochromatin-enriched *D.*
257 *melanogaster* genome (Chang and Larracunte, 2019) for the 7 types of non-B-DNA: a phased DNA, G-
258 quadruplexes, H-DNA, slipped DNA, Short Tandem Repeats (STR), Triplex Forming Oligonucleotides (TFO), and Z-
259 DNA. The packages output likelihoods for each nucleotide from a range of one to three asterisks representing
260 the likelihood of non-B-DNA formations. For those that did not output a likelihood, we used 2 asterisks as the
261 default likelihood value. We then analyzed the data by combining all likelihoods for the 7 types of non-B-DNA for
262 a respective sequence such that if there were overlaps in likelihoods of two different non-B-DNA types, we
263 added those likelihoods together. This results in an array where each position is a summation of all likelihoods
264 for a particular base pair.

266 **Identifying relative amounts of non-B-DNA using Gquad**

267 Using the Gquad R package, we ran the package on the heterochromatin-enriched *D. melanogaster* genome
268 (Chang and Larracunte, 2019) for the 7 types of non-B-DNA as similar to above. We then added all the positions
269 predicted to form non-B-DNA for each of the 7 types and created a pie chart. To determine significance of
270 prevalence between specific types of non-B-DNA in the centromere versus the controls, we used the one sample
271 t-test on the average centromeric value and the control values for each respective non-B-DNA.

273 **Prediction of G-Quadruplexes using G4Hunter**

274 G4Hunter (<https://www.bioinformatics.nl/cgi-bin/emboss/help/palindrome>) was run using a window size of
275 25 base pairs and threshold values of 1 and 1.5. The program outputs the positions of the nucleotides that are
276 predicted to form G-Quadruplexes. Using these positions, we calculated the density of G-Quadruplexes by taking
277 the total number of nucleotides predicted to form G-Quadruplexes and dividing them by the total number of
278 nucleotides in the respective sequence.

280 **Validating non-B-DNA predictions of Gquad using IPDs**

281 Publicly available PacBio sequencing reads from *D. melanogaster* (Kin et al 2014) were aligned to the
282 heterochromatin-enriched *D. melanogaster* genome (Chang et al. 2019) with pbalign (SMRT v7.0), and IPDs were
283 computed at nucleotide resolution with ipdSummary.py using the P5C3 chemistry
284 (<https://github.com/PacificBiosciences/kineticsTools/tree/master/kineticsTools>). This outputs an IPD value
285 which is an average of 3 IPD subheads values per nucleotide. All normalization of intermolecular variability and
286 trimming for outliers was done automatically. Then, using the positive strand, all regions predicted to be Z-DNA
287 by Gquad with a likelihood of two asterisks or higher were extracted in 300 base pair windows. The IPDs values
288 of these sequences were extracted such that the predicted sequence to form Z-DNA was centered. All windows
289 with no IPD values were filtered out, after which the IPD values of all sequences were averaged lengthwise and
290 plotted.

292 **Statistical tests**

293 The two-sample Kolmogorov–Smirnov test was used to compare distributions of SIST and GQuad likelihood
294 values. One sample t-test was used for both the dyad density and G4Hunter distributions.

295 References

296

- 297 1. Bedrat, A., L. Lacroix, and J.L. Mergny. 2016. Re-evaluation of G-quadruplex propensity with G4Hunter.
298 *Nucleic Acids Res.* 44:1746-1759.
- 299 2. Chang, C.H., A. Chavan, J. Palladino, X. Wei, N.M.C. Martins, B. Santinello, C.C. Chen, J. Erceg, B.J. Beliveau,
300 C.T. Wu, A.M. Larracuente, and B.G. Mellone. 2019. Islands of retroelements are major components of
301 *Drosophila* centromeres. *PLoS Biol.* 17:e3000241.
- 302 3. Chang, C.H., and A.M. Larracuente. 2019. Heterochromatin-Enriched Assemblies Reveal the Sequence and
303 Organization of the *Drosophila melanogaster* Y Chromosome. *Genetics.* 211:333-348.
- 304 4. Chen, C.C., M.L. Dechassa, E. Bettini, M.B. Ledoux, C. Belisario, P. Heun, K. Luger, and B.G. Mellone. 2014.
305 CAL1 is the *Drosophila* CENP-A assembly factor. *J Cell Biol.* 204:313-329.
- 306 5. Garavis, M., N. Escaja, V. Gabelica, A. Villasante, and C. Gonzalez. 2015a. Centromeric Alpha-Satellite DNA
307 Adopts Dimeric i-Motif Structures Capped by AT Hoogsteen Base Pairs. *Chemistry.* 21:9816-9824.
- 308 6. Garavis, M., M. Mendez-Lago, V. Gabelica, S.L. Whitehead, C. Gonzalez, and A. Villasante. 2015b. The
309 structure of an endogenous *Drosophila* centromere reveals the prevalence of tandemly repeated sequences
310 able to form i-motifs. *Sci Rep.* 5:13307.
- 311 7. Guiblet, W.M., M.A. Cremona, M. Cechova, R.S. Harris, I. Kejnovska, E. Kejnovsky, K. Eckert, F. Chiaromonte,
312 and K.D. Makova. 2018. Long-read sequencing technology indicates genome-wide effects of non-B DNA on
313 polymerization speed and error rate. *Genome Res.* 28:1767-1778.
- 314 8. Hamer, D.H., and C.A. Thomas, Jr. 1974. Palindrome theory. *J Mol Biol.* 84:139-144.
- 315 9. Henikoff, S., K. Ahmad, and H.S. Malik. 2001. The centromere paradox: stable inheritance with rapidly
316 evolving DNA. *Science.* 293:1098-1102.
- 317 10. Jonstrup, A.T., T. Thomsen, Y. Wang, B.R. Knudsen, J. Koch, and A.H. Andersen. 2008. Hairpin structures
318 formed by alpha satellite DNA of human centromeres are cleaved by human topoisomerase IIalpha. *Nucleic*
319 *Acids Res.* 36:6165-6174.
- 320 11. Kabeche, L., H.D. Nguyen, R. Buisson, and L. Zou. 2018. A mitosis-specific and R loop-driven ATR pathway
321 promotes faithful chromosome segregation. *Science.* 359:108-114.
- 322 12. Kasinathan, S., and S. Henikoff. 2017. Non-B-Form DNA Is Enriched at Centromeres. *Mol Biol Evol.* 35:949-
323 962.
- 324 13. Leach, D.R. 1994. Long DNA palindromes, cruciform structures, genetic instability and secondary structure
325 repair. *Bioessays.* 16:893-900.
- 326 14. Malik, H.S., and S. Henikoff. 2009. Major evolutionary transitions in centromere complexity. *Cell.* 138:1067-
327 1082.
- 328 15. McKinley, K.L., and I.M. Cheeseman. 2016. The molecular basis for centromere identity and function. *Nat*
329 *Rev Mol Cell Biol.* 17:16-29.
- 330 16. Mellone, B.G., and D. Fachinetti. 2021. Diverse mechanisms of centromere specification. *Curr Biol.*
331 31:R1491-R1504.
- 332 17. Melters, D.P., K.R. Bradnam, H.A. Young, N. Telis, M.R. May, J.G. Ruby, R. Sebra, P. Peluso, J. Eid, D. Rank, J.F.
333 Garcia, J.L. DeRisi, T. Smith, C. Tobias, J. Ross-Ibarra, I. Korf, and S.W. Chan. 2013. Comparative analysis of
334 tandem repeats from hundreds of species reveals unique insights into centromere evolution. *Genome Biol.*
335 14:R10.
- 336 18. Mendiburo, M.J., J. Padeken, S. Fulop, A. Schepers, and P. Heun. 2011. *Drosophila* CENH3 is sufficient for
337 centromere formation. *Science.* 334:686-690.
- 338 19. Palladino, J., A. Chavan, A. Sposato, T.D. Mason, and B.G. Mellone. 2020. Targeted De Novo Centromere
339 Formation in *Drosophila* Reveals Plasticity and Maintenance Potential of CENP-A Chromatin. *Dev Cell.*
340 53:129.
- 341 20. Sinden, R.R., M.J. Pytlos-Sinden, and V.N. Potaman. 2007. Slipped strand DNA structures. *Front Biosci.*
342 12:4788-4799.
- 343 21. Talbert, P.B., and S. Henikoff. 2020. What makes a centromere? *Exp Cell Res.* 389:111895.

- 344 22. Tanaka, Y., O. Nureki, H. Kurumizaka, S. Fukai, S. Kawaguchi, M. Ikuta, J. Iwahara, T. Okazaki, and S.
345 Yokoyama. 2001. Crystal structure of the CENP-B protein-DNA complex: the DNA-binding domains of CENP-B
346 induce kinks in the CENP-B box DNA. *EMBO J.* 20:6612-6618.
- 347 23. Verdaasdonk, J.S., and K. Bloom. 2011. Centromeres: unique chromatin structures that drive chromosome
348 segregation. *Nat Rev Mol Cell Biol.* 12:320-332.
- 349 24. Zhabinskaya, D., S. Madden, and C.J. Benham. 2015. SIST: stress-induced structural transitions in
350 superhelical DNA. *Bioinformatics.* 31:421-422.

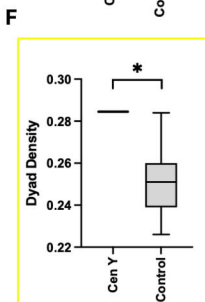
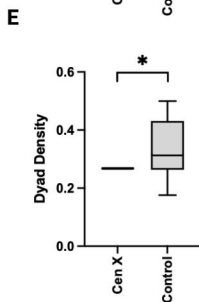
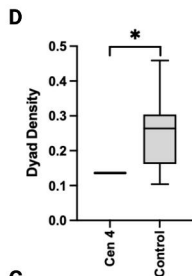
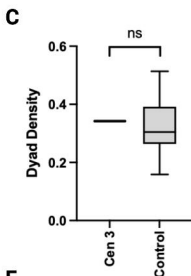
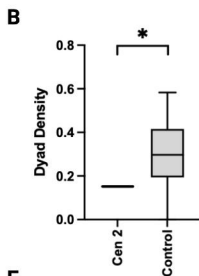
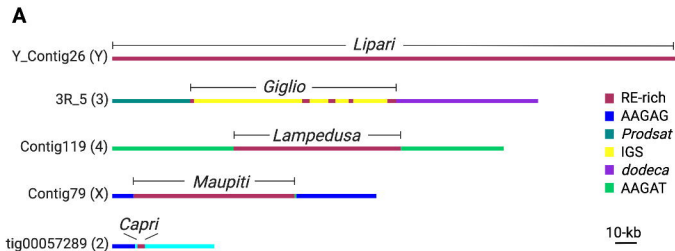


Figure 1

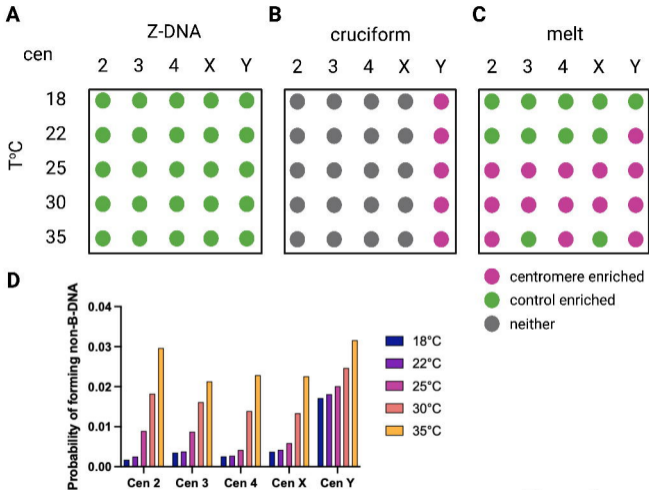


Figure 2

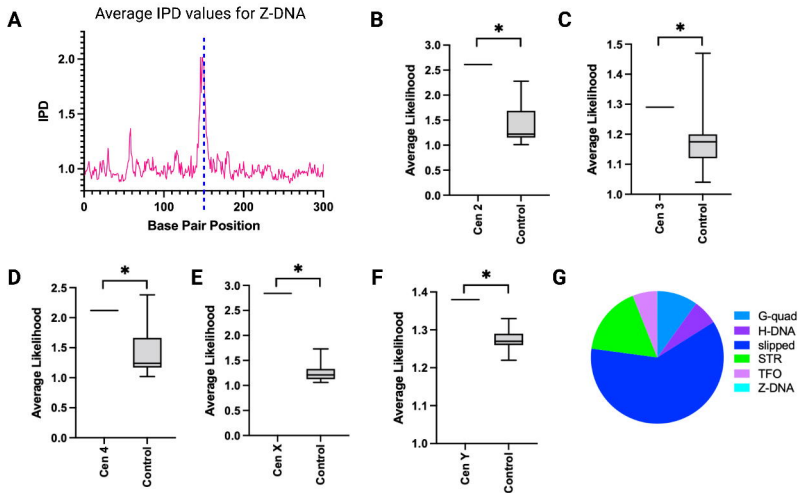


Figure 3

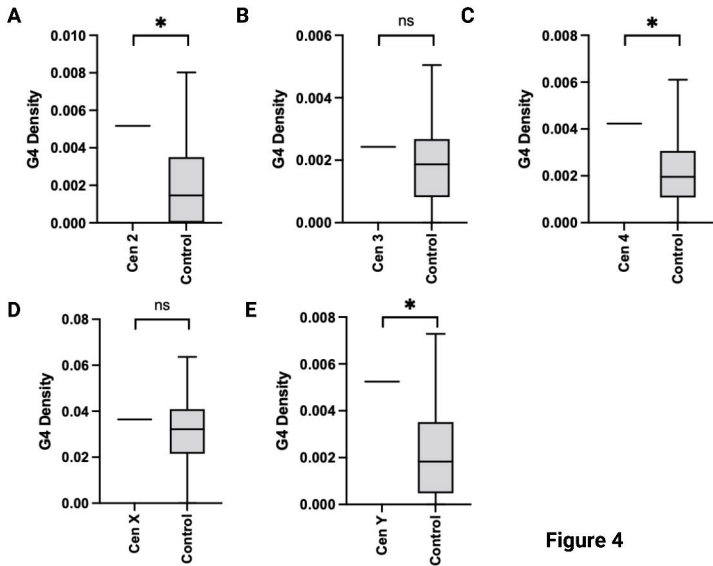


Figure 4