



GC content but not nucleosome positioning directly contributes to intron-splicing efficiency in *Paramecium*

Stefano Gnan, Melody Matelot, Marion Weiman, et al.

Genome Res. published online March 9, 2022

Access the most recent version at doi:[10.1101/gr.276125.121](https://doi.org/10.1101/gr.276125.121)

P<P	Published online March 9, 2022 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

An advertisement banner with a teal background. On the left, the text reads "CRISPR and RNAi Genetic Screening. Your new superpower." In the center, there is a white box with the words "LEARN MORE" in green. On the right, there is a photograph of a woman wearing a red mask and a white cape, and the Cellecta logo, which consists of a green molecular structure and the word "CELLECTA" in white.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **GC content but not nucleosome positioning directly contributes to intron-splicing**
2 **efficiency in *Paramecium***

3

4 Stefano Gnan^{1,4}, Mélody Matelot^{2,4}, Marion Weiman³, Olivier Arnaiz³, Frédéric Guérin², Linda
5 Sperling³, Mireille Bétermier³, Claude Thermes³, Chun-Long Chen^{1*} and Sandra
6 Duharcourt^{2*}

7

8 ¹ Institut Curie, Université PSL, Sorbonne Université, CNRS UMR3244, Dynamics of Genetic
9 Information, Paris, 75005, France

10 ² Université de Paris, CNRS, Institut Jacques Monod, F-75013, Paris, France

11 ³ Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC),
12 91198, Gif-sur-Yvette, France

13 ⁴ joint First Authors

14 * corresponding authors

15 Sandra Duharcourt, Tel: [+33-(0)157278009]; Email: [sandra.duharcourt@ijm.fr]

16 ORCID: <https://orcid.org/0000-0002-8913-8799>

17 Chun-Long Chen, Tel: [+33-(0)156246205]; Fax: [+33-(0)156246674]; Email:
18 [chunlong.chen@curie.fr]

19 ORCID: <https://orcid.org/0000-0002-4795-0295>

20

21 Present Address:

22 [Mélody Matelot], IGBMC - CNRS UMR 7104 - Inserm U 1258, 1 rue Laurent Fries, BP
23 10142, 67404 Illkirch CEDEX, France

24 [Marion Weiman], OncoDNA group, IntegraGen, Evry, 91000, France

25 [Frédéric Guérin], Scipio bioscience, 92120 Montrouge, France

26

27 Running Title: nucleosome position in intron splicing efficiency

28 **Keywords:** Nucleosome positioning, intron splicing, Nonsense Mediated Decay, MNase-seq

29

30 **ABSTRACT**

31

32 Eukaryotic genes are interrupted by introns that must be accurately spliced from mRNA
33 precursors. With an average length of 25 nt, the >90,000 introns of *Paramecium tetraurelia*
34 stand among the shortest introns reported in eukaryotes. The mechanisms specifying the
35 correct recognition of these tiny introns remain poorly understood. Splicing can occur co-
36 transcriptionally and it has been proposed that chromatin structure might influence splice site
37 recognition. To investigate the roles of nucleosome positioning in intron recognition, we
38 determined the nucleosome occupancy along the *P. tetraurelia* genome. We showed that *P.*
39 *tetraurelia* displays a regular nucleosome array with a nucleosome repeat length of ~151 bp,
40 amongst the smallest periodicities reported. Our analysis revealed that introns are frequently
41 associated with inter-nucleosomal DNA, pointing to an evolutionary constraint favoring
42 introns at the AT-rich nucleosome edge sequences. Using accurate splicing efficiency data
43 from cells depleted for nonsense-mediated decay effectors, we showed that introns located
44 at the edge of nucleosomes display higher splicing efficiency than those at the center.
45 However, multiple regression analysis indicated that the low GC content of introns, rather
46 than nucleosome positioning, is associated with high splicing efficiency. Our data reveal a
47 complex link between GC content, nucleosome positioning and intron evolution in
48 *Paramecium*.

49

50 INTRODUCTION

51

52 In eukaryotes, genomic DNA is compacted by histones into chromatin. The basic unit of
53 chromatin is the nucleosome, which comprises a histone octamer made of the four core
54 histones (H2A, H2B, H3 and H4) and 146-147 bp of DNA wrapped around it (Luger et al.
55 1997; Parmar and Padinhateeri 2020). Nucleosomes are not randomly located along the
56 genome but positioned with respect to DNA sequence. The affinity of DNA for histone
57 octamers and the energy needed to bend different DNA fragments around the histone
58 octamer are influenced by the primary DNA sequence, which is therefore an important
59 determinant of nucleosome positioning along the genome (Tillo and Hughes 2009; Tillo et al.
60 2010; Segal et al. 2006; Lorch et al. 2014; Lyer and Struhl 1995; Peckham et al. 2007;
61 Vaillant et al. 2010). Nucleosome positioning is highly dynamic, and its regulation is crucial
62 to control chromatin accessibility, recruitment of chromatin modifiers and transcription factors
63 (Bartholomew 2014; Prendergast and Semple 2011; Kornberg and Lorch 2020; Jiang and
64 Pugh 2009). In most genomes, genes display a Nucleosome Free Region (NFR) at the 5' of
65 their Transcription Start Site (TSS) due to the formation of complexes made by transcription
66 factors around promoter regions (Bernstein et al. 2004; Yuan et al. 2005; Jiang and Pugh
67 2009). A second NFR is present at Transcription Termination Sites (TTS), likely due to the
68 adverse nucleotide composition of the poly(A) signal (Fan et al. 2010; Chereji et al. 2016).
69 Nucleosomes are organized in regular arrays with a periodic distance called the nucleosome
70 repeat length (NRL). Such periodicity is especially evident over gene bodies, and it is
71 species and cell-type specific (Beshnova et al. 2014; Allan et al. 2013). Genome-wide
72 studies have shown that nucleosomes are preferentially positioned in exons compared to
73 introns in diverse organisms including *Schizosaccharomyces pombe*, *Drosophila*, worms and
74 human (Andersson et al. 2009; Schwartz et al. 2009; Tilgner et al. 2009; Nahkuri et al. 2009;
75 Spies et al. 2009; Iannone et al. 2015). Several lines of evidence indicated that a well-
76 positioned nucleosome might slow down RNA polymerase II and favor exon inclusion and
77 alternative splicing (Jonkers et al. 2014; Wilhelm et al. 2011), suggesting a functional role of

78 nucleosome arrays during mRNA maturation. This is in agreement with recent studies
79 showing that intron splicing can occur in a co-transcriptional manner (Brody et al. 2011;
80 Herzel et al. 2017). Some studies have suggested that GC richness at exons, and not
81 nucleosome positioning per se, is important for intron splicing (Amit et al. 2012; Gelfman et
82 al. 2013). Yet, the contribution of nucleosome positioning to intron splicing efficiency has not
83 been investigated thoroughly. The nonsense-mediated decay (NMD) machinery recognizes
84 and degrades transcripts containing premature termination codons (Lykke-Andersen and
85 Jensen 2015; Kurosaki et al. 2019). Therefore, most of the mis-splicing or un-splicing events
86 are removed rapidly by this powerful surveillance mechanism to avoid the production of
87 erroneous proteins. To date, most studies estimated splicing efficiency from NMD-proficient
88 cells, which eliminate most mis-splicing or un-splicing events, and therefore cannot provide a
89 solid evaluation of the intrinsic efficiency of intron splicing.

90 The ciliate *Paramecium tetraurelia* is a unicellular eukaryotic model organism. Like all ciliates,
91 two distinct types of nuclei co-exist within the same cytoplasm in *P. tetraurelia* (Aury et al.
92 2006). The diploid germline micronucleus (MIC) is transcriptionally silent during vegetative
93 growth and transmits the germline genome to sexual progeny through meiosis, while the
94 highly polyploid somatic macronucleus (MAC) is responsible for gene expression (Bétermier
95 and Duharcourt 2014). The >90,000 introns annotated in the MAC genome are among the
96 shortest reported in eukaryotes (18 to 33 nt, 25 nt on average) (Jaillon et al. 2008). How
97 such a large number of tiny introns can be efficiently spliced is not known. In *P. tetraurelia*,
98 no exon skipping has been reported so far (Saudemont et al. 2017; Jaillon et al. 2008).
99 Introns are associated with weak splice signals, as shown by the very low information
100 content of 5' and 3' splice sites, with only the first and last three bases of introns being highly
101 constrained (Jaillon et al. 2008). A strong counter-selection for introns that cannot be
102 detected by the NMD machinery was previously shown, suggesting that introns rely on NMD
103 to compensate for suboptimal splicing efficiency and accuracy (Jaillon et al. 2008;
104 Saudemont et al. 2017). Whether nucleosome positioning or other factors, such as GC

105 content, can regulate splicing efficiency and shape intron evolution in *Paramecium* has not
106 been studied so far.

107 Here, we investigated a possible role of nucleosome positioning in the recognition of introns
108 in *P. tetraurelia*. We mapped the nucleosome occupancy in the somatic macronucleus (MAC)
109 through paired-end MNase-seq. We compared the positioning of nucleosomes with that of
110 introns, whose accurate splicing efficiency data was determined from NMD-depleted cells.

111

112 **RESULTS**

113 **Genome-wide nucleosome position profiling along the *Paramecium* somatic genome**

114 Using MNase-seq, we derived a first nucleosome positioning profile of the macronuclear
115 (MAC) genome of *P. tetraurelia* during vegetative growth. Both chromatin samples and
116 naked MAC DNA controls were digested to mono-nucleosome size (~150 bp, Fig 1A-B and
117 Supplemental Fig S1A). The results obtained from two biological replicates were highly
118 reproducible (Pearson's correlation $R = 0.94$, Supplemental Fig S1B). We therefore
119 combined data from both biological replicates for downstream analyses. All the data
120 presented in the main figures were obtained with the average of two chromatin samples and
121 two naked DNA controls, respectively. The results of each individual sample are reported in
122 the Supplemental Figures. Using the gene annotation, together with the Transcription Start
123 Sites (TSSs) identified by 5' CAP-seq and Transcription Termination Sites (TTSs) identified
124 by poly(A) detection (Arnaiz et al. 2017), we investigated the nucleosome occupancy along
125 transcription units and around their extremities. As described in other eukaryotes, *P.*
126 *tetraurelia* presents an enriched nucleosome density over the transcription units compared to
127 the flanking regions, showing regular arrays of nucleosomes over transcription units (Fig 1C-
128 D and Supplemental Fig S1C-D). As expected, we were able to identify Nucleosome Free
129 Regions (NFRs) upstream of the TSSs of *Paramecium* genes, followed by an array of well-
130 positioned nucleosomes (Fig 1C-D and Supplemental Fig S1C). The analysis of TTSs shows
131 regions with very low nucleosome occupancy downstream of the TTSs and a weakly

132 organized array towards the gene body (Fig 1D and Supplemental Fig S1D). We further
133 separated gene pairs into 3 groups based on their relative orientation: tandem (n=20,233),
134 convergent (n=8,876) and divergent (n=8,867) (Fig 1E and Supplemental Fig S1E-G). We
135 found that nucleosome arrays are clearly visible upstream of the TTSs only when genes are
136 positioned in tandem (Fig 1E and Supplemental Fig S1F), but not in convergent pairs (Fig 1E
137 and Supplemental Fig S1G). This observation suggests that the nucleosome positioning at
138 TTS observed for tandem genes might be due to the downstream TSS, as suggested for
139 *Saccharomyces cerevisiae* (Chereji et al. 2017). Alternatively, convergent genes might be
140 influenced by the transcription readthrough of the gene in the opposite orientation.

141 Based on our nucleosome position calling and using only well-positioned nucleosomes
142 identified in both replicates (see Methods), we calculated the nucleosome repeat length
143 (NRL) (Methods). We found that *P. tetraurelia* displays one of the smallest NRL reported in
144 eukaryotes (150.89 ± 0.57 bp on average, Fig 1F-G and Supplemental Fig S1H-I), close to
145 the 156 ± 2 bp of *S. pombe* (Godde and Widom 1992), which is much smaller than the 167 bp
146 of *S. cerevisiae* (Vaillant et al. 2010) (see Discussion). In human, the NRL within gene
147 bodies is smaller than outside (Valouev et al. 2011). We performed a similar analysis sub-
148 setting nucleosomes based on whether their centers overlap with gene bodies or not. We
149 found a negligible difference between the NRL within gene bodies (151.00 ± 0.94 bp, more
150 than 80% of the analyzed sequences) and those outside of genes (150.29 ± 1.29 bp)
151 (Supplemental Fig S1J).

152

153 **The tiny introns of *Paramecium* genes are frequently associated with inter-** 154 **nucleosomal DNA**

155 We then analyzed nucleosome positioning over gene bodies. In *P. tetraurelia*, exons range
156 from several nucleotides to a few kilobases (Fig 2A, Supplemental Fig S2A for transcription
157 units identified by 5' CAP-seq and poly(A) detection) and are interspersed with tiny introns,
158 the majority spanning between 20 and 35 bp with a median size of 25 bp (Fig 2B). The

159 distribution of exon size shows a peak around 150 bp close to the size of nucleosomes in *P.*
160 *tetraurelia*, which is smaller than the simulated exon size by assuming uniform exon sizes
161 within each gene (Fig 2A, Supplemental Fig S2A). By visual inspection of the nucleosome
162 occupancy profiles, we noticed a tendency of the MNase signal to be stronger over exons
163 leaving the introns preferentially between two nucleosome peaks (Fig 2C, Supplemental Fig
164 S2B for each MNase-digested chromatin sample). This was especially visible when we
165 examined the nucleosome density over introns sorted by the distance of each intron center
166 to the closest nucleosome center (Fig 2D, Supplemental Fig S2C). This distance is
167 significantly higher than what we would expect by calculating the distance of random
168 positions inside gene bodies to the closest nucleosome center (p value $< 10^{-10}$ calculated
169 with Mann-Whitney U test, one sided, alternative H_1 : Intron distance from the closest
170 nucleosome is higher than random chance. Fig 2E). Using this distance, we grouped introns
171 into 3 categories: central, proximal and distal (as illustrated in Fig 2F and Supplemental Fig
172 S2C). We calculated their distribution and compared it with that of exons smaller than 300 bp
173 (roughly the same sample size) categorized in the same way (Fig 2G and Supplemental Fig
174 S2D). Introns were found enriched at distal positions, i.e. located in the regions between two
175 neighbor nucleosomes, compared to exons (45% vs 22%, respectively). In contrast, exons
176 were more enriched in central positions compared to introns (46% vs 28%, respectively).
177 These distributions are statistically significantly different: p value $< 10^{-10}$ calculated with a χ^2
178 test (Fig 2G). Moreover, *P. tetraurelia* exons seem to favor mono-nucleosome length sizes
179 with 35% of exons having sizes comprised between 100 bp and 200 bp. Such a size
180 distribution is significantly shorter than what would be expected if we simulated exon sizes
181 as uniformly distributed within each transcript, in which case only 24% of the exons would
182 fall in this range ($p < 10^{-10}$, Mann-Whitney U test, one sided, alternative H_1 : simulated exons
183 are bigger than real exons) (Fig 2A). Similar results were obtained using only exons of
184 transcription units whose extremities are identified by both 5' CAP-seq and poly(A) detection
185 (Supplemental Fig S2A). This distribution of exon sizes might reflect some selective

186 constraint keeping introns in phase in distal position, i.e. at the edge of the nucleosome.

187

188

189 **Higher splicing efficiency for introns at the edges of nucleosomes**

190 Previous studies have described the effect of nucleosome positioning on mRNA maturation
191 in multiple organisms (Andersson et al. 2009; Schwartz et al. 2009; Tilgner et al. 2009;
192 Nahkuri et al. 2009; Spies et al. 2009; Iannone et al. 2015). To address whether nucleosome
193 positioning affects intron splicing in *P. tetraurelia*, we examined the relationship between
194 nucleosome positioning and intron splicing efficiency, using published datasets from both
195 wild-type (WT) and NMD-depleted cells, which provide a measurement of the splicing
196 efficiency of *P. tetraurelia* introns (Saudemont et al. 2017). Since NMD has been shown to
197 play an important role in removing mis-spliced transcripts and different evolutionary
198 constraints have been observed for NMD-sensitive (presence of a premature termination
199 codon, PTC, after retention) and NMD-insensitive (absence of a PTC after retention) introns
200 (Saudemont et al. 2017), we further divided our 3 positional categories (central, proximal,
201 distal) of introns into NMD-sensitive or NMD-insensitive groups (Fig 3A).

202 First, we observed that the proportion of distal introns is higher in NMD-insensitive introns
203 compared to NMD-sensitive ones, independent of the introduction of a frameshift (3n versus
204 non-3n introns) (Fig 3A). We could not observe statistically significant differences between
205 3n and non-3n NMD-insensitive intron distributions ($p=0.38$, χ^2 test), and only a minor
206 significant increase of distal introns at the expense of central introns and proximal introns
207 can be detected between 3n and non-3n NMD-sensitive introns ($p<10^{-3}$, χ^2 test) (Fig 3A).
208 Since no major differences in the intron distribution between 3n and non-3n introns were
209 observed, we decided to consider only the NMD state for subsequent analyses.

210 According to previous reports, a PTC is more likely to be recognized by the NMD system if it
211 is located far away from the actual termination codon at the 3' end of the gene (Brognna and
212 Wen 2009; Vitali et al. 2019). We reasoned that an NMD-sensitive intron close to the TSS

213 has a higher probability to induce a PTC far away from the actual termination codon.
214 Therefore, we analyzed the distribution of intron positional categories with regard to
215 nucleosomes (distal, central, proximal) as a function of their relative position within genes
216 and of their NMD sensitivity. We found that, for NMD-insensitive introns, the proportion of
217 distal introns is much higher than that of central introns for all distance classes, with only a
218 slight increase of distal intron percentage toward the gene 3' end (Fig 3B and Supplemental
219 Fig S3A). However, for the NMD-sensitive introns, we observed a linear increase of the
220 percentage of distal introns toward the gene 3' end (Fig 3B and Supplemental Fig S3A). This
221 indicates that introns close to the TTS are more frequently associated with distal positions,
222 i.e. at the edge of the nucleosome. To assess whether these introns close to the TTS are
223 less sensitive to NMD, we monitored intron retention rates for the different intron groups.
224 This confirmed that i) the NMD pathway is more efficient for NMD-sensitive introns close to
225 the TSS, i.e. located at the beginning of a gene (Fig 3C), and ii) much higher retention rates
226 in NMD-depleted cells are observed for NMD-sensitive introns located near a TSS compared
227 to those near a TTS (Fig 3C left panel). As expected, no difference can be observed for
228 NMD-insensitive introns (Fig 3C right panel). For NMD-sensitive introns, we observed a
229 higher splicing efficiency (i.e. lower retention rate) for introns located in distal positions
230 compared to those in central and proximal positions independent of their relative position
231 within a gene (Fig 3C, left and Supplemental Fig S3B). We conclude that NMD-sensitive
232 introns located at distal positions, i.e. at the edge of nucleosomes, are more efficiently
233 spliced.

234 As shown in (Saudemont et al. 2017), the intron retention rate is inversely correlated with the
235 gene expression level and is higher for introns that can be detected by the NMD machinery
236 than for those that cannot. In WT cells, both NMD-sensitive and NMD-insensitive introns
237 showed similar retention rates, with higher retention rates for genes with lower expression
238 levels (Fig 3D). The retention rate of NMD-sensitive introns increased significantly upon
239 NMD depletion, while it did not for NMD-insensitive introns (Fig 3D). We extended this
240 analysis to our intron positional categories. As expected, NMD-insensitive introns showed

241 similar splicing efficiency for all intron classes in both WT and NMD-depleted cells (Fig 3D
242 right panel and Supplemental Fig S3C). We found that the retention rate of NMD-sensitive
243 introns is lower for distal introns compared to the other two categories (Fig 3D left and
244 Supplemental Fig S3C), indicating again that NMD-sensitive introns located at the edges of
245 nucleosomes are more efficiently spliced. This can already be observed in WT cells, while in
246 NMD-depleted cells, where nonsense mRNAs are no longer degraded, this difference is
247 much stronger (Fig 3D left panel and Supplemental Fig S3C). For the low-expressed genes
248 ($\text{RPKM} \leq 1$) the retention rate of central introns is 36.6% higher than that of distal introns,
249 and it drops to 24.6% and 13.8% for the mid-expressed ($1 < \text{RPKM} \leq 10$) and highly-
250 expressed ($\text{RPKM} > 10$) genes, respectively (Fig 3D left panel and Supplemental Fig S3C).
251 We further analyzed the proportion of intron positional categories within genes with different
252 expression levels and found similar proportions for all expression classes (Supplemental Fig
253 S3D). Similar results were also observed for genes issued from the last whole-genome
254 duplication (Aury et al. 2006) that have different expression levels (Supplemental Fig S3E).
255 Finally, after controlling both gene expression levels and the relative distance of the intron to
256 the TSS, we still observed that the distal introns have a higher splicing efficiency than the
257 central and proximal ones (Supplemental Fig S3F).

258 It has been shown that the splicing efficiency of *P. tetraurelia* introns depends on the
259 sequences at the donor and acceptor sites (Jaillon et al. 2008). We thus assessed whether
260 this difference in splicing efficiency between our nucleosome-positional classes could be
261 explained by a different distribution of stronger donor (5' GTA) and/or stronger acceptor (3'
262 TAG) sites (Supplemental Fig S3G) within different intron groups. As expected, NMD-
263 insensitive introns were more frequently associated with both stronger donors and acceptors
264 whatever the distance of the intron to the closest nucleosome center (Fig 3E). In contrast,
265 we found a minor increase, for the NMD-sensitive introns, in the association of distal introns
266 with "weaker donor and acceptor" (0.82%) and "stronger acceptor only" (4.46%) intron
267 groups compared to central introns (0.50% and 3.80%, respectively) with respectively 64%

268 and 17% increase (Fig 3E). This slight increase was not associated with a higher retention
269 rate for distal introns compared to central introns. Instead, we did observe a reduced
270 retention rate in distal introns (Fig 3C-D). We conclude that the reduced retention rate in
271 distal introns is not due to a difference of donor/acceptor signals in this class.

272

273 **GC content related to nucleosome positioning contributes to intron-splicing efficiency** 274 **at the edges of nucleosomes**

275 It is well known that nucleosome positioning is highly associated with GC content:
276 nucleosome centers show higher GC than distal regions (Tillo and Hughes 2009; Lorch et al.
277 2014; Lyer and Struhl 1995; Peckham et al. 2007; Vaillant et al. 2010). In *P. tetraurelia*, we
278 observed that NMD-sensitive introns have a higher GC content (18.9%, 18.4% and 15.7%
279 for central, proximal and distal introns, respectively) than their NMD-insensitive counterparts
280 (16.3%, 16.1% and 13.2% for central, proximal and distal introns, respectively) (Fig 4A).
281 Moreover, as we would expect, the central introns have the highest GC content followed by
282 proximal and distal introns (Fig 4A). We therefore analyzed the impact of GC content on
283 intron retention rates. We found a direct correlation between GC percentage and retention
284 rate in NMD-sensitive introns, yet no statistically significant difference could be observed
285 between different intron groups (Fig 4B and Supplemental Fig S4A). This suggests that GC-
286 content anti-correlates with intron splicing efficiency. To further evaluate how different
287 parameters, such as GC content, gene expression level, intron relative position within genes,
288 nucleosome positioning and RNA secondary structure prediction (Supplemental Table S1)
289 affect intron splicing efficiency, we first filtered the parameters by trying to lower the variance
290 inflation factor below 5, and then used the resulting parameters to train a multivariate
291 regression model as previously described (Chen et al. 2010). Only the parameters with a
292 statistically significant contribution were retained (Methods). The final fitted model has an
293 $R=0.62$, which explains 39% of the variation in intron splicing efficiency measured in NMD-

294 depleted cells. The model allowed us to estimate the contribution of each parameter (full list
295 of parameters in Supplemental Table S1).

296 The highest contribution came from the level of gene expression that accounts for about 46%
297 of the model (Fig 4C). The GC content of the intron accounts for 15%, and together with
298 other parameters associated with the base composition of the introns (e.g. TC% accounts for
299 6.2%), the total contribution of base composition reaches ~22%. Although the difference
300 between GC content of an intron and the flanking exons (Δ GC) has been previously reported
301 to be linked to intron splicing efficiency (Amit et al. 2012), Δ GC was not retained in our final
302 model (Supplementary Table S1). As expected, GC content and Δ GC (introns – flanking
303 exons) are highly correlated (Supplementary Fig S4B) and forcing the usage of the latter
304 does not improve the model. Splicing signals account for 9.5% of the model and the
305 parameters associated with the size and base composition of the transcript account for
306 about 8.7%. Intron size and position in the transcript account for 7.1% of the model, followed
307 by intron sensitivity to the NMD pathway (2.5%), the size and base composition of the
308 flanking exons (2.4%), whether an intron is 3n or not (0.66%) and the parameters associated
309 with the formation of secondary structures (0.35%). All the parameters relative to
310 nucleosome positioning account for only 0.43% of the model (Fig 4C-D and Supplemental
311 Table S1). Moreover, if we divide the introns based on their NMD sensitivity, our model can
312 explain 40% of the variation in intron splicing efficiency for the NMD-sensitive introns, while
313 only 27% for the NMD-insensitive ones (Supplementary Fig S4C). We therefore conclude that
314 the GC content, which is tightly linked to nucleosome positioning, contributes to intron-
315 splicing efficiency: a high GC content, which is correlated with high nucleosome occupancy,
316 is associated with low splicing efficiency.

317

318 **DISCUSSION**

319 We have performed the first nucleosome position profiling in the *P. tetraurelia* MAC genome
320 during vegetative growth. Despite its high AT richness (72% AT), the *P. tetraurelia* MAC

321 genome displays a very regular nucleosome positioning pattern as observed in other
322 eukaryotes: NFRs at the TSSs and TTSs, and a regular nucleosome array along genes. An
323 independent study reached the same conclusions (Drews et al. 2021). Unlike *Tetrahymena*
324 *thermophila*, another AT-rich ciliate (78% AT) (with an NRL of 199 bp) (Beh et al. 2015), the
325 NRL in the *P. tetraurelia* MAC genome presents a smaller periodicity (151 ± 1 bp), very similar
326 to that of *S. pombe* (156 bp) (Godde and Widom 1992) and of *Plasmodium falciparum* (155
327 bp) (>80%AT) (Kensche et al. 2015; Silberhorn et al. 2016). This short NRL means that the
328 naked “linker” DNA between nucleosomes in *Paramecium* is extremely small (only a few bp)
329 compared to that of most other eukaryotic genomes, at least tens of bp or even larger
330 (Arceci and Gross 1980). A higher H1/core-histone ratio has been previously reported being
331 associated with a longer NRL (Fan et al. 2003, 2005; Woodcock et al. 2006). For the three
332 eukaryotes with the smallest NRL, *P. tetraurelia*, *P. falciparum* and *S. pombe*, no orthologue
333 of histone H1 has been identified so far. This strongly suggests that the absence of H1 might
334 contribute to the extremely short NRL observed in *Paramecium* chromatin organization in the
335 somatic MAC genome.

336 In yeast and human, actively transcribed genes tend to have shorter NRL than
337 transcriptionally inactive genes, partially due to the binding of H1 generating inaccessible
338 chromatin at inactive genes (Correll et al. 2012; Barbier et al. 2021; Valouev et al. 2011).
339 With the separation of the germline MIC and the somatic MAC genomes in two distinct nuclei,
340 the *Paramecium* MAC genome is characterized by very high coding density. Indeed, >80%
341 of the MAC is covered by annotated genes and 65% of the coding genes are expressed
342 (RNA-seq coverage of at least 1 RPKM) during vegetative growth (Aury et al. 2006; Arnaiz et
343 al. 2017), which might explain the extremely short length and narrow distribution of NRL. A
344 significant difference in the nucleosome organization between MAC and MIC genomes has
345 been reported for *T. thermophila* (Xiong et al. 2016). How nucleosomes are organized in the
346 *Paramecium* MIC genome is unknown. At each sexual cycle of *Paramecium*, the parental
347 MAC is destroyed and the new MIC and MAC are generated from the parental germline MIC

348 (Bétermier and Duhaucourt 2014). During new MAC development, at least 30% of the
349 germline DNA is eliminated during massive genome rearrangements (Guérin et al. 2017;
350 Sellis et al. 2021). A large amount of extremely short (26 to ~1000 bp) non-coding germline
351 sequences, called IESs (Internal Eliminated Sequences), need to be precisely excised to
352 correctly assemble functional genes in the new MAC genome of *Paramecium* species (Sellis
353 et al. 2021). How nucleosome positioning is organized in the germline MIC genome relative
354 to IESs and whether nucleosome positioning and/or GC content might play a role in IES
355 excision are open questions (Coyne et al. 2012; Lhuillier-Akakpo et al. 2014).

356 In multicellular eukaryotes, long introns are recognized through exon definition and
357 nucleosomes positioned along exons might contribute to the exon-intron architecture,
358 possibly pointing to a function in exon definition (Andersson et al. 2009; Schwartz et al. 2009;
359 Tilgner et al. 2009; Nahkuri et al. 2009; Spies et al. 2009; Iannone et al. 2015). By contrast,
360 short introns are recognized through intron definition. With an average length of 25 nt,
361 introns of *P. tetraurelia* are among the shortest reported in eukaryotes (Jaillon et al. 2008).
362 The large number of introns (>90,000) are associated with weak splicing signals. In the
363 current study, we examined the role of nucleosome positioning in intron splicing. We found a
364 regular nucleosome array associated with intron positions within genes, with exons wrapped
365 around nucleosomes and introns frequently located at the edge of nucleosomes. By using
366 the accurate splicing efficiency data obtained from NMD-depleted cells (Saudemont et al.
367 2017), we performed a thorough investigation on the effect of nucleosome positioning on
368 splicing efficiency. We showed that the NMD-sensitive introns located at the edge of
369 nucleosomes display higher splicing efficiency than those at the nucleosome centers.
370 However, we found that this higher splicing efficiency is due to the fact that the introns
371 located at the edges of nucleosomes display lower GC content. Our multiple regression
372 analysis indicated that the nucleosome positioning has a minimal contribution (0.43%) to the
373 intron splicing efficiency (Supplementary Fig S4C and Supplemental Table S1). Our results
374 strongly indicate that GC content, and more broadly intron base composition, rather than

375 nucleosome positioning, directly influences intron splicing efficiency in *Paramecium*. This
376 conclusion may pave the way for future mechanistic studies to decipher how GC content
377 impinges on intron splicing efficiency. Whether the effect of GC content and nucleosome
378 positioning on intron splicing efficiency observed in *Paramecium* can be extended to other
379 eukaryotes remains an open question.

380 We also observed that during evolution, nucleosome positioning has been displaced relative
381 to introns, frequently locating the AT-rich intron sequences at the edge of nucleosomes (Fig
382 4A). Although both NMD-sensitive and NMD-insensitive introns present a higher proportion
383 of distal positions, NMD-insensitive introns show a significantly higher proportion (50% for 3n
384 and 48% for non-3n introns) than NMD-sensitive introns (40% for 3n and 44% for non-3n
385 introns) (Fig 3A). This strongly suggests that the NMD-insensitive introns not located at the
386 AT-rich nucleosome edges, whose retention in transcripts cannot be cleaned up by the NMD
387 pathway, are counter-selected during evolution. Whether introns in *Paramecium* might play a
388 functional role is still unclear. These introns do not seem to contribute to alternative splicing
389 to generate protein diversity or to encode ncRNAs as in large other genomes with long and
390 abundant introns (Lee and Rio 2015; Chen et al. 2003; Ruby et al. 2007). Due to their
391 extremely small size, it seems unlikely that these introns play a role in regulating
392 transcription rate as suggested in recent publications (Fong et al. 2014; Aslanzadeh et al.
393 2018; Alexander et al. 2010). As the parameters analyzed in this study only explain ~40% of
394 the variation in intron splicing efficiency, other parameters remain to be identified and
395 perhaps other models would be necessary to fully understand what intron properties
396 determine splicing efficiency. How such a large number of tiny introns in *Paramecium* is
397 maintained during evolution and how these introns can be efficiently spliced need to be
398 further investigated.

399

400 **METHODS**

401

402 ***Paramecium* strains, cultivation and autogamy**

403 All experiments were carried out with the entirely homozygous wild type strain 51 of *P.*
404 *tetraurelia*. Cells were grown at 27°C in wheat grass powder (WGP) infusion medium
405 bacterized the day before use with *Klebsiella pneumoniae* and supplemented with 0.8
406 mg/mL β -sitosterol (Beisson et al. 2010a, 2010b).

407

408 **Macronuclei preparation**

409 Cells were exponentially grown for 12 divisions then cultures at 1,000 cells/mL were filtered
410 through eight layers of sterile gauze. Cells were collected by low-speed centrifugation (550 g
411 for 1 min) and washed once with 10 mM Tris-HCl pH 7.4. The pellet was diluted 3-fold by
412 addition of lysis buffer (0.25 M sucrose, 10 mM MgCl₂, 10 mM Tris pH 6.8, 0.2% Nonidet P-
413 40) and processed at 4°C as described in (Arnaiz et al. 2012) with some modifications.
414 Briefly, cells were lysed with 10 strokes of a Dounce homogenizer. Particular care was taken
415 to make sure that macronuclei were still intact under the microscope. Washing buffer (0.25
416 M sucrose, 10 mM MgCl₂, 10 mM Tris-HCl pH 7.4) was added to a final volume of 10 times
417 the initial pellet. Macronuclei were collected by centrifugation at 2,000 g for 1 min and
418 washed once in washing buffer. The pellet was diluted 2-fold in 2.1 M sucrose, 10 mM MgCl₂,
419 10 mM Tris pH 7.4 and loaded on top of a 3-mL sucrose layer (2.1 M sucrose, 10 mM MgCl₂,
420 10 mM Tris-HCl pH 7.4) and centrifuged in a swinging rotor for 1 hr at 210,000 g. The
421 macronuclear pellet was washed once, centrifuged at 2,000 g for 1 min and resuspended in
422 washing buffer at 10⁷ nuclei/mL. The macronuclei recovery is quite low, of the order of 10-
423 20%.

424

425 **MNase digestion on chromatin isolated from macronuclei**

426 Samples containing 10⁵ macronuclei were incubated in the digestion buffer (0.25 M sucrose,
427 10 mM MgCl₂, 10 mM Tris pH 7.4, 1 mM CaCl₂) with increasing amounts (0, 0.5, 1, 2, 5, 7.5,
428 10 U) of MNase (Sigma-Aldrich) at 30°C for 10 min. Reactions were stopped by the addition
429 of 3 volumes of 0.5 M EDTA pH 9.0, 1% N-laurylsarcosine (Sigma-Aldrich), 1% SDS, 1

430 mg/mL Proteinase K (Merck) and incubated at 55°C overnight. DNA from each sample was
431 gently extracted once with phenol, and dialyzed twice against TE (10 mM Tris-HCl, 1 mM
432 EDTA at pH 8.0) containing 25% ethanol, and once against TE. Samples were then treated
433 with RNase A and DNA was quantified with a NanoDrop spectrophotometer (Thermo Fisher
434 Scientific) and separated on a 1.2% agarose gel. The reactions containing mostly mono-
435 nucleosomal DNA fragments (see Fig 1) were selected and mono-nucleosomal DNA
436 fragments were purified from 3% low melting-temperature agarose gels and treated with β -
437 agarase (Sigma-Aldrich) for sequencing.

438

439 **MNase digestion on naked DNA**

440 Following purification on a sucrose layer, the macronuclear pellet was washed once,
441 centrifuged at 2,000 g for 1 min, and was resuspended in three volumes of lysis solution (0.5
442 M EDTA at pH 9.0, 1% SDS, 1% N-laurylsarcosine (Sigma-Aldrich), 1 mg/mL of Proteinase
443 K (Merck) then incubated at 55°C overnight. DNA was gently extracted with phenol, and
444 dialyzed twice against TE (10 mM Tris-HCl, 1 mM EDTA at pH 8.0) containing 20% ethanol,
445 and once against Tris 10mM pH 8.0. 1.6 μ g of DNA was digested with increasing amounts of
446 MNase (0 to 1×10^{-3} U) in the digestion buffer at 30°C for 10 min. The reactions were stopped
447 with 250 mM EDTA. The samples were analyzed on a 1.2 % agarose gel and reactions
448 containing fragments of 100-200 bp were gel-purified for DNA sequencing (see
449 Supplemental Fig S1).

450

451 **MNase library preparation and sequencing**

452 Sequencing libraries were generated using the sequencing kit: TruSeq SBS Kit v5 – GA (36
453 Cycle) (FC-104-5001, Illumina). Samples were then sequenced on an Illumina GA-IIx
454 sequencer using paired-end (PE) 74 bp setting. The MNase-seq datasets used in this study
455 are from (Hardy et al. 2021) and are available under accession number PRJEB39679 at the
456 European Nucleotide Archive (ENA: <https://www.ebi.ac.uk/ena>)

457 Alignment was performed using Bowtie 2 (v2.3.3 --local and other default parameters)
458 (Langmead and Salzberg 2012) and mapping to the MAC genome of strain 51 v1.0
459 (ptetraurelia_mac_51.fa), available at ParameciumDB ([https://paramecium.i2bc.paris-
460 saclay.fr/](https://paramecium.i2bc.paris-saclay.fr/)) (Arnaiz et al. 2019).

461

462 **Nucleosome positioning calling**

463 After aligning reads to the reference MAC genome, PCR duplicates with the same start and
464 end positions were removed. Only reads mapped in proper pair with a mapping quality score
465 equal or higher than 30 were kept. Filtering, sorting and filling of mate related flags were
466 performed using Samtools, version 1.9 (Danecek et al. 2021). Bamfiles were converted into
467 bed using BEDTools, version 2.29.2 (Quinlan and Hall 2010) and a customized script. We
468 aimed to use only reads deriving from mono-nucleosomes, therefore, read pairs longer than
469 150 bp and shorter than 75 bp were excluded. We used only the data within the scaffolds
470 larger than 200 kb. A nucleosome score was calculated using the central 75 bp of each read
471 pair. Signal was then smoothed with a gaussian filter and a sigma of 30 over 90 bp for visual
472 assessment of nucleosome position calling. Local maxima and local minima were identified
473 by convoluting the nucleosome score with a first derivative of a gaussian (sigma 30 over ± 90
474 bp). The points of inflection were identified by convoluting the nucleosome score with a
475 kernel containing the second derivative of a gaussian (sigma 30 over ± 90 bp). Peaks were
476 called as a local maximum between two inflection points with opposite inclination. Peaks
477 were called independently in the two chromatin samples, and then a list of well-positioned
478 nucleosomes was compiled using those nucleosomes whose dyad (i.e. center) differs by
479 less than 10 bp between the two biological replicates (about 75% of all nucleosomes). These
480 well-positioned nucleosomes were used for downstream analyses.

481

482 **Computation of Nucleosome Repeat Length**

483 To compute the Nucleosome Repeat Length (NRL), we first calculated the distance of each
484 nucleosome to all the other nucleosomes on the same scaffold, then used the distances

485 obtained to generate the density distribution. This density distribution was then smoothed
486 using a gaussian filter (sigma=10 over ± 30 bp) and local maxima identified convoluting the
487 density distribution with the first derivative of a gaussian (sigma=10 over ± 30 bp). The first n
488 local maxima were then ordered by increasing distances and fitted using a linear model. The
489 slope of the fitted model corresponds to the NRL.

490

491 **Nucleosome distribution calculation**

492 Gene annotation v2.0 of MAC was from (Arnaiz et al. 2019), and the transcription start sites
493 and transcription termination sites were from (Arnaiz et al. 2017). The gene annotations and
494 RNA-seq data are available at ParameciumDB (<https://paramecium.i2bc.paris-saclay.fr/>). To
495 compare with the distribution of real exon sizes, a set of simulated exons was created
496 assuming uniform exon sizes within each gene, i.e. for a given gene with n exons, we
497 divided its total exon length by n to get the length of n simulated exons of the corresponding
498 gene. The NMD data were obtained from (Saudemont et al. 2017), splicing efficiency of each
499 intron was calculated as the Splicing events / Total number of observations (i.e. spliced +
500 unspliced reads). The mean profiles and heatmaps were drawn using a customized script
501 and plotting using Matplotlib (version 3.1.0) (Hunter 2007). All statistical analyses were
502 performed with Python (version 3.7.4, <https://docs.python.org/release/3.7.4/>) using
503 statsmodels (version 0.10.1)(Seabold and Perktold 2010) and SciPy (version 1.3.1)(Virtanen
504 et al. 2020) modules.

505

506 **Multilinear regression**

507 The starting parameters used for the multiple linear regression can be found in
508 Supplemental Table S1. Parameters were transformed using appropriate functions in order
509 to maximize their linearity with intron splicing efficiency, e.g. log transformation of expression
510 levels. Values were then standardized. Variance inflation factors (VIF) were calculated for
511 the whole pool of parameters. If the parameter with highest VIF exceeded the threshold of 5
512 it was excluded from the pool. Parameters VIF was then recalculated and the process

513 repeated until VIF was greater than 5. Parameters from this first selection were used to fit
514 our linear regression model. A randomly selected set of introns (10% of all introns) was kept
515 from the multilinear regression model fitting, and used as a test dataset to evaluate the
516 model performance. We performed a two-sided Z-test for each coefficient with $H_0: C=0$ and
517 $H_1: C \neq 0$. Statistically significant coefficients were then retained and the linear model was
518 trained again with the associated parameters. This step was repeated until the number of
519 variables is stabilized. Estimation of the contribution of each parameter is calculated as in
520 (Chen et al. 2010), which is based on the absolute value of the product of each coefficient
521 and the Pearson's correlation value of its parameter with the splicing efficiency.
522 Contributions were then converted to percentages. Using the intron test dataset, we
523 calculated the Pearson's correlation between real and predicted data. To calculate the
524 Pearson's correlation between prediction and real data divided by NMD-sensitive and NMD-
525 insensitive, all the introns belonging to either group were used. For this part, Python (version
526 3.7.4, <https://docs.python.org/release/3.7.4/>) was used with scikit-learn (version 0.21.3)
527 (Pedregosa et al. 2011), statsmodels (version 0.10.1) (Seabold and Perktold 2010) and
528 SciPy (version 1.3.1) (Virtanen et al. 2020) modules. The full list of parameters can be found
529 in Supplemental Table S1.

530

531 **DATA ACCESS**

532 The customized script and Jupyter notebooks used for this study are available as a
533 Supplemental Code File and on our GitHub page ([https://github.com/CL-CHEN-
534 Lab/Nucleosome](https://github.com/CL-CHEN-Lab/Nucleosome)).

535

536 **COMPETING INTEREST STATEMENT**

537 The authors declare no competing interest.

538

539 **ACKNOWLEDGMENTS**

540 The authors would like to thank Laurent Duret for useful suggestions and discussion, Laurent
541 Duret and Eric Meyer for sharing with us the NMD data, and to acknowledge the high-
542 throughput sequencing facility of I2BC for its sequencing and bioinformatics expertise.

543 LS, MB, CT, CLC and SD conceived and planned the study. MM, FG and SD conducted the
544 experiments. SG, MW, OA and CLC performed the bioinformatics analyses. CT and CC
545 supervised the bioinformatics analyses. SG, CLC and SD wrote the manuscript, and all the
546 authors reviewed it.

547

548 REFERENCES

- 549 Alexander RD, Innocente SA, Barrass JD, Beggs JD. 2010. Splicing-Dependent RNA
550 polymerase pausing in yeast. *Mol Cell* **40**: 582–593.
- 551 Allan J, Fraser RM, Owen-Hughes T, Docherty K, Singh V. 2013. A comparison of in vitro
552 nucleosome positioning mapped with chicken, frog and a variety of yeast core histones.
553 *J Mol Biol* **425**: 4206–4222.
- 554 Amit M, Donyo M, Hollander D, Goren A, Kim E, Gelfman S, Lev-Maor G, Burstein D,
555 Schwartz S, Postolsky B, et al. 2012. Differential GC Content between Exons and
556 Introns Establishes Distinct Strategies of Splice-Site Recognition. *Cell Rep* **1**: 543–556.
- 557 Andersson R, Enroth S, Rada-Iglesias A, Wadelius C, Komorowski J. 2009. Nucleosomes
558 are well positioned in exons and carry characteristic histone modifications. *Genome*
559 *Res* **19**: 1732–41.
- 560 Arceci RJ, Gross PR. 1980. Sea urchin sperm chromatin structure as probed by pancreatic
561 DNase I: Evidence for a novel cutting periodicity. *Dev Biol* **80**: 210–224.
- 562 Arnaiz O, Mathy N, Baudry C, Malinsky S, Aury J-M, Denby Wilkes C, Garnier O, Labadie K,
563 Lauderdale BE, Le Mouël A, et al. 2012. The Paramecium germline genome provides a
564 niche for intragenic parasitic DNA: evolutionary dynamics of internal eliminated
565 sequences. *PLoS Genet* **8**: e1002984.
- 566 Arnaiz O, Meyer E, Sperling L. 2019. ParameciumDB 2019: integrating genomic data across
567 the genus for functional and evolutionary biology. *Nucleic Acids Res* **48**: D599–D605.

- 568 Arnaiz O, Van Dijk E, Bétermier M, Lhuillier-Akakpo M, de Vanssay A, Duharcourt S, Sallet
569 E, Gouzy J, Sperling L. 2017. Improved methods and resources for paramecium
570 genomics: Transcription units, gene annotation and gene expression. *BMC Genomics*
571 **18**: 1–12.
- 572 Aslanzadeh V, Huang Y, Sanguinetti G, Beggs JD. 2018. Transcription rate strongly affects
573 splicing fidelity and cotranscriptionality in budding yeast. *Genome Res* **28**: 203–213.
- 574 Aury J-M, Jaillon O, Duret L, Noel B, Jubin C, Porcel BM, Ségurens B, Daubin V, Anthouard
575 V, Aiach N, et al. 2006. Global trends of whole-genome duplications revealed by the
576 ciliate *Paramecium tetraurelia*. *Nature* **444**: 171–178.
- 577 Barbier J, Vaillant C, Volff J-NN, Brune FG, Audit B, Brunet F, Audit B. 2021. Coupling
578 Between Sequence-Mediated Nucleosome Organization and Genome Evolution.
579 *Genes (Basel)* **12**: 1–23.
- 580 Bartholomew B. 2014. Regulating the Chromatin Landscape: Structural and Mechanistic
581 Perspectives. *Annu Rev Biochem* **83**: 671–696.
- 582 Beh LY, Müller MM, Muir TW, Kaplan N, Landweber LF. 2015. DNA-guided establishment of
583 nucleosome patterns within coding regions of a eukaryotic genome. *Genome Res* **25**:
584 1727–38.
- 585 Beisson J, Bétermier M, Bré MH, Cohen J, Duharcourt S, Duret L, Kung C, Malinsky S,
586 Meyer E, Preer JR, et al. 2010a. Maintaining clonal *paramecium tetraurelia* cell lines of
587 controlled age through daily reisolation. *Cold Spring Harb Protoc* **5**: pdb.prot5361.
- 588 Beisson J, Bétermier M, Bré MH, Cohen J, Duharcourt S, Duret L, Kung C, Malinsky S,
589 Meyer E, Preer JR, et al. 2010b. Mass culture of *paramecium tetraurelia*. *Cold Spring*
590 *Harb Protoc* **5**: pdb.prot5362.
- 591 Bernstein BE, Liu CL, Humphrey EL, Perlstein EO, Schreiber SL. 2004. Global nucleosome
592 occupancy in yeast. *Genome Biol* **5**: R62.
- 593 Beshnova DA, Cherstvy AG, Vainshtein Y, Teif VB. 2014. Regulation of the Nucleosome
594 Repeat Length In Vivo by the DNA Sequence, Protein Concentrations and Long-Range
595 Interactions ed. Ioshikhes. *PLoS Comput Biol* **10**: e1003698.

- 596 Bétermier M, Duharcourt S. 2014. Programmed Rearrangement in Ciliates: Paramecium.
597 *Microbiol Spectr* **2**: MDNA3-0035–2014.
- 598 Brody Y, Neufeld N, Bieberstein N, Causse SZ, Böhnlein EM, Neugebauer KM, Darzacq X,
599 Shav-Tal Y. 2011. The in vivo kinetics of RNA polymerase II elongation during co-
600 transcriptional splicing. *PLoS Biol* **9**: e1000573.
- 601 Brogna S, Wen J. 2009. Nonsense-mediated mRNA decay (NMD) mechanisms. *Nat Struct*
602 *Mol Biol* **16**: 107–113.
- 603 Chen C-L, Rappailles A, Duquenne L, Huvet M, Guilbaud G, Farinelli L, Audit B,
604 D'Aubenton-Carafa Y, Arneodo A, Hyrien O, et al. 2010. Impact of replication timing on
605 non-CpG and CpG substitution rates in mammalian genomes. *Genome Res* **20**: 447–57.
- 606 Chen CL, Liang D, Zhou H, Zhuo M, Chen YQ, Qu LH. 2003. The high diversity of snoRNAs
607 in plants: Identification and comparative study of 120 snoRNA genes from *Oryza sativa*.
608 *Nucleic Acids Res* **31**: 2601–2613.
- 609 Chereji R V., Kan TW, Grudniewska MK, Romashchenko A V., Berezikov E, Zhimulev IF,
610 Guryev V, Morozov A V., Moshkin YM. 2016. Genome-wide profiling of nucleosome
611 sensitivity and chromatin accessibility in *Drosophila melanogaster*. *Nucleic Acids Res*
612 **44**: 1036–1051.
- 613 Chereji R V., Ocampo J, Clark DJ. 2017. MNase-Sensitive Complexes in Yeast:
614 Nucleosomes and Non-histone Barriers. *Mol Cell* **65**: 565-577.e3.
- 615 Correll SJ, Schubert MH, Grigoryev SA. 2012. Short nucleosome repeats impose rotational
616 modulations on chromatin fibre folding. *EMBO J* **31**: 2416–2426.
- 617 Coyne RS, Lhuillier-Akakpo M, Duharcourt S. 2012. RNA-guided DNA rearrangements in
618 ciliates: is the best genome defence a good offence? *Biol cell* **104**: 309–25.
- 619 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
620 McCarthy SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools.
621 *Gigascience* **10**: 1–4.
- 622 Drews F, Salhab A, Karunanithi S, Cheaib M, Jung M, Schulz MH, Simon M. 2021. Broad
623 domains of histone marks in the highly compact *Paramecium* macronuclear genome.

- 624 *bioRxiv* 2021.08.05.454756.
- 625 Fan X, Moqtaderi Z, Jin Y, Zhang Y, Liu XS, Struhl K. 2010. Nucleosome depletion at yeast
626 terminators is not intrinsic and can occur by a transcriptional mechanism linked to 3'-
627 end formation. *Proc Natl Acad Sci U S A* **107**: 17945–17950.
- 628 Fan Y, Nikitina T, Morin-Kensicki EM, Zhao J, Magnuson TR, Woodcock CL, Skoultchi AI.
629 2003. H1 Linker Histones Are Essential for Mouse Development and Affect
630 Nucleosome Spacing In Vivo. *Mol Cell Biol* **23**: 4559–4572.
- 631 Fan Y, Nikitina T, Zhao J, Fleury TJ, Bhattacharyya R, Bouhassira EE, Stein A, Woodcock
632 CL, Skoultchi AI. 2005. Histone H1 depletion in mammals alters global chromatin
633 structure but causes specific changes in gene regulation. *Cell* **123**: 1199–1212.
- 634 Fong N, Kim H, Zhou Y, Ji X, Qiu J, Saldi T, Diener K, Jones K, Fu XD, Bentley DL. 2014.
635 Pre-mRNA splicing is facilitated by an optimal RNA polymerase II elongation rate.
636 *Genes Dev* **28**: 2663–2676.
- 637 Gelfman S, Cohen N, Yearim A, Ast G. 2013. DNA-methylation effect on cotranscriptional
638 splicing is dependent on GC architecture of the exon-intron structure. *Genome Res* **23**:
639 789–799.
- 640 Godde JS, Widom J. 1992. Chromatin structure of *Schizosaccharomyces pombe*. A
641 nucleosome repeat length that is shorter than the chromatosomal DNA length. *J Mol*
642 *Biol* **226**: 1009–1025.
- 643 Guérin F, Arnaiz O, Boggetto N, Denby Wilkes C, Meyer E, Sperling L, Duharcourt S. 2017.
644 Flow cytometry sorting of nuclei enables the first global characterization of *Paramecium*
645 germline DNA and transposable elements. *BMC Genomics* **18**: 327.
- 646 Hardy A, Matelot M, Touzeau A, Klopp C, Lopez-Roques C, Duharcourt S, Defrance M.
647 2021. DNAModAnnot: A R toolbox for DNA modification filtering and annotation ed. P.
648 Robinson. *Bioinformatics* **37**: 2738–2740.
- 649 Herzel L, Ottoz DSM, Alpert T, Neugebauer KM. 2017. Splicing and transcription touch base:
650 Co-transcriptional spliceosome assembly and function. *Nat Rev Mol Cell Biol* **18**: 637–
651 650.

- 652 Hunter JD. 2007. Matplotlib: A 2D graphics environment. *Comput Sci Eng* **9**: 90–95.
- 653 Iannone C, Pohl A, Papasaikas P, Soronellas D, Vicent GP, Beato M, Valcárcel J. 2015.
654 Relationship between nucleosome positioning and progesterone-induced alternative
655 splicing in breast cancer cells. *RNA* **21**: 360–374.
- 656 Jaillon O, Bouhouche K, Gout JF, Aury JM, Noel B, Saudemont B, Nowacki M, Serrano V,
657 Porcel BM, Ségurens B, et al. 2008. Translational control of intron splicing in
658 eukaryotes. *Nature* **451**: 359–362.
- 659 Jiang C, Pugh BF. 2009. Nucleosome positioning and gene regulation: Advances through
660 genomics. *Nat Rev Genet* **10**: 161–172.
- 661 Jonkers I, Kwak H, Lis JT. 2014. Genome-wide dynamics of Pol II elongation and its
662 interplay with promoter proximal pausing, chromatin, and exons. *Elife* **2014**: e02407.
- 663 Kensche PR, Hoeijmakers WAM, Toenhake CG, Bras M, Chappell L, Berriman M, Bártfai R.
664 2015. The nucleosome landscape of *Plasmodium falciparum* reveals chromatin
665 architecture and dynamics of regulatory sequences. *Nucleic Acids Res* **44**: 2110–24.
- 666 Kornberg RD, Lorch Y. 2020. Primary Role of the Nucleosome. *Mol Cell* **79**: 371–375.
- 667 Kurosaki T, Popp MW, Maquat LE. 2019. Quality and quantity control of gene expression by
668 nonsense-mediated mRNA decay. *Nat Rev Mol Cell Biol* **20**: 406–420.
- 669 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
670 *2012* **9**: 357–359.
- 671 Lee Y, Rio DC. 2015. Mechanisms and regulation of alternative Pre-mRNA splicing. *Annu*
672 *Rev Biochem* **84**: 291–323.
- 673 Lhuillier-Akakpo M, Frapporti A, Denby Wilkes C, Matelot M, Vervoort M, Sperling L,
674 Duharcourt S. 2014. Local effect of enhancer of zeste-like reveals cooperation of
675 epigenetic and cis-acting determinants for zygotic genome rearrangements. *PLoS*
676 *Genet* **10**: e1004665.
- 677 Lorch Y, Maier-Davis B, Kornberg RD. 2014. Role of DNA sequence in chromatin
678 remodeling and the formation of nucleosome-free regions. *Genes Dev* **28**: 2492–2497.
- 679 Luger K, Mäder AW, Richmond RK, Sargent DF, Richmond TJ, Mäder AW, Richmond RK,

- 680 Sargent DF, Richmond TJ. 1997. Crystal structure of the nucleosome core particle at
681 2.8 Å resolution. *Nature* **389**: 251–260.
- 682 Lyer V, Struhl K. 1995. Poly(dA:dT), a ubiquitous promoter element that stimulates
683 transcription via its intrinsic DNA structure. *EMBO J* **14**: 2570–2579.
- 684 Lykke-Andersen S, Jensen TH. 2015. Nonsense-mediated mRNA decay: An intricate
685 machinery that shapes transcriptomes. *Nat Rev Mol Cell Biol* **16**: 665–677.
- 686 Nahkuri S, Taft RJ, Mattick JS. 2009. Nucleosomes are preferentially positioned at exons in
687 somatic and sperm cells. *Cell Cycle* **8**: 3420–3424.
- 688 Parmar JJ, Padinhateeri R. 2020. Nucleosome positioning and chromatin organization. *Curr*
689 *Opin Struct Biol* **64**: 111–118.
- 690 Peckham HE, Thurman RE, Fu Y, Stamatoyannopoulos JA, Noble WS, Struhl K, Weng Z.
691 2007. Nucleosome positioning signals in genomic DNA. *Genome Res* **17**: 1170–1177.
- 692 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M,
693 Prettenhofer P, Weiss R, Dubourg V, et al. 2011. Scikit-learn: Machine Learning in
694 Python. *J Mach Learn Res* **12**: 2825–2830.
- 695 Prendergast JGD, Semple CAM. 2011. Widespread signatures of recent selection linked to
696 nucleosome positioning in the human lineage. *Genome Res* **21**: 1777–1787.
- 697 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
698 features. *Bioinformatics* **26**: 841–842.
- 699 Ruby JG, Jan CH, Bartel DP. 2007. Intronic microRNA precursors that bypass Drosha
700 processing. *Nature* **448**: 83–86.
- 701 Saudemont B, Popa A, Parmley JL, Rocher V, Blugeon C, Necsulea A, Meyer E, Duret L.
702 2017. The fitness cost of mis-splicing is the main determinant of alternative splicing
703 patterns. *Genome Biol* **18**: 1–15.
- 704 Schwartz S, Meshorer E, Ast G. 2009. Chromatin organization marks exon-intron structure.
705 *Nat Struct Mol Biol* **16**: 990–5.
- 706 Seabold S, Perktold J. 2010. Statsmodels: Econometric and Statistical Modeling with Python.
707 In *Proceedings of the 9th Python in Science Conference* (eds. S. Van der Walt and J.

- 708 Millman), Vol. SCIPY 2010 of, pp. 92–96.
- 709 Segal E, Fondufe-Mittendorf Y, Chen L, Thåström A, Field Y, Moore IK, Wang J-PZ, Widom
710 J. 2006. A genomic code for nucleosome positioning. *Nature* **442**: 772–778.
- 711 Sellis D, Guérin F, Arnaiz O, Pett W, Lerat E, Boggetto N, Krenek S, Berendonk T, Couloux
712 A, Aury J-M, et al. 2021. Massive colonization of protein-coding exons by selfish
713 genetic elements in Paramecium germline genomes ed. H.S. Malik. *PLOS Biol* **19**:
714 e3001309.
- 715 Silberhorn E, Schwartz U, Löffler P, Schmitz S, Symelka A, de Koning-Ward T, Merkl R,
716 Längst G. 2016. Plasmodium falciparum Nucleosomes Exhibit Reduced Stability and
717 Lost Sequence Dependent Nucleosome Positioning ed. K.W. Deitsch. *PLOS Pathog* **12**:
718 e1006080.
- 719 Spies N, Nielsen CB, Padgett RA, Burge CB. 2009. Biased Chromatin Signatures around
720 Polyadenylation Sites and Exons. *Mol Cell* **36**: 245–254.
- 721 Tilgner H, Nikolaou C, Althammer S, Sammeth M, Beato M, Valcárcel J, Guigó R. 2009.
722 Nucleosome positioning as a determinant of exon recognition. *Nat Struct Mol Biol* **16**:
723 996–1001.
- 724 Tillo D, Hughes TR. 2009. G+C content dominates intrinsic nucleosome occupancy. *BMC*
725 *Bioinformatics* **10**: 442.
- 726 Tillo D, Kaplan N, Moore IK, Fondufe-Mittendorf Y, Gossett AJ. 2010. High Nucleosome
727 Occupancy Is Encoded at Human Regulatory Sequences. *PLoS One* **5**: 9129.
- 728 Vaillant C, Palmeira L, Chevereau G, Audit B, D'Aubenton-Carafa Y, Thermes C, Arneodo A.
729 2010. A novel strategy of transcription regulation by intragenic nucleosome ordering.
730 *Genome Res* **20**: 59–67.
- 731 Valouev A, Johnson SM, Boyd SD, Smith CL, Fire AZ, Sidow A. 2011. Determinants of
732 nucleosome organization in primary human cells. *Nature* **474**: 516–20.
- 733 Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E,
734 Peterson P, Weckesser W, Bright J, et al. 2020. SciPy 1.0: fundamental algorithms for
735 scientific computing in Python. *Nat Methods* **17**: 261–272.

- 736 Vitali V, Hagen R, Catania F. 2019. Environmentally induced plasticity of programmed DNA
 737 elimination boosts somatic variability in *Paramecium tetraurelia*. *Genome Res* **29**:
 738 1693–1704.
- 739 Wilhelm BT, Marguerat S, Aligianni S, Codlin S, Watt S, Bähler J. 2011. Differential patterns
 740 of intronic and exonic DNA regions with respect to RNA polymerase II occupancy,
 741 nucleosome density and H3K36me3 marking in fission yeast. *Genome Biol* **12**: R82.
- 742 Woodcock CL, Skoultchi AI, Fan Y. 2006. Role of linker histone in chromatin structure and
 743 function: H1 stoichiometry and nucleosome repeat length. *Chromosom Res* **14**: 17–25.
- 744 Xiong J, Gao S, Dui W, Yang W, Chen X, Taverna SD, Pearlman RE, Ashlock W, Miao W,
 745 Liu Y. 2016. Dissecting relative contributions of cis-and trans-determinants to
 746 nucleosome distribution by comparing *Tetrahymena* macronuclear and micronuclear
 747 chromatin. *Nucleic Acids Res* **44**: 10091–10105.
- 748 Yuan G-C, Liu Y-J, Dion MF, Slack MD, Wu LF, Altschuler SJ, Rando OJ. 2005. Genome-
 749 Scale Identification of Nucleosome Positions in *S. cerevisiae*. *Science (80-)* **309**: 626–
 750 630.

751
 752

753 **FIGURE LEGENDS**

754 **Figure 1. Nucleosome occupancy along the *Paramecium* MAC genome. (A)** Schematic
 755 representation of the MNase-seq experiment. **(B)** MNase digestion of MAC chromatin with
 756 increasing MNase enzyme concentration. **(C)** Heatmap showing nucleosome occupancy \pm 1
 757 kb around the center of each gene ordered by gene size (small genes on top and large
 758 genes at the bottom) for 38,143 genes located on scaffolds that are at least 200 kb long. Left
 759 panel: average of two chromatin treated samples (Chromatin). Right panel: average of two
 760 naked DNA control samples (Naked DNA). **(D)** Average nucleosome occupancy around
 761 Transcription Start Sites (TSSs) identified by 5' CAP-seq on the left, and Transcription
 762 Termination Sites (TTSs) identified by poly(A) detection on the right: in green, the average
 763 profile of chromatin treated sample (Chromatin); in blue, average profile of naked DNA
 764 treated sample (Naked DNA); and in magenta the Chromatin/Naked DNA ratio, enrichment
 765 of which is shown on the second axis on the right (red axis). **(E)** Average nucleosome
 766 occupancy \pm 1 kb around the center of intergenic regions: same color code as in panel D.

767 Intragenic regions have been divided into three groups based on the relative positions of
 768 gene pairs: tandem (left), convergent (middle) or divergent (right). **(F)** Inter-center distance
 769 between well-positioned nucleosomes (Methods) on the same scaffold. In blue, distance
 770 distributions from actual data (from 1 bp to 2 kb, binning=1); and in orange, the gaussian
 771 smoothed signal. Black dashed lines indicate the local maxima (peak centers) of the
 772 smoothed data (Methods). **(G)** In orange, the first 8 local maxima from panel F ordered by
 773 increasing distance, and in blue the linear fitted model. At the bottom right, information about
 774 linear fitting and estimated NRL (Mean \pm SD) is given. P-value is calculated using a two-sided
 775 Z-test.

776

777 **Figure 2. Inter-nucleosomal DNA is frequently associated with intron position. (A)**
 778 Histogram showing exon size distribution (bin size = 25 bp): in blue, real exons; and in
 779 orange, simulated exons created assuming uniform exon sizes within each gene. **(B)**
 780 Histogram showing intron size distribution (bin size = 1 bp). **(C)** Example track reporting
 781 nucleosome occupancy over genes with intron locations indicated by vertical dashed lines.
 782 We can observe nucleosome free regions (NFRs) around the gene promoters and introns
 783 frequently associated with inter-nucleosomal DNA. **(D)** Heatmap showing nucleosome
 784 occupancy \pm 200 bp around intron centers. Introns are ordered based on increasing
 785 distances from their center to the closest nucleosome center, from top to bottom. The
 786 average of the chromatin samples is shown on the left and the average of the naked DNA
 787 samples on the right, with the same color code as in Fig 1C. Vertical black dashed lines
 788 delineate the average size of an intron (25 bp). Individual samples are displayed in
 789 Supplemental Fig S2C. **(E)** Histogram reporting the distance of an intron center to the
 790 closest nucleosome center (red). For each intron, a random position inside the
 791 corresponding gene body was selected and the distance to its closest nucleosome center is
 792 reported (green). Bin size= 5 bp. **(F)** Schematic representation of the criteria to assign
 793 features for each intron (or exon) into one of the 3 classes, based on the distance (d)
 794 between its center and the closest nucleosome center position: central, $d \leq 25$ bp; proximal,
 795 $25 \text{ bp} < d < 50$ bp; distal $50 \text{ bp} \leq d \leq 75$ bp. **(G)** Relative distribution of introns, exons and
 796 both features over categories defined in panel F for the introns overlapping with a fixed
 797 nucleosome (about 70% of all introns, see Methods) and exons with a size below 300 bp
 798 overlapping with fixed peaks. See Supplemental Fig S2D including the features with $d > 75$
 799 bp.

800

801 **Figure 3. Nucleosome positioning is associated with intron-splicing efficiency. (A)**
 802 Relative distribution of different classes of introns. Introns are grouped based on their length
 803 (3n or non-3n) and whether their retention causes a premature termination codon making

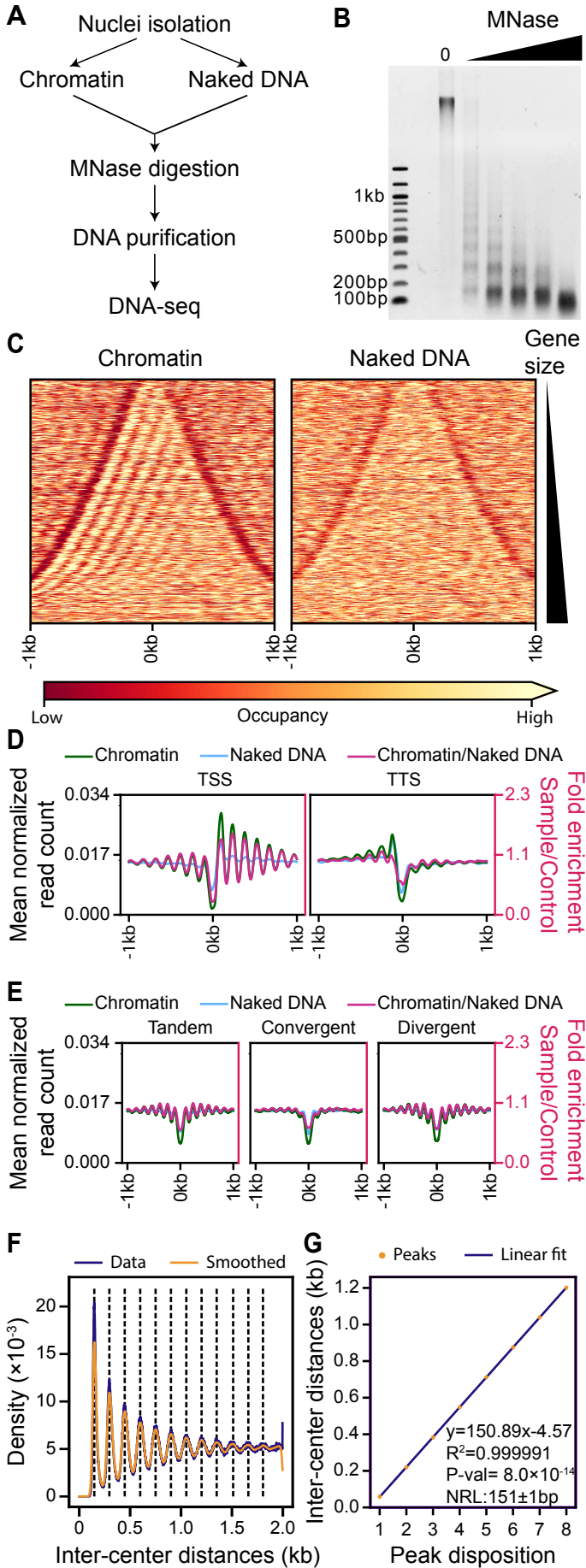
804 them sensitive to the nonsense-mediated decay mechanism (NMD-sensitive) or not (NMD-
 805 insensitive). Within each group, introns are classified based on the distance to the closest
 806 nucleosome center as in Fig 2G. P-values are calculated using the χ^2 test and only the
 807 significant ones are indicated. **(B)** Intron repartition according to the categories defined in Fig
 808 2F as a function of their relative position within a gene. Introns are grouped based on their
 809 NMD sensitivity. Bin size= 20%. A barplot representation with relative p-values is displayed
 810 in Supplemental Fig S3A. **(C)** The retention rate of introns in WT (dashed lines) and in NMD-
 811 depleted (NMD^{KD}, solid lines) cells as a function of their relative position within a gene.
 812 Introns are grouped as in panel B. Error bars represent the standard error of the mean. P-
 813 values calculated using Mann-Whitely U test, and adjusted with false discovery rate (5%),
 814 are displayed in Supplemental Fig S3B. Bin size= 20%. **(D)** The retention rate of introns in
 815 WT (dashed lines) and in NMD-depleted (NMD^{KD}, solid lines) cells as a function of gene
 816 expression levels. Error bars represent the standard error of the mean. Colors and groups
 817 are as in panel B. P-values calculated using Mann-Whitely U test, and adjusted with false
 818 discovery rate (5%), are displayed in Supplemental Fig S3C. **(E)** Relative characterization of
 819 introns, within the same categories as in panel B, based on the strength of splicing acceptor
 820 and donor sites. P-values are calculated using the χ^2 test and adjusted with false discovery
 821 rate (5%). Tests were run between introns belonging to the same positional group or
 822 between introns belonging to the same NMD group. (P-value in all the plots * <0.05 , ** $< 10^{-2}$,
 823 *** $< 10^{-3}$, **** $< 10^{-4}$, ***** $< 10^{-5}$, ***** $< 10^{-6}$).

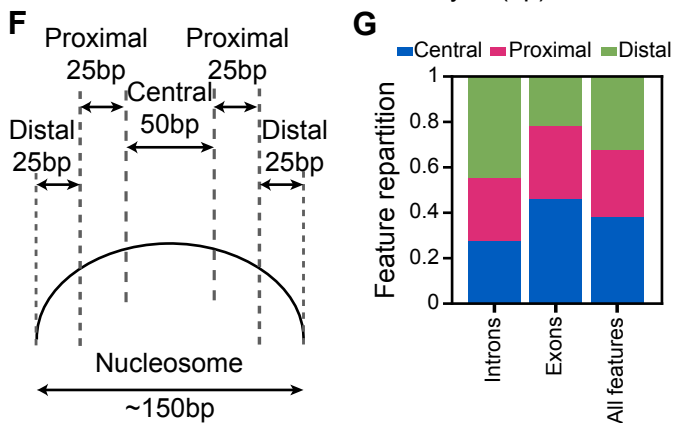
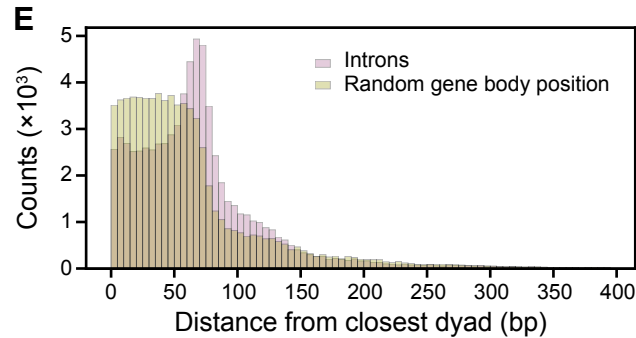
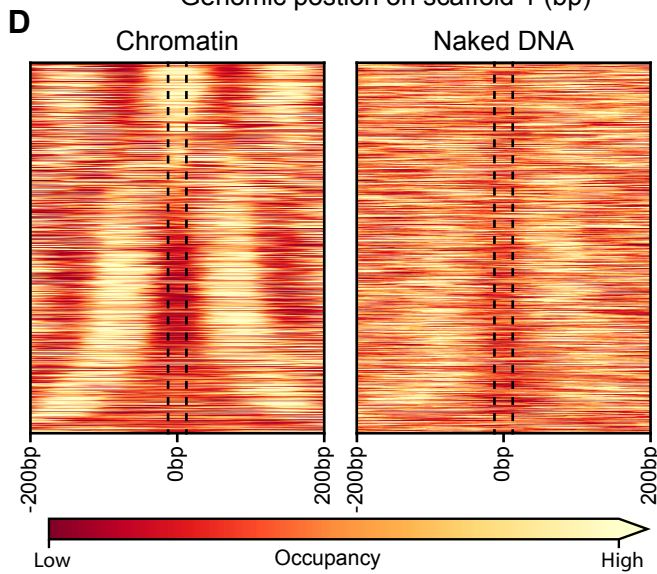
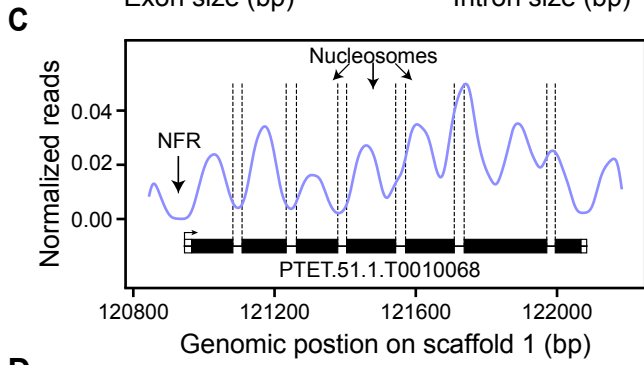
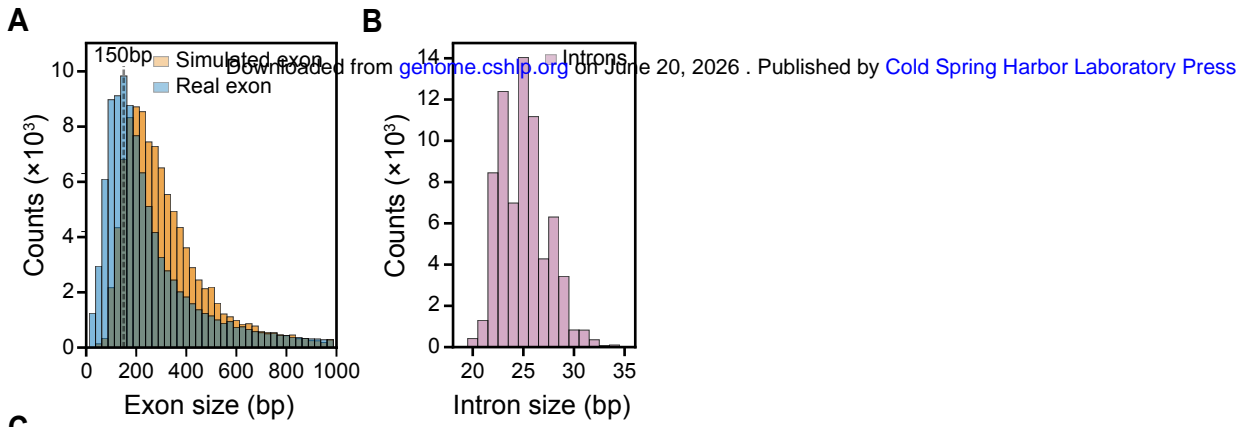
824

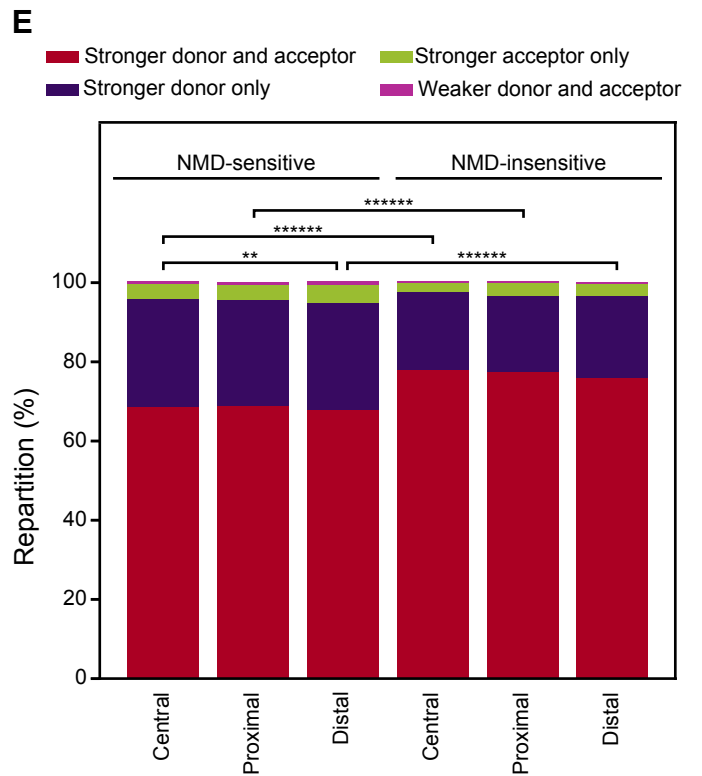
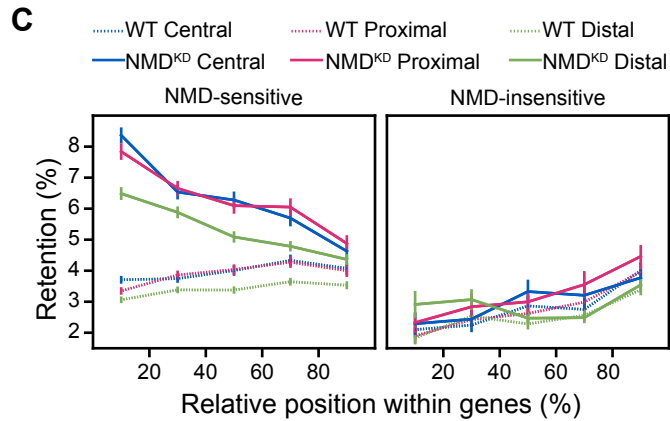
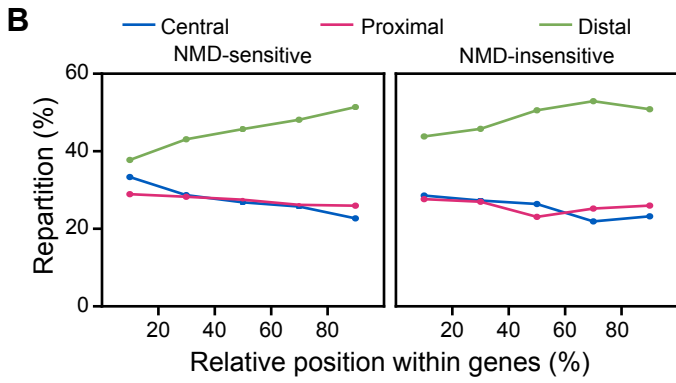
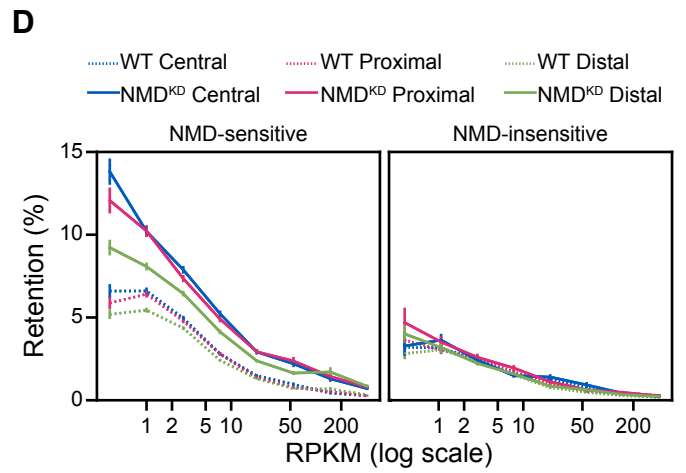
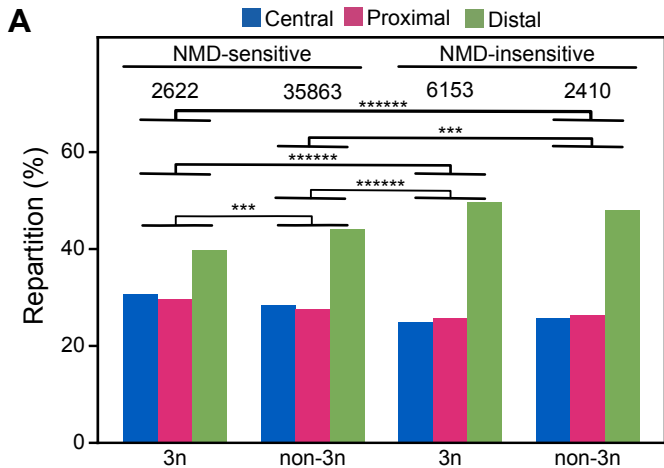
825 **Figure 4. GC content related to nucleosome positioning contributes to intron-splicing**
 826 **efficiency. (A)** GC content (%) distribution of introns based on the distance to the closest
 827 nucleosome center and NMD sensitivity. Mean and standard deviation for each group is
 828 reported at the bottom. P-values were calculated using the Mann-Whitney U test and
 829 adjusted using the false discovery rate (5%). Tests were run between introns belonging to
 830 the same positional group or between introns belonging to the same NMD group (P-value
 831 * <0.05 , ** $< 10^{-2}$, *** $< 10^{-3}$, **** $< 10^{-4}$, ***** $< 10^{-5}$, ***** $< 10^{-6}$). **(B)** The retention rate of introns
 832 in WT and NMD-depleted (NMD^{KD}) cells as a function of their GC content (excluded GT and
 833 AG dinucleotides at both extremities). Introns are classified based on their distance to the
 834 closest nucleosome center and on whether they are NMD sensitive or not. Binning = 10%.
 835 Error bars represent the standard error of the mean. P-values calculated using the Mann-
 836 Whitney U test, and adjusted using the false discovery rate (5%), are displayed in
 837 Supplemental Fig S4A. **(C)** Modelling Splicing Efficiency (SE) in NMD-depleted cells: the pie
 838 chart reports the contribution of each parameter or group of parameters used in the final
 839 model. The full list of retained parameters, reporting their contribution and their statistical
 840 significance, is displayed in Supplemental Table S1 as well as in Supplemental Fig S4C. **(D)**

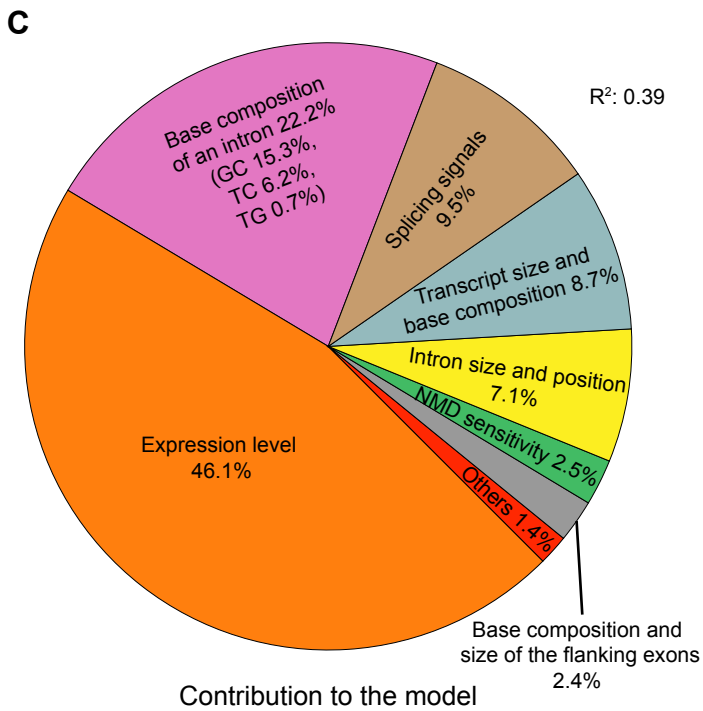
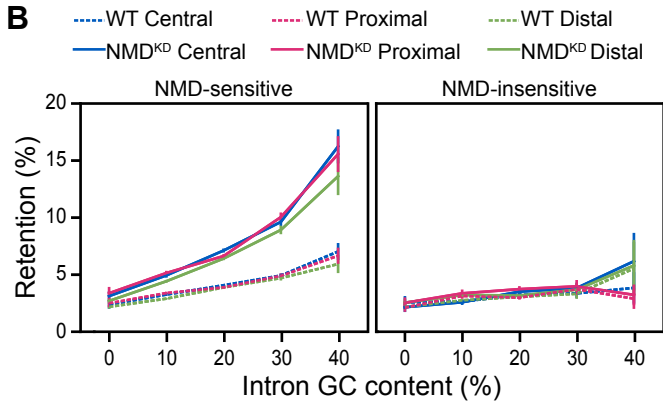
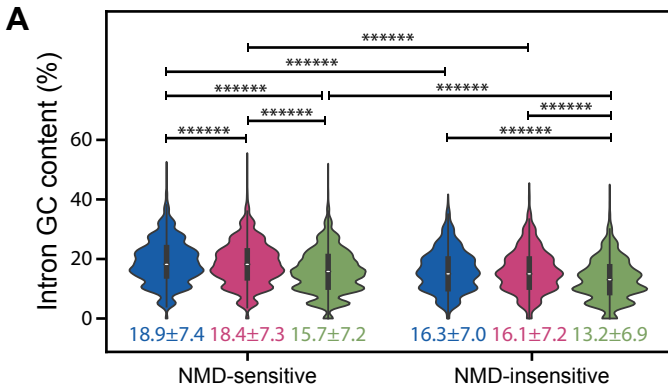
841 The full fitted model in explaining intron splicing efficiency, indicating whether each
842 parameter is positively or negatively correlated with splicing efficiency. The parameter
843 abbreviations are explained in Supplemental Table S1.

844









D

$$e^{SE} = 0.295 \times \log_{10}(EL) - 0.170 \times GC + 0.098 \times SA + 0.109 \times TC + 0.09 \times N + 0.054 \times D_{TSS} + 0.052 \times L_I - 0.037 \times NMD + 0.067 \times SD + 0.074 \times L_T - 0.03 \times GC_T - 0.023 \times L_{FE} - 0.024 \times TG_I + 0.017 \times 3N + 0.032 \times TG_T - 0.032 \times TC_T + 0.026 \times GC_{FE} + 0.013 \times GC_{PE} + 0.023 \times TG_{FE} - 0.019 \times TC_{FE} + 0.011 \times D_N - 0.021 \times \Delta G_{p50} + 0.01 \times M_I + 0.018 \times \Delta G_{p50-I} + 0.022 \times \Delta G_{p150-I} - 0.023 \times TG_{PE} + 0.009 \times M_{FE}$$