

1 **1. Extended Data**

Figure #	Figure title One sentence only	Filename This should be the name the file is saved as when it is uploaded to our system. Please include the file extension. i.e.: <i>Smith_ED_Fig1.jpg</i>	Figure Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
----------	-----------------------------------	---	---

2

3 **2. Supplementary Information:**

4

5 **A. Flat Files**

Item	Present?	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	Supplementary_Information.docx	Supplementary Notes 1 - 21
Reporting Summary	Yes	ReportingSummary_Basu	
Peer Review Information	Yes	<i>TPRFile_Basu</i>	

6

7 **B. Additional Supplementary Files**

Type	Number If there are multiple files of the same type this should be the numerical indicator. i.e. "1" for Video 1, "2" for Video 2, etc.	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_Supplementary_Video_1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Table	1	Supplementary_Table.xlsx	Supplementary Tables 1 - 21

8

9

3. Source Data

Parent Figure or Table	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_SourceData_Fig1.xls</i> , or <i>Smith_Unmodified_Gels_Fig1.pdf</i>	Data description i.e.: Unprocessed Western Blots and/or gels, Statistical Source Data, etc.
Source Data Fig. 1	Fig1_source_data.xlsx	Source data for Fig 1
Source Data Fig. 2	Fig2_source_data.xlsx	Source data for Fig 2
Source Data Fig. 3	Fig3_source_data.xlsx	Source data for Fig 3
Source Data Fig. 4	Fig4_source_data.xlsx	Source data for Fig 4

10

11

Deciphering the mechanical code of the genome and epigenome

12 Aakash Basu^{1,2,a}, Dmitriy G. Bobrovnikov², Basilio Cieza³, Juan Pablo Arcon⁴, Zan Qureshi³, Modesto
13 Orozco^{4,5} and Taekjip Ha^{2,3,6,7,a}.

14

15 ¹Department of Biosciences, Durham University, Durham, DH7 3LE, UK.

16 ²Department of Biophysics and Biophysical Chemistry, Johns Hopkins University School of Medicine,
17 Baltimore, MD 21205, USA.

18 ³Department of Biophysics, Johns Hopkins University, Baltimore, MD 21218, USA.

19 ⁴Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science and
20 Technology, 08028, Barcelona, Spain.

21 ⁵Department of Biochemistry and Biomedicine. Universitat de Barcelona. 08028, Barcelona, Spain.

22 ⁶Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21205.

23 ⁷Howard Hughes Medical Institute, Baltimore, MD 21205, USA.

24

25 ^aTo whom correspondence should be addressed (e-mail: aakash.basu@durham.ac.uk, tjha@jhu.edu)

26

27 **Abstract**

28 Diverse DNA-deforming processes are impacted by the local mechanical and structural properties of DNA,
29 which in turn depend on local sequence and epigenetic modifications. Deciphering this ‘mechanical code’,
30 i.e., this dependence, has been challenging due to the lack of high-throughput experimental methods. Here
31 we present a comprehensive characterization of the mechanical code. Utilizing high-throughput
32 measurements of DNA bendability via loop-seq, we quantitatively established how the occurrence and
33 spatial distribution of dinucleotide, tetranucleotides, and methylated CpG, impact DNA bendability. We
34 use our measurements to develop a physical model for the sequence and methylation dependence of DNA
35 bendability. We validate the model by performing loop-seq on mouse genomic sequences around
36 transcription start sites and CTCF binding sites. We apply our model to test the predictions of all atom
37 molecular dynamics simulations, and to demonstrate that sequence and epigenetic modifications can
38 mechanically encode regulatory information in diverse contexts.

39

40 **Introduction**

41 The local structural, shape, and mechanical properties of DNA can regulate diverse phenomena
42 that require DNA deformation. Various techniques have revealed that these properties depend on local
43 sequence and epigenetic modifications. For example, structural data⁶ have been used to predict various
44 shape parameters (twist/ roll/ tilt) associated with all dinucleotides. These models predict plectoneme
45 density along genomic DNA¹, and to some extent, nucleosome organization from DNA mechanics,
46 although they do not capture key features of nucleosome occupancy data². All atom Monte Carlo (MC)
47 simulations have predicted the sequence-dependence of local DNA shape^{3,4}, while molecular dynamics
48 (MD) simulations have observed surprising polymorphic behaviors associated with the shape parameters
49 describing B-DNA⁵⁻⁸. However, experimental validation of such predictions is extremely limited⁶.

50 Direct experimental methods to measure how DNA sequence impacts its mechanical properties
51 have mainly included low-throughput techniques such as measuring cyclization rates of short DNA

52 molecules⁹, DNase I accessibility of DNA minicircles^{10,11} and DNA stretching¹² via AFM or hydrodynamic
53 drag. Low-throughput has restricted their use to primarily studying the mechanics of isolated examples of
54 promoter sequences^{13,14}, nucleosomal sequences¹⁵⁻¹⁷, or sequences involved in transcription activation¹⁸.
55 High-throughput SELEX-based methods have identified sequence features that are more enriched or de-
56 enriched in highly loopable molecules by performing several rounds of selection from an initial random
57 pool of sequences^{19,20}. However, quantitative bendability values have not been assigned to any sequence,
58 and as a result, SELEX data have not led to quantitative measures of DNA mechanics, or quantitative
59 models.

60 Cytosine methylation in the CpG contexts is a form of epigenetic modification in mammalian
61 genomes, and is capable of altering gene expression during development and in disease through preventing
62 transcription factor binding and recruiting specialized proteins^{21,22}. CpG methylation has also been
63 suggested to impact transcription by altering the mechanical and structural properties of DNA and
64 chromatin: CpG methylation increases the rupture force necessary for DNA strand separation²³, reduces
65 DNA bendability¹⁷, alters the stability of nucleosomes^{17,24,25}, and enhances polyamine-mediated DNA-DNA
66 association²⁶. However, methylation-induced changes in DNA physics have not been characterized in
67 diverse sequence contexts.

68 We recently reported a high-throughput experimental method called loop-seq to determine the
69 ‘intrinsic cyclizability’ - a measure of cyclization propensity or bendability - of as many as ~90,000
70 different 50 bp DNA sequences at a time²⁷. Here we apply loop-seq to a large library of randomly specified
71 sequences, correlate intrinsic cyclizability with sequence features and CpG methylation state, and develop
72 a comprehensive characterization of the mechanical code. Application of our predictive model in various
73 contexts suggests the functional relevance of DNA mechanics. Loop-seq provides a high-throughput
74 experimental platform to test the predictions of MD simulations regarding distortion energy penalty and
75 structural polymorphism.

76

77 **Results**

78 We previously reported²⁷ the intrinsic cyclizabilities of 12,472 50 bp DNA fragments (this set of
79 randomly chosen sequences was termed the 'random library' (supplementary note 1)). Here, we measured
80 the intrinsic cyclizability of all sequences in the random library after we enzymatically methylated all CpG
81 dinucleotides for a side-by-side comparison. Overall, DNA becomes more rigid upon CpG methylation,
82 thereby suggesting that in addition to sequence, DNA mechanics depends on epigenetic modifications (Fig.
83 1a, Supplementary Note 2). We therefore use these two datasets as a starting point for understanding how
84 DNA sequence and epigenetic modifications influence DNA mechanics.

85

86 **Sequence features that influence DNA cyclizability**

87 The overall G/C content is uncorrelated with intrinsic cyclizability except for its increase at very high G/C
88 content above 65% (Fig. 1b). To test if certain dinucleotides influence intrinsic cyclizability, we sorted the
89 random library members according to intrinsic cyclizability and divided them into 12 bins. The number of
90 times each dinucleotide occurs in individual bins after appropriate normalization (Supplementary Note 3)
91 (Fig. 1c) varied monotonically with bin number (there are some exceptions we will discuss later). We used
92 a linear fit and defined its slope as the bendability quotient of each dinucleotide (Fig. 1d, Supplementary
93 Note 3). For example, the TA dinucleotide (TpA) has the most positive bendability quotient, meaning that
94 having more TpAs makes DNA more bendable on average. Conversely, the CG dinucleotide (CpG) has the
95 most negative bendability quotient. Previous studies have suggested that TpA is particularly flexible^{9,28-31}
96 and is associated with highly bent DNA in nucleosomes^{32,33}. Less is known about the contribution of CpG
97 to DNA mechanics. A SELEX-based study found that CpGs are depleted among the highly loopable
98 sequences, but they did not find specific enrichment of TA dinucleotides over other T/A dinucleotides¹⁹.

99 TT and AA dinucleotides are represented highly in both the highest and the lowest intrinsic
100 cyclizability bins (Fig. 1c), suggesting that their positional distribution could play a greater role in
101 determining intrinsic cyclizability than their overall content. Short stretches of dA:dTs bend DNA towards
102 the minor groove. These bend directions add in phase or cancel out if repeated at the helical or half-helical

103 period, resulting in globally curved or straight molecules that would have high or low intrinsic cyclizability,
104 respectively³⁴⁻³⁶.

105 The availability of high-throughput data now permits us to quantitatively evaluate how phased
106 repeats of any sequence motifs (not just short dA:dT stretches) influence DNA mechanics. There are 136
107 NN-NN dinucleotide pairs (such as AA-TT, AC-GG, etc.). For a given pair in a sequence, we define the
108 pairwise distance distribution function, $\rho(i)$, as the number of times that the two dinucleotides appear
109 separated by a distance i , appropriately normalized (Supplementary Note 4). For example, the AA-TT
110 dinucleotide pair shows striking oscillations in $\rho(i)$ when averaged over thousands of sequences: these
111 dinucleotides tend to more likely be separated by 10, 20 or 30 bp (i.e. multiples of the helical repeat) and
112 less likely be separated by 5, 15 or 25 bp (i.e. odd multiples of half the helical repeat) among the 1,000 most
113 cyclizable sequences (Fig. 1e, red curve) and vice versa among the 1,000 least cyclizable sequences (blue
114 curve). These observations are consistent with the known fact that short dA:dT stretches repeated every 10
115 bp result in curved molecules^{37,38}.

116 We discovered that other dinucleotide pairs can also influence intrinsic cyclizability depending on
117 their mutual distance (Fig. 1e). The CG – GC pair shows a similar behavior to AA-TT pair. The TT – GC
118 pair shows the reverse effects: sequences with high cyclizability have more instances of TT – GC separation
119 at 5, 15, or 25 bp, and vice versa. Several other pairs, such as CA – CA, show no oscillatory patterns. To
120 quantify these effects, we define a metric called the helical repetition index (HRI) of a given dinucleotide
121 pair in a given sequence as follows:

$$122 \quad HRI = \sum_{i=10,20,30} \max(\rho(i-1), \rho(i), \rho(i+1)) - \sum_{i=5,15,25} \max(\rho(i-1), \rho(i), \rho(i+1))$$

123 A high HRI of a dinucleotide in a sequence implies that, in that sequence, the dinucleotide pair in
124 question is more often separated by multiples of the helical repeat than by odd multiples of half the helical
125 repeat. HRI of each dinucleotide pair showed a linear dependence on intrinsic cyclizability, and we defined
126 the slope as a measure of the pair's contribution towards intrinsic cyclizability (Fig. 1f). Dinucleotide pairs

127 with all As or Ts make a positive contribution to intrinsic cyclizability (Fig. 1f). A similar effect was
128 observed for some of the dinucleotide pairs made of Gs or Cs only (Fig. 1f) and this is not a result of
129 exclusion of A/Ts (supplementary note 7), but is in fact consistent with an earlier suggestion that G/C
130 dinucleotides tend to bend DNA towards the major groove⁹. These curvatures might similarly add in phase
131 or cancel out in sequences where G/C rich dinucleotides occur every helical or half helical repeat,
132 respectively. In the case of hybrid pairs, where one dinucleotide is purely A/T containing and the other
133 purely G/C containing, DNA sequence with high HRI is less cyclizable. Consecutive bend directions
134 towards the minor and major groove every half-helical repeat or helical repeat would add in phase and
135 cancel out respectively, presumably contributing to this effect. These findings are consistent with SELEX
136 experiments which noted that A/T rich stretches separated by 5 bp from G/C rich stretches, on average, are
137 more represented in highly loopable molecules¹⁹. In addition, there are several hybrid pairs, such as AG-
138 GG, or GT-AG, which make significant positive or negative contributions to intrinsic cyclizability, which
139 could not have been predicted from prior literature.

140 Both dynamic flexibility as well as static bends may contribute to intrinsic cyclizability as noted
141 previously²⁷. Speculatively, the helical repetition index of dinucleotides determines intrinsic cyclizability
142 mainly by altering the static bent structure of DNA and determining whether alternate static bend directions
143 add in phase or cancel out. On the other hand, bendability quotients, by virtue of being independent of
144 nucleotide distribution, may reflect how dinucleotides contribute to dynamic flexibility.

145 We next asked whether the sequence features linked to high or low intrinsic cyclizability are found
146 in genomic landmarks. Genes are characterized by a well-ordered array of downstream nucleosomes
147 (labeled +1, +2, etc, depending on their distance from the TSS) and an upstream Nucleosome Depleted
148 Region (NDR)^{39,40}. Via loop-seq on *S. cerevisiae*, we discovered a narrow well-defined region of unusually
149 low intrinsic cyclizability co-centric with the NDR and regions of high local intrinsic cyclizabilities at the
150 locations of the downstream nucleosomes²⁷.

151 Averaged over 4,912 genes in *S. cerevisiae* that have previously been annotated as ORFs and both
152 ends mapped with high confidence⁴¹, G/C content shows no special features around the TSSs, other than a

153 general dip near the NDR (Fig. 1g), likely because NDRs are A/T rich. However, CpG content peaks at the
154 NDR and dips at nucleosomal dyads (Fig. 1g, supplementary note 8), contrary to the overall G/C content
155 distribution, but consistent with CpG having the most negative bendability quotient (Fig. 1d). Likewise,
156 TpA content peaks at +/-1 nucleosomal dyads and dips near the NDR (Fig. 1g, supplementary note 8),
157 consistent with TpA having the highest bendability quotient (Fig. 1d). The sum of the helical repetition
158 index of the 10 dinucleotide pairs that make the most negative contribution to intrinsic cyclizability peaks
159 near the NDRs and dips near the dyads of the +1 nucleosomes (Fig. 1g). In other words, such pairs tend to
160 more often be separated at the helical repeat (or less often at half the helical repeat) near the NDRs of genes,
161 and vice versa near the dyads of +1 nucleosomes. Likewise, the sum of the helical repetition index of the
162 10 dinucleotide pairs that make the most positive contribution to intrinsic cyclizability (such as AA-TT,
163 GC-GC, etc, Fig. 1f) averaged over all genes, peaks at the dyads of +1 nucleosomes and dips near the NDRs
164 (Fig. 1g). Therefore, the sequence features for high or low intrinsic cyclizability, identified from the random
165 library, are indeed found in genomic regions where we have independently measured²⁷ intrinsic cyclizability
166 to be high or low, respectively. The low intrinsic cyclizability at the NDRs, and other shape properties of
167 DNA in this region, may contribute to nucleosome depletion by impacting the nucleosome sliding activities
168 of chromatin remodelers, as suggested earlier^{27,42,43}. The high rigidity of DNA may also reduce nucleosome
169 occupancy by energetically disfavoring extensive DNA bending^{32,44-48}, although its effect may be
170 overshadowed by the action of chromatin remodelers.

171

172 **Intrinsic cyclizability is related to DNA shape**

173 DNA shape plays a fundamental role in determining the extent of protein-DNA interactions⁴⁹ and
174 the binding specificities of transcription factors⁵⁰ and nucleosome formation⁵¹. It is described by parameters
175 such as twist, roll, propeller twist, etc⁵². Propeller twist is the angle between the planes of the two bases in
176 a base pair (Supplementary Note 9). High propeller twist has been suggested to be an attribute of rigid
177 DNA⁵³. Indeed, we found that sequences with low measured intrinsic cyclizability have high predicted
178 propeller twist values (Fig. 2a).

179 Twist is the angle by which one base pair rotates relative to the adjacent base pair around the DNA
180 duplex axis. The CpG step transitioned between high and low twist conformations in MD simulation with
181 their relative populations depending on its tetranucleotides context⁵ and even within a defined
182 tetranucleotide context their relative weights depended on the identities of the flanking nucleotides^{5,6}. To
183 provide an experimental test of these predictions⁵⁴ we first estimated the bendability quotients of all 256
184 NNNN tetranucleotides (Fig. 2b, supplementary note 10). NCGN and NTAN tetranucleotides tend to have
185 low and high bendability quotients, respectively, consistent with CpG and TpA having the lowest and
186 highest bendability quotients, respectively (Fig. 2b). NCGNs with stronger weight toward the low twist
187 conformation in MD simulations also showed low measured bendability quotients and vice versa (Fig. 2c),
188 lending experimental support to the polymorphic twist behavior of CpGs⁶ and suggesting that low
189 bendability of CpGs arises from their low twist conformation.

190 Roll measures the rotation of a base pair plane around its long axis and may influence DNA
191 bending⁵⁵. All-atom MD simulations have reported the equilibrium roll angles associated with all NNNN
192 tetranucleotides^{7,8}. However, when a tetranucleotide occurs as part of a cyclized molecule, its roll angle has
193 to deviate from the unstressed value (supplementary note 11). We tested if the large variation in bendability
194 even for those sharing the same central dinucleotides is caused by variations in roll energy penalty incurred
195 upon cyclization. Indeed, for 16 different NCGN tetranucleotides, the energy penalty calculated from MD
196 simulations was higher if the experimentally determined bendability quotient is lower ($r_{Pearson} = -0.83$)
197 (Supplementary Note 11) (Fig. 2d). Expanding the analysis to all 16 dinucleotides (NNs) showed several
198 other NNs where the 16 bendability quotients of tetranucleotides that contain the NN in the middle are
199 significantly anti-correlated with the corresponding roll energy values (Supplementary Note 12). Therefore,
200 roll energy penalty is a major determinant of DNA bendability in a broad sequence context.

201 Overall, MD simulations provided atomic level insights into the sequence-dependence of DNA
202 bendability derived from loop-seq analysis.

203

204 **CpG methylations stiffens DNA in all sequence contexts**

205

206 CpG methylation has been proposed to affect transcription by altering the relevant physical
207 properties of chromatin^{17,23-25}. However, there has been no high-throughput characterization of how DNA
208 mechanics is impacted by CpG methylation in a broad sequence context. Recently, DNA cyclization
209 measurements and MD simulations suggested that CpG methylation can make DNA more rigid^{17,56}. These
210 experiments, however, were performed on a single DNA sequence, and require reevaluation because they
211 did not account for the surface tethering position that sets the phase of preferred bending²⁷. Intrinsic
212 cyclizability, which is determined from loop-seq data with three different tethering positions, showed that
213 CpG methylation, averaged over all sequence contexts, makes DNA more rigid (Fig. 1a). CpG methylation
214 decreases the bendability quotients of all tetranucleotides that contain at least one CG (Fig. 2e), indicating
215 that CpG methylation stiffens DNA in all sequence contexts.

216 We next examine the effect of CpG methylation on how the spatial distribution of dinucleotides
217 influences intrinsic cyclizability. Higher helical repetition index of unmethylated CG-NN pairs make DNA
218 more cyclizable if N = G or C, and more rigid if N = A or T (Fig. 1f). CpG methylation de-emphasize the
219 importance of phased dinucleotide pairs: CG-NN pairs with G/C rich and A/T rich NNs decrease and
220 increase, respectively, their contribution to intrinsic cyclizability as a result of CpG methylation (Fig. 2f).
221 Because contributions of dinucleotide pairs towards intrinsic cyclizability likely arise from the directions
222 of consecutive static bending every helical repeat either adding in phase or canceling out, we propose that
223 CpG methylation suppresses DNA bending towards the major groove at CpG steps. It has indeed been
224 suggested, based on crystal structures, that methylation increases steric hindrance at the major groove,
225 thereby widening the major groove and suppressing bending⁵⁷. Overall, loop-seq provides evidence that
226 CpG methylation decreases dynamic flexibility and buffers sequence-dependent static curving of DNA.

227

228 **Models to predict intrinsic cyclizability**

229

230 After identifying sequence features that impact intrinsic cyclizability, we aimed to build a unified
231 physical model that predicts DNA mechanics from sequence by incorporating the effects of NN content
232 and NN-NN helical repetition indices (HRIs) (Fig. 1). A 50 bp DNA sequence can be described by a set of
233 16 NN contents (the number of times each of the 16 dinucleotides occurs in the sequence) and 136 NN-NN
234 HRIs. We built linear models that treat intrinsic cyclizability as a linear combination of firstly just the NN
235 contents, then just the NN-NN HRIs, and finally all the 152 parameters (16 NN contents and 136 NN-NN
236 HRIs (supplementary note 13, Figs. 3a-b). The best fit coefficients are obtained through multivariate linear
237 regression, using the Tiling Library as the training dataset. We tested the models by predicting the intrinsic
238 cyclizabilities of sequences in the Random, ChrV and *Cerevisiae* Nucleosomal Libraries (Fig. 3c).
239 Correlation coefficients between measured and predicted values obtained by the combined linear model are
240 in the range 0.52 – 0.54 (Fig. 3c), which are significantly higher than those obtained using NN contents
241 alone or NN-NN HRIs alone and are close to what we obtained using neural nets (Supplementary Note 14).

242 Predicted intrinsic cyclizability around 576 genes in *S. cerevisiae* captures all essential features that
243 were identified using the measured intrinsic cyclizability values – a well-defined region of unusually low
244 intrinsic cyclizability co-centric with the NDR, and local peaks in intrinsic cyclizability around the dyads
245 of downstream nucleosomes (Fig. 3d). The absolute values do not match between prediction and
246 measurement because the measured intrinsic cyclizability is defined only up to an additive constant which
247 depends on other sequences present in the library²⁷.

248

249 **The impact of DNA mechanics in promoters and CTCF sites**

250 To explore relations between DNA mechanics and nucleosome organization, we compared
251 predicted intrinsic cyclizability with measured nucleosome occupancy around the TSSs in four different
252 organisms (Fig. 4a). For *S. cerevisiae*, predicted intrinsic cyclizability patterns broadly agree with what was
253 measured from a subset of genes (576 genes)²⁷ (Fig. 3d) – there is a well-defined region of rigid DNA
254 upstream of the +1 dyad, which coincides with the NDR (Fig. 4a). Additionally, flexibility peaks are visible
255 at the location of the +1 and -1 nucleosomes, and to a lesser extent, the +2 nucleosome. For *S. pombe*, we

256 observed two marked differences⁵⁸ (Fig. 4a). First, there is no special region of rigid DNA upstream of the
257 +1 nucleosome, and second, intrinsic cyclizability shows much stronger oscillations in phase with the
258 nucleosome occupancy oscillation, at least up to +7 nucleosome. It is possible that in the absence of a well-
259 defined region of rigid DNA at the NDR, paused polymerases⁵⁹, transcription factors⁶⁰ or other bound
260 factors may serve as a similar barrier to nucleosome sliding in *S. pombe*. The more prominent oscillations
261 in intrinsic cyclizability are indicative of a greater importance of DNA mechanics in evenly spacing gene
262 body nucleosomes in *S. pombe* than in *S. cerevisiae*, although the overall contribution of such an effect in
263 either species could be dwarfed by the action of chromatin remodelers. Future experiments involving
264 altering DNA flexibility via genome editing would be necessary to understand differences in how either
265 species translates genomic information into nucleosome positioning. In drosophila and mouse, our
266 predictive models suggest that, as in *S. cerevisiae*, there exists a sharply-defined region of rigid DNA near
267 TSSs. Further, in drosophila, small peaks in intrinsic cyclizability are also visible at the locations of the +1
268 and +2 nucleosomes.

269 We next predicted how CpG methylation would change the mechanical landscape around TSSs. In
270 *S. cerevisiae*, *S. pombe*, and *D. melanogaster*, CpG methylation does not cause any significant change in
271 the overall predicted landscape of intrinsic cyclizability around TSSs (Fig. 4a). CpG methylation also does
272 not occur in these organisms. In contrast, our model predicts that methylation significantly alters the
273 landscape around TSSs in mouse (Fig. 4a): the characteristic sharply-defined region of rigid DNA near
274 TSSs becomes broadened and less well-defined. To test the prediction, we measured the intrinsic
275 cyclizabilities near the TSS of 649 randomly selected genes in mouse with or without CpG methylation.
276 The measured mechanical landscape indeed confirmed that CpG methylation results in the broadening and
277 dilution of the rigid DNA region (Fig. 4b, Supplementary Note 17). It is therefore possible that methylation-
278 induced changes in DNA mechanics near mouse promoters contributes to transcriptional silencing by
279 altering promoter-proximal nucleosome positioning and the action of chromatin remodelers which we have
280 shown to depend on DNA mechanics²⁷.

281 Nucleosomes are also well organized around other genomic landmarks, one notable example being
282 the binding sites of the transcription factor CTCF⁶¹. CTCF sites are surrounded by well-ordered
283 nucleosomes and the sites themselves are either free of nucleosomes or bound by fragile nucleosomes^{62,63}.
284 Our model predicts that CTCF binding sites in mouse embryonic stem cells^{62,64} are characterized by a
285 prominent peak in intrinsic cyclizability (Fig. 4c). The width of the peak extends significantly beyond the
286 20 bp CTCF consensus motif, and this is not caused by averaging (supplementary note 18), consistent with
287 the report that a structured DNA region exceeds the CTCF consensus motif⁶⁵. The sharp peak in intrinsic
288 cyclizable may aid in the known positioning of nucleosomes at CTCF sites. Indeed, our model predicts that
289 this peak cyclizability value, as well as the contrast in cyclizability between the peak and the surrounding
290 DNA is higher for sites with higher nucleosome occupancy (Fig. 4d). G/C content alone cannot explain
291 these observations (Supplementary Note 20).

292 We tested these model predictions on CTCT sites using loop-seq. The measured mechanical
293 landscape shows a prominent peak in cyclizability at CTCF sites which extends beyond the 20 bp consensus
294 motif (Fig. 4e). Further, we observed two flanking peaks in cyclizability which were not captured by the
295 predictive model and align with nucleosome occupancy⁶⁵, suggesting that they might aid in the positioning
296 of the flanking nucleosomes (Fig. 4e). The 5' upstream flanking peak in cyclizability is somewhat further
297 from the CTCF motif than the corresponding peak in nucleosome occupancy, consistent with the
298 requirement of the chromatin remodeler SNF2H for reduced nucleosome repeat length⁶⁵. We further
299 experimentally verified the prediction that CTCF sites with higher nucleosome occupancy have a higher
300 contrast in intrinsic cyclizability (Fig. 4f).

301

302 **Impact of DNA mechanics in non-nucleosomal contexts**

303 Next, we examined DNA mechanics in several non-nucleosomal contexts.

304 Groups of A_n/T_n nucleotides present in a periodic fashion, extending for hundreds of nucleotides,
305 has been reported in eukaryotes genomes^{66,67}. Such hyperperiodic DNA is particularly pronounced in the

306 *C. elegans* genome at ~10 bp periodicity, and recent structural studies have revealed that it is curved⁶⁸. The
307 intrinsic cyclability predicted across the hyperperiodic segment using the NN-NN phased repeat model
308 indeed shows much higher values than when the nucleotide order in the ~900 bp region is scrambled,
309 suggesting that hyperperiodic phase repeats are responsible for the curved DNA observed (Fig. 4g). The
310 functional consequences of hyperperiodic DNA extending for hundreds of nucleotides are unclear⁶⁸ and
311 require future experiments.

312 TATA binding protein (TBP) binds to promoter DNA, about 30 bp upstream of the TSS in *S.*
313 *cerevisiae*⁶⁹, during transcription initiation⁷⁰. The TATA-box – a consensus sequence to which TBP binds,
314 is present in only 15% of promoters, although TBP binding is universally required for transcription⁷¹. TBP
315 binding to the TATA-box induces DNA bending by ~90°⁷². Consistent with this observation, sequence-
316 dependent curvature of DNA is linked to transcription efficiency⁷³, and TBP binds less efficiently to poly-
317 A tracts. Our predictive model, averaging across all categories of TSSs⁷⁴, showed a local peak about 25 bp
318 upstream of the TSSs (Fig. 4h). This coincides with the location of the TBP⁶⁹ and of the TATA box in
319 TATA-containing promoters. Our findings suggest that genome-wide, DNA mechanics might influence
320 DNA bending by TBP regardless of the presence of the TATA-box.

321 DNA bending is also ubiquitous in bacteria. The *E. coli* RNA polymerase $\sigma 70$ (RNAP) open
322 promoter complex has 300° wrapping of about ~90 bp of DNA around the RNAP⁷⁵. Our predicted intrinsic
323 cyclizability averaged around *E. coli* TSSs⁷⁶ show a sharp peak from -35 to 0 position relative to the TSS,
324 and an overall region of elevated intrinsic cyclizability extending from about position -100 to about +150
325 (Fig. 4i). It is possible that such an extended region of flexible DNA might aid in extensive DNA bending
326 required to form the open promoter complex and for DNA wrapping around RNAP although verification
327 requires future experiments.

328 Negative supercoiling of DNA by DNA gyrase is critical for bacterial life. Gyrase bends DNA at
329 its cleavage site and extensively wraps DNA up to 80 bp away⁷⁷. Strong Gyrase Sites (SGSs) in some
330 genomes have been proposed to allow extensive DNA bending⁷⁸. FluMu is a prophage whose 35kb genome
331 is found inserted within the *Haemophilus influenzae* genome^{79,80}, and contains one SGS⁸¹. Predicted intrinsic

332 cyclizability along the entire genome has three most prominent peaks all within 200 bp of the cleavage site
333 of the SGS (Fig. 3j). One of the peaks coincides with the cleavage site, where gyrase must bend DNA. The
334 second peak lies within 80 bp of the cleavage site – this region is where the wrapping of DNA occurs, in
335 turn requiring extensive bending. Thus the mechanical properties of DNA around SGSs may have evolved
336 to facilitate gyrase activity.

337 Looping of 100 – 140 bp of DNA is required for efficient communication between enhancers and
338 σ^{54} -dependent promoters in bacteria⁸². The bacterial protein IHF extensively bends DNA⁸³. Low-
339 throughput DNA cyclization experiments have revealed that in IHF-dependent σ^{54} -promoters, DNA cannot
340 intrinsically bend to support this loop formation, although it can in IHF-independent promoters. Applying
341 our predictive model, we found that indeed among the eight σ^{54} -dependent promoters on which cyclization
342 experiments were performed earlier⁸², the two that did not contain an IHF-binding site were more cyclizable
343 than the rest in the region between the enhancer and the promoter (Fig. 4k). This is consistent with the
344 proposal⁸² that a function of IHF might be to overcome a communication block set up by the rigid
345 mechanical properties of the IHF binding site. In certain promoters, intrinsically flexible DNA may take
346 over the role of IHF in facilitating DNA bending.

347

348 **Discussion**

349 We obtained a comprehensive characterization of the mechanical code of the genome and
350 epigenome using high-throughput experimental measurements. Using high-throughput loop-seq
351 measurements, we quantified the contributions of short sequence features (such as all dinucleotide contents
352 and HRIs) to intrinsic cyclizability and measured, in a broad sequence context, how these contributions are
353 altered by CpG methylation. We demonstrate loop-seq to be a powerful experimental platform for testing
354 the predictions of structural studies and MD simulations in a broad sequence context, and making new
355 testable predictions about the structural dynamics of nucleic acids. We integrated our data into regression
356 and machine learning based models for predicting how local sequence and epigenetic modifications impact

357 DNA mechanics. Applications of our predictive model, and its confirmation via new loop-seq
358 measurements, broadly suggests that DNA mechanics may have evolved to facilitate processes that require
359 DNA bending in both eukaryotes and bacteria.

360 Physical deformations of DNA are not restricted to bending, but include twisting, supercoiling, or
361 unzipping, during, for example, transcription initiation⁸⁴ and elongation⁸⁵, the binding of various
362 transcription factors⁸⁶, or 3D folding of the genome at local⁸⁷ and global scales⁸⁸. Future high-throughput
363 experimental and computational methods to understand how DNA sequence impacts such diverse modes
364 of DNA deformability will expand our understanding of how regulatory information is encoded via the
365 mechanical code, and how this information can be evolved, adapted, or controlled by processes that
366 chemically alter DNA via sequence mutations, epigenetic modifications, or DNA damage.

367

368 **Acknowledgements:**

369 A.B. and T.H. would like to thank Jun S. Song for suggestions and insights pertaining to developing the
370 linear predictive models. This work was supported by the Royal Society URF\R21\211659 (to A.B.), the
371 Royal Society RF\ERE\210288 (to A.B.), funding from Durham University (to A.B.), the National
372 Science Foundation Grants PHY-1430124, EFMA 1933303 (to T.H.), the National Institutes of Health
373 Grant GM122569 (to T.H.), the European Union's Horizon 2020 research and innovation programme
374 under the Marie Skłodowska-Curie grant agreement No. 754510 (to J.P.A), the Spanish Ministry of
375 Science RTI2018-096704-B-100 and AGAUR, Generalitat de Catalunya, Grups de Reserca Consolidats
376 2017 SGR 1110 (to M.O.). T.H. is an Investigator with the Howard Hughes Medical Institute. A.B. is a
377 Royal Society University Research Fellow.

378 **Author contributions:**

379 AB and TH designed research. AB performed research, analyzed data, and built the predictive models. AB
380 and TH wrote the paper. DGB compiled nucleosome occupancy data from various organisms. BC

381 investigated pairwise correlation among NN-NN dinucleotide pairs in highly loopable and rigid sequences.
382 ZQ assisted with library preparation loop-seq experiments pertaining to the random library. JPA and MO
383 related cyclizabilities to DNA shape parameters.

384

385 **Competing Interests:**

386 The authors declare no competing interests.

387

388 **Figure Legends:**

389 **Figure 1:** (a) The ordinate depicts mean intrinsic cyclizabilities of sequences in the random library (blue)
390 and in the methylated random library (red) that have at least the number of CpGs as specified by the
391 abscissa. Intrinsic cyclizability values of the methylated random library have been adjusted to allow for
392 comparison with those of the random library (Supplementary Note 2). For each color, from left to right, N
393 = 12,472, 12,036, 10,451, 7,620, 4,516, 2,147, 864, 282. Error bars are s.e.m. (b) Mean intrinsic
394 cyclizability of sequences in the random library as a function of G/C content. Mean is calculated only over
395 those sequences in the random library whose G/C content is as specified by the abscissa. Error bars are
396 s.e.m. From left to right, N = 30, 77, 110, 214, 324, 500, 746, 984, 1239, 1353, 1403, 1316, 1210, 982, 749,
397 487, 307, 201, 120, 53, 27. (c) The 12,472 sequences in the random library were sorted according to
398 increasing intrinsic cyclizability and grouped into 12 bins with 1,039 sequences each. 4 remaining
399 sequences were ignored. Within each bin, the normalized number of times each of the 16 dinucleotides
400 occur is color coded and depicted (see supplementary note 3). (d) Bendability quotient for a dinucleotide is
401 defined as the slope of the linear fit to the plot of the normalized number of times it occurs among sequences
402 in each of the 12 bins in panel c, vs the mean intrinsic cyclizability of sequences in the 12 bins
403 (supplementary note 3) (e) Pairwise distance distribution functions vs separation distance, averaged over
404 the 1,000 sequences in the random library that have the most (red) or least (blue) values of intrinsic
405 cyclizability, for four different NN-NN pairs. See Supplementary Note 5 for plotting details. (f) For a given

406 NN-NN dinucleotide pair, the best fit linear relationship between its helical separation extent in a sequence
407 and the intrinsic cyclizability of that sequence, for all sequences in the random library, was obtained. The
408 heatmap here depicts the slopes of these linear relationships for all 136 NN-NN pairs. See Supplementary
409 Note 6 for details. (g) Nucleosome occupancy and various sequence parameters as functions of position
410 from the dyad of the +1 nucleosome, averaged over all identified 4,904 genes in *S. cerevisiae* (see
411 supplementary note 8 for plotting details). To calculate $\sum_{rigid} \sigma_{NN-NN}$, the ten NN-NN pairs that make
412 the most negative contribution to intrinsic cyclizability were identified (supplementary note 8). The sum of
413 the helical separation extents of these pairs over a 50 bp DNA fragment centered around the ordinate value
414 was calculated for each gene. The values were averaged to obtain the abscissa value. $\sum_{flexible} \sigma_{NN-NN}$ was
415 similarly calculated. See supplementary note 8 for heatmaps of TpA and CpG contents.

416 **Figure 2:** (a) Mean propeller twist as a function of position averaged over N = 1,000 50 bp DNA sequences
417 in the random library that had the most (red) and least (blue) values of intrinsic cyclizability. Sequences
418 were tiled as a series of pentamers and the associated propeller twist of the central base in the pentamer was
419 assigned on the basis of earlier reports⁴. Error bars are s.e.m. (b) The 12,472 sequences in the random library
420 were sorted according to increasing intrinsic cyclizability and grouped into 12 bins with 1,039 sequences
421 each. 4 remaining sequences were ignored. Within each bin, the normalized number of times each of the
422 256 tetranucleotides occur is color coded and depicted (Supplementary Note 10). Tetranucleotides
423 containing a CG in the middle (ie of the form NCGN) are indicated by a magenta circle, while those
424 containing a TA in the middle are indicated by a green circle. (c) Bendability quotients of all 16 NCGN
425 tetranucleotides as obtained from loop-seq of the random library (see Supplementary Note 10) vs the
426 weighted average of the high-twist and low-twist conformations that the central CG step samples, as
427 obtained from MD simulations⁶. Pearson's r value is shown. 95% confidence interval (CI) = 0.33, 0.89. P
428 value determined by t-test (two-sided). (d) Bendability quotients of all 16 NCGN tetranucleotides vs roll
429 energy. Roll energy is the energy penalty associated with the tetranucleotide having to deviate from its
430 equilibrium roll angle when adopting a constrained conformation after looping, i.e., the tetranucleotide

431 being part of a circle of circumference 110 bp (Supplementary Note 11). Pearson's r value is shown. 95%
432 confidence interval (CI) = -0.57, -0.94. P value determined by t-test (two-sided). (e) Bendability quotients
433 for all 256 NNNs obtained from the measured values of intrinsic cyclizability of the 12,472 sequences in
434 the random library vs those obtained from the measured values of intrinsic cyclizability of the sequences in
435 the methylated random library (which contain the identical set of 12,472 sequences, except all occurring
436 CpG are cytosine methylated). Dashed line represents $x=y$. NNNs are marked in red if at least one CG
437 occurs in it (such as ACGA, CGCG, CGAC, etc). Other NNNs (such as AAGC, GGGC, etc) are marked
438 in blue. (f) Heatmap representing the contributions of all NN-CG dinucleotide pairs towards intrinsic
439 cyclizability, obtained by considering the intrinsic cyclizability values of sequences in the random library
440 (first column, identical to the 15th column in figure 1f except for the color scale), and obtained from
441 measurements on the methylated random library (second column). Contribution towards intrinsic
442 cyclizability of a NN-CG pair is calculated as done in the case of figure 1f.

443 **Figure 3:** (a) 2-D histogram of the scatter plot between measured intrinsic cyclizabilities of sequences in
444 the random library and the associated predicted intrinsic cyclizability. Here, prediction was made via a
445 model where intrinsic cyclizability of a 50 bp sequence is a linear combination of the number of times each
446 of the 16 dinucleotides occur in the sequence and a constant term (supplementary note 13). Best fit
447 coefficients of the linear model were derived by training the model using the measured intrinsic cyclizability
448 values of sequences in the tiling library (Supplementary Note 1). Pearson's r value is shown. 95%
449 confidence interval (CI) = 0.34, 0.37. $P = 0.0000$ (determined by two-sided t-test). (b) Same as panel a,
450 except that prediction was made using a linear model where intrinsic cyclizability of a 50 bp sequence is a
451 linear combination of the 136 helical separation extent values in the sequence of the 136 NN-NN pairs, and
452 a constant term (supplementary note 13). Pearson's r value is shown. 95% confidence interval (CI) = 0.36,
453 0.39. $P = 0.0000$ (determined by two-sided t-test). (c) 2-D histogram of the scatter plots between measured
454 and predicted intrinsic cyclizability values of sequences in the random, chrV, and cerevisiae nucleosomal
455 libraries (Supplementary Note 1). Here, prediction was made using a model where intrinsic cyclizability of

456 a 50 bp sequence is a linear combination of a constant term and the 16 dinucleotide contents (subject to the
457 constraint that their sum = 49) and 136 helical separation extents that describe the sequence⁸⁹
458 (supplementary note 13). Coefficients were derived by training the model against the tiling library.
459 Pearson's r values are shown. 95% confidence intervals (CIs) are (0.51, 0.53), (0.52, 0.55), and (0.52, 0.55).
460 P = 0.0000 (determined by two-sided t-test) in all cases. (d) Measured intrinsic cyclizability, predicted
461 intrinsic cyclizability, and nucleosome occupancy as functions of position from the dyad of the +1
462 nucleosome, averaged over 576 genes in *S. cerevisiae*. Mean occupancy values⁹⁰ and positions³³ were as
463 reported earlier. Prediction was performed using the linear physical model trained using the measured
464 intrinsic cyclizability values of the random library. See supplementary note 15 for details.

465 **Figure 4:** (a) Nucleosome occupancy, predicted intrinsic cyclizability in absence of CpG methylation,
466 predicted intrinsic cyclizability in presence of CpG methylation, and G/C content, as a function of distance
467 from the dyad of the +1 nucleosome (in the case of *S. cerevisiae* and *S. pombe*) or from the TSS (in the case
468 of drosophila and mouse), averaged over a large number of genes in these four organisms. Nucleosome
469 occupancy metrics were obtained from previous publications. See supplementary note 16 for details. (b)
470 Measured intrinsic cyclizability as a function of position from the TSS, averaged over 629 mouse genes, in
471 absence and presence of CpG methylation. See supplementary note 17 for details. (c) Top panel: Predicted
472 intrinsic cyclizability (predicted by using the linear physical model trained against the random library) as a
473 function of position from the start of the 20 bp CTCF consensus motif, averaged over 19,900 reported
474 CTCF binding sites in mouse embryonic stem cells⁶⁴. Bottom panel: same as the top panel, except DNA
475 outside the 20 bp consensus sequence motif was chosen at random and not obtained from the mouse
476 genome. See Supplementary Note 18 for details. (d) Predicted intrinsic cyclizability as a function of position
477 (where 0 is the start of the CTCF consensus motif), averaged over the two groups of 1,000 sites that have
478 the highest (blue) and lowest (red) nucleosome Center Weighted Occupancy (CWO⁶²) at the CTCF motif.
479 See supplementary note 19 for details. (e) Measured intrinsic cyclizability and nucleosome occupancy vs
480 position from the edge of the CTCF motif, averaged over 433 randomly selected CTCF binding sites. See

481 supplementary note 21 for details. Nucleosome occupancy is the nucleosome CWO⁶². (f) Measured mean
482 intrinsic cyclizability vs position from the edge of the CTCF motif, averaged over two sets of 433 CTCF
483 binding sites that have the least and greatest mean nucleosome CWO⁶² at the CTCF motif. See
484 supplementary note 21 for details. (g) Predicted intrinsic cyclizability as a function of position along the
485 923 bp Ω 4 region of the *C. elegans* genome⁶⁸ (blue), and along three more 923 bp DNA sequences obtained
486 by randomizing the order of nucleotides that occur along the native Ω 4 sequence (black, green, red). (h)
487 Predicted intrinsic cyclizability as a function of distance from TSSs in *S. cerevisiae*. Data is averaged over
488 all 11,102 annotated TSSs⁷⁴ that were more than 500 bp away from chromosome edges. The red arrow
489 points to a local peak in intrinsic cyclizability coinciding with the known location of the TATA box. The
490 black arrow indicates the flexibility peak associated with the +1 nucleosome dyad, as reported earlier²⁷ (i)
491 Predicted intrinsic cyclizability as a function of distance from TSSs in *E. coli*. Data is averaged over 14,860
492 annotated TSSs⁷⁶. (j) Predicted intrinsic cyclizability along the 35 kb genome of the MU-like prophage
493 FluMu found inserted within the *H. influenza* genome. Arrow points to the region of the gyrase cleavage
494 site. Right panel: zoomed view of a 1 kb region around the strong gyrase binding site in the FluMu genome.
495 The red dashed line corresponds to the gyrase cleavage site. The dashed grey lines demarcate 80 bp from
496 the cleavage site. (k) Predicted intrinsic cyclizability as a function of position along the 220 bp DNA
497 fragment encompassing the enhancer and promoter of the eight σ 54-promoters on which earlier DNA
498 cyclization experiments have been reported. In bold are the two promoters (*K. pneumonia nifLA* promoter
499 and *E. coli glnAp2* promoter) among the eight which lack an IHF binding site. Plotted on the x-axis is the
500 position along the 220 bp fragment from the enhancer to the promoter.

501

502 **References**

503

504 1 Kim, S. H. *et al.* DNA sequence encodes the position of DNA supercoils. *Elife* **7**, e36557
505 (2018).

506 2 Morozov, A. V. *et al.* Using DNA mechanics to predict in vitro nucleosome positions and
507 formation energies. *Nucleic acids research* **37**, 4707-4722 (2009).

508 3 Rohs, R., Sklenar, H. & Shakked, Z. Structural and energetic origins of sequence-specific
509 DNA bending: Monte Carlo simulations of papillomavirus E2-DNA binding sites.
510 *Structure* **13**, 1499-1509 (2005).

511 4 Chiu, T.-P. *et al.* GBshape: a genome browser database for DNA shape annotations.
512 *Nucleic acids research* **43**, D103-D109 (2015).

513 5 Pasi, M. *et al.* μ ABC: a systematic microsecond molecular dynamics study of
514 tetranucleotide sequence effects in B-DNA. *Nucleic acids research* **42**, 12272-12283
515 (2014).

516 6 Dans, P. D. *et al.* Unraveling the sequence-dependent polymorphic behavior of d (CpG)
517 steps in B-DNA. *Nucleic acids research* **42**, 11304-11320 (2014).

518 7 Dans, P. D. *et al.* The static and dynamic structural heterogeneities of B-DNA: extending
519 Calladine–Dickerson rules. *Nucleic acids research* **47**, 11090-11102 (2019).

520 8 Walther, J. *et al.* A multi-modal coarse grained model of DNA flexibility mappable to the
521 atomistic level. *Nucleic acids research* **48**, e29-e29 (2020).

522 9 Geggier, S. & Vologodskii, A. Sequence dependence of DNA bending rigidity.
523 *Proceedings of the National Academy of Sciences* **107**, 15421-15426 (2010).

524 10 Brukner, I., Jurukovski, V. & Savic, A. Sequence-dependent structural variations of DNA
525 revealed by DNase I. *Nucleic acids research* **18**, 891-894 (1990).

526 11 Brukner, I., Sanchez, R., Suck, D. & Pongor, S. Sequence-dependent bending propensity
527 of DNA as revealed by DNase I: parameters for trinucleotides. *The EMBO journal* **14**,
528 1812-1818 (1995).

529 12 Rief, M., Clausen-Schaumann, H. & Gaub, H. E. Sequence-dependent mechanics of single
530 DNA molecules. *Nature structural biology* **6**, 346-349 (1999).

531 13 Davis, N. A., Majee, S. S. & Kahn, J. D. TATA box DNA deformation with and without
532 the TATA box-binding protein. *Journal of molecular biology* **291**, 249-265 (1999).

533 14 Parvin, J. D., McCormick, R. J., Sharp, P. A. & Fisher, D. E. Pre-bending of a promoter
534 sequence enhances affinity for the TATA-binding factor. *Nature* **373**, 724-727 (1995).

535 15 Satchwell, S. C., Drew, H. R. & Travers, A. A. Sequence periodicities in chicken
536 nucleosome core DNA. *Journal of molecular biology* **191**, 659-675 (1986).

537 16 Ngo, T. T., Zhang, Q., Zhou, R., Yodh, J. G. & Ha, T. Asymmetric unwrapping of
538 nucleosomes under tension directed by DNA local flexibility. *Cell* **160**, 1135-1144 (2015).

539 17 Ngo, T. *et al.* Effects of cytosine modifications on DNA flexibility and nucleosome
540 mechanical stability. *Nature communications* **7**, 1-9 (2016).

541 18 Bracco, L., Kotlarz, D., Kolb, A., Diekmann, S. & Buc, H. Synthetic curved DNA
542 sequences can act as transcriptional activators in Escherichia coli. *The EMBO journal* **8**,
543 4289-4296 (1989).

544 19 Rosanio, G., Widom, J. & Uhlenbeck, O. C. In vitro selection of DNA s with an increased
545 propensity to form small circles. *Biopolymers* **103**, 303-320 (2015).

546 20 Beutel, B. A. & Gold, L. In vitro evolution of intrinsically bent DNA. *Journal of molecular*
547 *biology* **228**, 803-812 (1992).

548 21 Greenberg, M. V. & Bourc'his, D. The diverse roles of DNA methylation in mammalian
549 development and disease. *Nature reviews Molecular cell biology* **20**, 590-607 (2019).

550 22 Jones, P. A. Functions of DNA methylation: islands, start sites, gene bodies and beyond.
551 *Nature Reviews Genetics* **13**, 484-492 (2012).

- 552 23 Severin, P. M., Zou, X., Gaub, H. E. & Schulten, K. Cytosine methylation alters DNA
553 mechanical properties. *Nucleic acids research* **39**, 8740-8751 (2011).
- 554 24 Lee, J. Y. & Lee, T.-H. Effects of DNA methylation on the structure of nucleosomes.
555 *Journal of the American Chemical Society* **134**, 173-175 (2012).
- 556 25 Keshet, I., Lieman-Hurwitz, J. & Cedar, H. DNA methylation affects the formation of
557 active chromatin. *Cell* **44**, 535-543 (1986).
- 558 26 Yoo, J., Kim, H., Aksimentiev, A. & Ha, T. Direct evidence for sequence-dependent
559 attraction between double-stranded DNA controlled by methylation. *Nature*
560 *communications* **7**, 1-7 (2016).
- 561 27 Basu, A. *et al.* Measuring DNA mechanics on the genome scale. *Nature* **589**, 462-467
562 (2021).
- 563 28 Protozanova, E., Yakovchuk, P. & Frank-Kamenetskii, M. D. Stacked–unstacked
564 equilibrium at the nick site of DNA. *Journal of molecular biology* **342**, 775-785 (2004).
- 565 29 Okonogi, T., Alley, S., Reese, A., Hopkins, P. & Robinson, B. Sequence-dependent
566 dynamics of duplex DNA: the applicability of a dinucleotide model. *Biophysical Journal*
567 **83**, 3446-3459 (2002).
- 568 30 Olson, W. K., Gorin, A. A., Lu, X.-J., Hock, L. M. & Zhurkin, V. B. DNA sequence-
569 dependent deformability deduced from protein–DNA crystal complexes. *Proceedings of*
570 *the National Academy of Sciences* **95**, 11163-11168 (1998).
- 571 31 El Hassan, M. & Calladine, C. Conformational characteristics of DNA: empirical
572 classifications and a hypothesis for the conformational behaviour of dinucleotide steps.
573 *Philosophical Transactions of the Royal Society of London. Series A: Mathematical,*
574 *Physical and Engineering Sciences* **355**, 43-100 (1997).
- 575 32 Lowary, P. & Widom, J. New DNA sequence rules for high affinity binding to histone
576 octamer and sequence-directed nucleosome positioning. *Journal of molecular biology* **276**,
577 19-42 (1998).
- 578 33 Brogaard, K., Xi, L., Wang, J.-P. & Widom, J. A map of nucleosome positions in yeast at
579 base-pair resolution. *Nature* **486**, 496-501 (2012).
- 580 34 Crothers, D. M., Haran, T. E. & Nadeau, J. G. Intrinsically bent DNA. *J. Biol. Chem* **265**,
581 7093-7096 (1990).
- 582 35 Koo, H.-S., Wu, H.-M. & Crothers, D. M. DNA bending at adenine· thymine tracts. *Nature*
583 **320**, 501-506 (1986).
- 584 36 Hagerman, P. J. Sequence dependence of the curvature of DNA: a test of the phasing
585 hypothesis. *Biochemistry* **24**, 7033-7037 (1985).
- 586 37 Wu, H.-M. & Crothers, D. M. The locus of sequence-directed and protein-induced DNA
587 bending. *Nature* **308**, 509-513 (1984).
- 588 38 Stefl, R., Wu, H., Ravindranathan, S., Sklenář, V. & Feigon, J. DNA A-tract bending in
589 three dimensions: solving the dA4T4 vs. dT4A4 conundrum. *Proceedings of the National*
590 *Academy of Sciences* **101**, 1177-1182 (2004).
- 591 39 Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory
592 epigenome. *Nature Reviews Genetics* **20**, 207-220 (2019).
- 593 40 Jiang, C. & Pugh, B. F. Nucleosome positioning and gene regulation: advances through
594 genomics. *Nature Reviews Genetics* **10**, 161-172 (2009).
- 595 41 Xu, Z. *et al.* Bidirectional promoters generate pervasive transcription in yeast. *Nature* **457**,
596 1033-1037 (2009).

597 42 Krietenstein, N. *et al.* Genomic nucleosome organization reconstituted with pure proteins.
598 *Cell* **167**, 709-721. e712 (2016).

599 43 Oberbeckmann, E. *et al.* Genome information processing by the INO80 chromatin
600 remodeler positions nucleosomes. *Nature communications* **12**, 1-19 (2021).

601 44 Cloutier, T. E. & Widom, J. Spontaneous sharp bending of double-stranded DNA.
602 *Molecular cell* **14**, 355-362 (2004).

603 45 Drew, H. R. & Travers, A. A. DNA bending and its relation to nucleosome positioning.
604 *Journal of molecular biology* **186**, 773-790 (1985).

605 46 Hayes, J. J., Tullius, T. D. & Wolffe, A. P. The structure of DNA in a nucleosome.
606 *Proceedings of the National Academy of Sciences* **87**, 7405-7409 (1990).

607 47 Widlund, H. R. *et al.* Nucleosome structural features and intrinsic properties of the
608 TATAACGCC repeat sequence. *Journal of Biological Chemistry* **274**, 31847-31852
609 (1999).

610 48 Shrader, T. E. & Crothers, D. M. Artificial nucleosome positioning sequences. *Proceedings*
611 *of the National Academy of Sciences* **86**, 7418-7422 (1989).

612 49 Rohs, R. *et al.* The role of DNA shape in protein–DNA recognition. *Nature* **461**, 1248-
613 1253 (2009).

614 50 Zhou, T. *et al.* Quantitative modeling of transcription factor binding specificities using
615 DNA shape. *Proceedings of the National Academy of Sciences* **112**, 4654-4659 (2015).

616 51 Barozzi, I. *et al.* Coregulation of transcription factor binding and nucleosome occupancy
617 through DNA features of mammalian enhancers. *Molecular cell* **54**, 844-857 (2014).

618 52 Li, J. *et al.* Expanding the repertoire of DNA shape features for genome-scale studies of
619 transcription factor binding. *Nucleic acids research* **45**, 12877-12887 (2017).

620 53 El Hassan, M. & Calladine, C. Propeller-twisting of base-pairs and the conformational
621 mobility of dinucleotide steps in DNA. *Journal of molecular biology* **259**, 95-103 (1996).

622 54 Dans, P. D., Perez, A., Faustino, I., Lavery, R. & Orozco, M. Exploring polymorphisms in
623 B-DNA helical conformations. *Nucleic acids research* **40**, 10668-10678 (2012).

624 55 Czapla, L., Swigon, D. & Olson, W. K. Sequence-dependent effects in the cyclization of
625 short DNA. *Journal of chemical theory and computation* **2**, 685-695 (2006).

626 56 Pérez, A. *et al.* Impact of methylation on the physical properties of DNA. *Biophysical*
627 *journal* **102**, 2140-2148 (2012).

628 57 Tippin, D. & Sundaralingam, M. Nine polymorphic crystal structures of d
629 (CCGGGCCCGG), d (CCGGGCCm5CGG), d (Cm5CGGGCCm5CGG) and d
630 (CCGGGCC (Br) 5CGG) in three different conformations: effects of spermine binding and
631 methylation on the bending and condensation of A-DNA. *Journal of molecular biology*
632 **267**, 1171-1185 (1997).

633 58 Moyle-Heyrman, G. *et al.* Chemical map of *Schizosaccharomyces pombe* reveals species-
634 specific features in nucleosome positioning. *Proceedings of the National Academy of*
635 *Sciences* **110**, 20158-20163 (2013).

636 59 Gilchrist, D. A. *et al.* Pausing of RNA polymerase II disrupts DNA-specified nucleosome
637 organization to enable precise gene regulation. *Cell* **143**, 540-551 (2010).

638 60 Garcia, H. G. *et al.* Biological consequences of tightly bent DNA: the other life of a
639 macromolecular celebrity. *Biopolymers: Original Research on Biomolecules* **85**, 115-130
640 (2007).

641 61 Braccioli, L. & de Wit, E. CTCF: a Swiss-army knife for genome organization and
642 transcription regulation. *Essays in biochemistry* **63**, 157-165 (2019).

643 62 Voong, L. N. *et al.* Insights into nucleosome organization in mouse embryonic stem cells
644 through chemical mapping. *Cell* **167**, 1555-1570. e1515 (2016).

645 63 Wiechens, N. *et al.* The chromatin remodelling enzymes SNF2H and SNF2L position
646 nucleosomes adjacent to CTCF and other transcription factors. *PLoS genetics* **12**, e1005940
647 (2016).

648 64 Chen, X. *et al.* Integration of external signaling pathways with the core transcriptional
649 network in embryonic stem cells. *Cell* **133**, 1106-1117 (2008).

650 65 Clarkson, C. T. *et al.* CTCF-dependent chromatin boundaries formed by asymmetric
651 nucleosome arrays with decreased linker length. *Nucleic acids research* **47**, 11181-11196
652 (2019).

653 66 Wang, J.-P. Z. & Widom, J. Improved alignment of nucleosome DNA sequences using a
654 mixture model. *Nucleic acids research* **33**, 6743-6755 (2005).

655 67 Fire, A., Alcazar, R. & Tan, F. Unusual DNA structures associated with germline genetic
656 activity in *Caenorhabditis elegans*. *Genetics* **173**, 1259-1273 (2006).

657 68 Moreno-Herrero, F., Seidel, R., Johnson, S. M., Fire, A. & Dekker, N. H. Structural
658 analysis of hyperperiodic DNA from *Caenorhabditis elegans*. *Nucleic Acids Research* **34**,
659 3057-3066 (2006).

660 69 Pugh, B. F. & Venters, B. J. Genomic organization of human transcription initiation
661 complexes. *PloS one* **11**, e0149339 (2016).

662 70 Kornberg, R. D. The molecular basis of eukaryotic transcription. *Proceedings of the*
663 *National Academy of Sciences* **104**, 12955-12961 (2007).

664 71 Cormack, B. P. & Struhl, K. The TATA-binding protein is required for transcription by all
665 three nuclear RNA polymerases in yeast cells. *Cell* **69**, 685-696 (1992).

666 72 Kim, Y., Geiger, J., Hahn, S. & Sigler, P. B. Crystal structure of a yeast TBP/TATA-box
667 complex. *Nature* **365**, 512-520 (1993).

668 73 Wu, J., Parkhurst, K. M., Powell, R. M., Brenowitz, M. & Parkhurst, L. J. DNA bends in
669 TATA-binding protein· TATA complexes in solution are DNA sequence-dependent.
670 *Journal of Biological Chemistry* **276**, 14614-14622 (2001).

671 74 Rossi, M. J. *et al.* A high-resolution protein architecture of the budding yeast genome.
672 *Nature* **592**, 309-314 (2021).

673 75 Rivetti, C., Guthold, M. & Bustamante, C. Wrapping of DNA around the *E. coli* RNA
674 polymerase open promoter complex. *The EMBO journal* **18**, 4464-4475 (1999).

675 76 Thomason, M. K. *et al.* Global transcriptional start site mapping using differential RNA
676 sequencing reveals novel antisense RNAs in *Escherichia coli*. *Journal of bacteriology* **197**,
677 18-28 (2015).

678 77 Basu, A. *et al.* Dynamic coupling between conformations and nucleotide states in DNA
679 gyrase. *Nature chemical biology* **14**, 565-574 (2018).

680 78 Oram, M., Travers, A. A., Howells, A. J., Maxwell, A. & Pato, M. L. Dissection of the
681 bacteriophage Mu strong gyrase site (SGS): significance of the SGS right arm in Mu
682 biology and DNA gyrase mechanism. *Journal of bacteriology* **188**, 619-632 (2006).

683 79 Oram, M. & Pato, M. L. Mu-like prophage strong gyrase site sequences: analysis of
684 properties required for promoting efficient Mu DNA replication. *Journal of bacteriology*
685 **186**, 4575-4584 (2004).

686 80 Fleischmann, R. D. *et al.* Whole-genome random sequencing and assembly of
687 *Haemophilus influenzae* Rd. *Science* **269**, 496-512 (1995).

- 688 81 Morgan, G. J., Hatfull, G. F., Casjens, S. & Hendrix, R. W. Bacteriophage Mu genome
689 sequence: analysis and comparison with Mu-like prophages in Haemophilus, Neisseria and
690 Deinococcus. *Journal of molecular biology* **317**, 337-359 (2002).
- 691 82 Huo, Y.-X. *et al.* IHF-binding sites inhibit DNA loop formation and transcription initiation.
692 *Nucleic acids research* **37**, 3878-3886 (2009).
- 693 83 Travers, A. DNA-protein interactions: IHF-the master bender. *Current Biology* **7**, R252-
694 R254 (1997).
- 695 84 Revyakin, A., Liu, C., Ebright, R. H. & Strick, T. R. Abortive initiation and productive
696 initiation by RNA polymerase involve DNA scrunching. *Science* **314**, 1139-1143 (2006).
- 697 85 Ma, J., Bai, L. & Wang, M. D. Transcription under torsion. *Science* **340**, 1580-1583 (2013).
- 698 86 Rohs, R. *et al.* Origins of specificity in protein-DNA recognition. *Annual review of*
699 *biochemistry* **79**, 233 (2010).
- 700 87 Ohno, M. *et al.* Sub-nucleosomal genome structure reveals distinct nucleosome folding
701 motifs. *Cell* **176**, 520-534. e525 (2019).
- 702 88 Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals
703 folding principles of the human genome. *science* **326**, 289-293 (2009).
- 704 89 Basu, A. (2022, Aug 29) aakashbasu2/Intrinsic-Cyclizability-Prediction-Codes: v1.0.0 [Computer
705 Software]. Zenodo. <https://doi.org/10.5281/zenodo.7031125>
- 706
- 707 90 Chereji, R. V., Ramachandran, S., Bryson, T. D. & Henikoff, S. Precise genome-wide
708 mapping of single nucleosomes and linkers in vivo. *Genome biology* **19**, 1-20 (2018).
- 709

710

711 **Methods:**

712 **Primer and adapter sequences for loop-seq**

713 25 nt left adapter: 5' - TTTCTTCACTTATCTCCCACCGTCC - 3'

714 25 nt right adapter: 5' - GGCAGAAGACAAGGGAACGAAATAG - 3'

715 Example sequence is a loop-seq library: 5' – Left adapter – 50 nt variable region – Right adapter – 3'

716 Primer P1: 5' – CAGAATCCGTGCGAAGAGCCTTATCTCCCACCGTCC – 3'

717 Primer P2: 5' – ACGGATTCTGCGAAGAGCTTCCCTTG/iBiodT/CTTCTGCC – 3'

718 Primer P2-Biotin26: 5' – ACGGATTCTGCGAAGAGCTTCCCTTG/iBiodT/CTTCTGCC – 3'

719 Primer P2-Biotin29: 5' – ACGGATTCTGCGAAGAGCTTCCCTTGTCT/iBiodT/CTGCC – 3'

720 Primer P2-Biotin31: 5' – ACGGATTCTGCGAAGAGCTTCCCTTGTCTTC/iBiodT/GCC – 3'

721 Primer SP1: 5' – TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCAGAATCCGTC

722 GAAGAGCCTTATCTCCCACCGTCC – 3'

723 Primer SP2: 5' – GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGACGGATTCTGC

724 GAAGAGCTTCCCTTGTCTTCTGCC – 3'

725

726 **Nucleotide sequences**

727 All DNA sequences that were part of various libraries where loop-seq was performed in this study,
728 or were sequences around genomic loci where intrinsic cyclizability values were predicted as part of this
729 study, are reported in Supplementary Tables 1 – 21. Please refer to the Key for a description of each table's
730 content.

731

732 **Loop-seq**

733 Loop-seq was performed according to protocol published earlier⁹¹. Libraries of ssDNA, containing ~90,000
734 different sequences, with the central 50 nt variable and flanked by 25 nt left and right adapters were obtained
735 (Genscript), amplified via emulsion PCR (Micellula DNA Emulsion and Purification Kit (CHIMERx))
736 using primers P1 and P2-Biotin26, using KAPA polymerase (Roche). These primers introduce a biotin
737 molecule, as well as the binding site for the site specific nicking enzyme Nt.BspQ1 (New England Biolabs
738 (NEB)). The amplified library DNA fragments were bound to streptavidin-coated magnetic beads
739 (Dynabeads MyOne Streptavidin T1, Thermo Fisher Scientific). Whenever buffer exchange was required,
740 the magnetic beads were pulled down with rare earth magnets. The library molecules immobilized on the
741 beads were nicked near the ends using Nt.BspQ1 to general 10 nt single-stranded complementary
742 overhangs. The nicking solution was replaced with buffer T2.5BSA (2.375 mM Sodium Chloride, 9.5 mM
743 tris HCl pH 8, 1 mg/ml BSA). Looping was initiated by replacing T2.5BSA with looping buffer (1M
744 Sodium Chloride, 20 mM tris pH 8, 1 mg/ml BSA). After 1 minute incubation in looping buffer, the beads
745 were incubated for 20 minutes in digestion buffer (1x NEB Buffer 4 (NEB), 1 mM ATP, 0.333 units/ml

746 RecBCD (NEB)), which digests those molecules that failed to loop in the prior looping step. Digestion
747 buffer was then replaced with looping buffer and prepared for sequencing. The above protocol was also
748 repeated for an identically prepared control where the digestion enzyme RecBCD was omitted.

749 The libraries (sample and control) were then prepared for sequencing. DNA was PCR amplified
750 off the beads using primers P1 and P2, followed by PCR amplification using primers SQ1 and SQ2,
751 followed by index primers available in the Nexters XT Index Kit (24 Indices 96 Samples) (Illumina).
752 Sequencing was performed on Illumina MiSeq (for libraries with less than 20,000 sequences) or HiSeq
753 platforms. Sequencing results were mapped to the known sequences in the library using Bowtie 1 and
754 SAMtools version 1.12. Cyclizability of a sequence was defined as the ratio of the relative population of
755 the sequence in the amplified library after loop-seq, to that in the control sample, where all steps were
756 identical except that unlooped molecules were not digested. These calculations were performed using
757 simple MATLAB (Matworks) versions 9.0, 9.2, 9.4, and 9.6 scripts that were developed previously⁹¹.

758 For any given library, loop-seq was repeated three times for three different biotin locations²⁷ (where
759 primers P2-Biotin29 and P2-Biotin31 were used in the initial amplification step instead of P2-Biotin26).
760 For each sequence, cyclizability was modeled to have a sinusoidal dependence on the position of the biotin,
761 with periodicity equal to that of the helical repeat. The constant term of the sinusoid was defined as the
762 intrinsic cyclizability²⁷.

763

764 **CpG Methylation of libraries**

765 Libraries to be methylated were first amplified via emulsion PCR as done as part of the regular loop-seq
766 protocol⁹¹. Prior to immobilization on streptavidin-coated magnetic beads, CpG methylation of the
767 amplified library was carried out using CpG Methyltransferase (New England Biolabs, product number
768 M0226S), by following the manufacturer's protocol. The methylation mixture involved amplified DNA
769 library (~120 ng), 10x methyltransferase reaction buffer (New England Biolabs) (5 μ l), 1600 μ M SAM (5
770 μ l), Methyltransferase (4 New England Biolabs units), and water (up to 50 μ l). The reaction mixture was

771 incubated at 37 °C for 1 hour. Heat inactivation was not performed in order to prevent denaturation and
772 incorrect annealing of the library.

773

774 **Statistics and reproducibility**

775 All presented loop-seq data in figures were compiled from a single run of loop-seq on the library in question.

776 The reproducibility of the loop-seq technique has been established earlier⁵³.

777 All Pearson's r, and associated P values, were calculated using the 'corrcoef' function in MATLAB

778 (Matworks), as described earlier⁵³.

779

780 **Data Availability:**

781 All new sequencing data reported as part of this study are deposited in the National Center for
782 Biotechnology Information (NCBI) Sequence Read Archive (SRA) under accession number
783 PRJNA746342. All measurements of intrinsic cyclizabilities obtained from our earlier study²⁷ were based
784 on sequencing data that is deposited in the National Center for Biotechnology Information (NCBI)
785 Sequence Read Archive (SRA) under accession number PRJNA667271.

786 The following datasets used in this study were downloaded from the NCBI Gene Expression
787 omnibus (GEO) using the following accession numbers: GSE97290 (nucleosome occupancy data around
788 +1 nucleosome dyads in *S. cerevisiae*); GSE46957 (nucleosome occupancy data around +1 nucleosome
789 dyads in *S. pombe*); GSE69336 (nucleosome occupancy data in *D. melanogaster* around TSSs); GSE82127
790 (nucleosome occupancy in *M. musculus* around TSSs and CTCF binding sites); GSE11431 (location of
791 CTCF binding sites in mouse embryonic stem cells); GSE147927 (location of TSSs in *S. cerevisiae*);
792 GSE55199 (TSS locations in *E. coli*);

793 *S. cerevisiae* (sacCer3), *S. pombe* (spo2), *D. melanogaster* (BDGP5/ dm3) and *M. musculus*

794 (mm9) genome sequences were downloaded from the UCSC Genome Browser

795 (<https://genome.ucsc.edu/cgi-bin/hgGateway>). The *E. coli* MG1655 genome was downloaded from NCBI

796 (accession number NC_000913.2). The *H. influenza* genome was downloaded from NCBI GeneBank
797 (L42023.1). The sequence of the $\Omega 4$ region of the *C. elegans* genome was downloaded from the
798 supplementary material of this publication: Moreno-Herrero, F., Seidel, R., Johnson, S. M., Fire, A.
799 & Dekker, N. H. Structural analysis of hyperperiodic DNA from *Caenorhabditis elegans*. *Nucleic*
800 *Acids Research* **34**, 3057-3066 (2006)⁶⁸.

801 Supplementary Tables 1 – 21 provide the following data: Sequences and intrinsic cyclizability
802 values in all libraries on which loop-seq was performed either in this study or earlier²⁷; All sequences and
803 predicted intrinsic cyclizability values around all genomic loci where we applied our predictive models to
804 predict intrinsic cyclizability; The values of all parameters that quantify the contribution of short sequence
805 features and their distributions to intrinsic cyclizability, as obtained in this study.

806

807 **Code availability:**

808 Custom MATLAB codes developed as part of this study for predicting intrinsic cyclizability based
809 on linear regression models or neural nets have been deposited in Zenodo⁸⁹: Basu, A. Intrinsic-
810 Cyclizability-Prediction-Codes: v1.0.0. Zenodo. <https://doi.org/10.5281/zenodo.7031125>. (2022, Aug
811 29).

812

813 **Methods only references:**

814 91 Basu, A. in *Methods in Enzymology* Vol. 661 305-326 (Elsevier, 2021).

815

816







