



## Genomic origin, fragmentomics, and transcriptional correlation of long cell-free DNA molecules in human plasma

Huiwen Che, Peiyong Jiang, Lois L Y Choy, et al.

*Genome Res.* published online February 26, 2024

Access the most recent version at doi:[10.1101/gr.278556.123](https://doi.org/10.1101/gr.278556.123)

---

<b>P&lt;P</b>	Published online February 26, 2024 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

---

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

---

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

---

Published by Cold Spring Harbor Laboratory Press

1 **Genomic origin, fragmentomics and transcriptional correlation of long cell-free DNA**  
2 **molecules in human plasma**

3  
4 **Running title:** Origin and characteristics of long cfDNA

5  
6 Huiwen Che<sup>1,2,3</sup>, Peiyong Jiang<sup>1,2,3,4</sup>, L.Y. Lois Choy<sup>1,2,3,4</sup>, Suk Hang Cheng<sup>1,2,3</sup>, Wenlei  
7 Peng<sup>1,2,3</sup>, Rebecca W.Y. Chan<sup>1,2,3</sup>, Jing Liu<sup>1,2,3</sup>, Qing Zhou<sup>1,2,3</sup>, W.K. Jacky Lam<sup>1,2,3,4</sup>,  
8 Stephanie C.Y. Yu<sup>1,2,3</sup>, So Ling Lau<sup>5</sup>, Tak Y. Leung<sup>5</sup>, John Wong<sup>6</sup>, Vincent Wai-Sun Wong<sup>7</sup>,  
9 Grace L.H. Wong<sup>7</sup>, Stephen L. Chan<sup>4,8</sup>, K.C. Allen Chan<sup>1,2,3,4</sup>, Y.M. Dennis Lo<sup>1,2,3,4,#</sup>

10

11 **AFFILIATIONS**

12 <sup>1</sup>Centre for Novostics, Hong Kong Science Park, Pak Shek Kok, Hong Kong SAR, China.

13 <sup>2</sup>Li Ka Shing Institute of Health Sciences, The Chinese University of Hong Kong, Shatin,  
14 Hong Kong SAR, China.

15 <sup>3</sup>Department of Chemical Pathology, Prince of Wales Hospital, The Chinese University of  
16 Hong Kong, Shatin, Hong Kong SAR, China.

17 <sup>4</sup>State Key Laboratory of Translational Oncology, The Chinese University of Hong Kong,  
18 Prince of Wales Hospital, Shatin, Hong Kong SAR, China.

19 <sup>5</sup>Department of Obstetrics and Gynaecology, Prince of Wales Hospital, The Chinese  
20 University of Hong Kong, Shatin, Hong Kong SAR, China.

21 <sup>6</sup>Department of Surgery, The Chinese University of Hong Kong, Prince of Wales Hospital,  
22 Shatin, Hong Kong SAR, China.

23 <sup>7</sup>Department of Medicine and Therapeutics, The Chinese University of Hong Kong, Prince of  
24 Wales Hospital, Shatin, Hong Kong SAR, China.

25 <sup>8</sup>Department of Clinical Oncology, Sir Y.K. Pao Centre for Cancer, The Chinese University of  
26 Hong Kong, Prince of Wales Hospital, Shatin, Hong Kong SAR, China.

27

28

29 **CORRESPONDENCE**

30 #Correspondence to Prof. Y.M. Dennis Lo

31 Address: Centre for Novostics, 1/F, 18 Science Park East Ave., Hong Kong Science Park,

32 Pak Shek Kok, New Territories, Hong Kong SAR, China.

33 Email: [loym@cuhk.edu.hk](mailto:loym@cuhk.edu.hk)

34 **ABSTRACT**

35 Recent studies have revealed an unexplored population of long cell-free DNA (cfDNA)  
36 molecules in human plasma using long-read sequencing technologies. However, the  
37 biological properties of long cfDNA molecules (> 500 bp) remain largely unknown. To this  
38 end, we investigated the origins of long cfDNA molecules from different genomic elements.  
39 Analysis of plasma cfDNA using long-read sequencing revealed uneven distribution of long  
40 molecules from across the genome. Long cfDNA molecules showed overrepresentation in  
41 euchromatic regions of the genome, in sharp contrast to short DNA molecules. We observed  
42 a stronger relationship between the abundance of long molecules and mRNA gene  
43 expression levels, compared with short molecules (Pearson's  $r = 0.71$  versus  $-0.14$ ).  
44 Moreover, long and short molecules demonstrated distinct fragmentation patterns  
45 surrounding CpG sites. Leveraging the cleavage preferences surrounding CpG sites, the  
46 combined cleavage ratios of long and short molecules could differentiate patients with  
47 hepatocellular carcinoma (HCC) from non-HCC subjects (AUC = 0.87). We further  
48 investigated knockout mice in which selected nuclease genes had been inactivated, in  
49 comparison with wild-type mice. The proportion of long molecules originating from  
50 transcription start sites were lower in *Dffb*-deficient mice but higher in *Dnase1/3*-deficient  
51 mice, compared to that of wild-type mice. This work thus provides new insights into the  
52 biological properties and potential clinical applications of long cfDNA molecules.

## 53 INTRODUCTION

54 Fragmentation patterns of cell-free DNA (cfDNA) are demonstrated to be linked to a myriad  
55 of biological characteristics, including nuclease activities (Serpas et al. 2019; Han et al.  
56 2020), DNA methylation (Zhou et al. 2022), nucleosome structures (Ivanov et al. 2015;  
57 Snyder et al. 2016; Sun et al. 2019), mRNA expression levels (Ulz et al. 2016), and DNA-  
58 binding transcription factor activity (Ulz et al. 2019). These characteristics have spurred  
59 much research efforts in understanding the underlying biological mechanisms. For example,  
60 size analysis of cfDNA molecules reflects characteristic patterns of cfDNA fragmentation.  
61 Through short-read sequencing (Illumina), a typical plasma cfDNA fragment size distribution  
62 could be depicted with a dominant peak around 166 base pairs (bp), with 10-bp periodicities  
63 below 150 bp, suggesting a nucleosomal origin (Lo et al. 2010; Ivanov et al. 2015). Thus,  
64 cfDNA is generally believed to be a population of short DNA molecules. Moreover, cfDNA  
65 fragments from different tissues may bear information tracing to their tissues of origin. As  
66 examples, fetal cfDNA molecules in pregnant women and tumoral cfDNA molecules in  
67 patients with cancer are generally shorter than background hematopoietic cell-derived  
68 cfDNA (Chan et al. 2004; Lo et al. 2010; Jiang et al. 2015). These findings have catalyzed  
69 the use of size information of cfDNA for disease detection (Yu et al. 2014; Jiang et al. 2015;  
70 Mouliere et al. 2018).

71

72 Single-molecule sequencing technologies, including Pacific Biosciences single-molecule  
73 real-time (SMRT) sequencing and Oxford Nanopore Technologies (ONT) sequencing, have  
74 opened up the possibility of detecting and characterizing long cfDNA molecules. Recent  
75 studies using these platforms have uncovered a population of long cfDNA molecules up to  
76 tens of kilobases, in the plasma DNA of healthy, pregnant individuals and patients with  
77 cancer (Yu et al. 2021; Choy et al. 2022; Yu et al. 2023b, 2023a). Depending on the  
78 sequencing platforms used, long molecules account for about a median of 15% and 5% of  
79 the total molecules from SMRT and ONT platforms, respectively (Yu et al. 2023b). Hence,

80 single-molecule, long-read sequencing technologies appear to be capable of sequencing a  
81 much wider spectrum of cfDNA molecules.

82

83 Studies have investigated the quantity, size distribution, methylation patterns and end motifs  
84 of cfDNA molecules generated by long-read sequencing (Yu et al. 2021; Choy et al. 2022;  
85 Katsman et al. 2022; Lau et al. 2023). As single-molecule sequencing technologies are able  
86 to measure methylation patterns (Tse et al. 2021; Lau et al. 2023), one can use methylation  
87 patterns to explore the tissues of origin of long cfDNA molecules. Additional biological  
88 properties of long molecules, however, remain largely unexplored. Hence, in this study, we  
89 investigated whether long and short cfDNA molecules might originate from different genomic  
90 elements, including euchromatin and heterochromatin. We aggregated cfDNA data  
91 generated by SMRT or ONT sequencing to profile the genomic representations and studied  
92 the correlation between long cfDNA molecules and gene transcriptional activity. We  
93 examined end frequencies of long molecules derived from regulatory regions such as DNase  
94 I hypersensitive sites (DHSs) and CCCTC-binding factors (CTCF). Moreover, we analyzed  
95 cleavage profiles surrounding cytosine-phosphate-guanine sites (CpGs) to investigate long  
96 cfDNA fragmentomics in the context of cancer detection. We generated additional cfDNA  
97 data from various nuclease-knockout mice using SMRT sequencing to gain biological  
98 insights regarding the role of nucleases in the generation of long cfDNA molecules.

99 **RESULTS**100 **Long cfDNA molecules originate unevenly across the human genome**

101 A preponderance of molecules with 5' end-nucleotide of base adenine (A), namely A-end,  
102 has been observed in cfDNA molecules longer than 500 bp (Yu et al. 2021). We therefore  
103 defined long cfDNA molecules as those of > 500 bp. We investigated whether long cfDNA  
104 molecules originated evenly across the genome. By pooling sequenced data from previous  
105 studies (Yu et al. 2021; Choy et al. 2022; Yu et al. 2023b) using the SMRT and ONT  
106 platforms (**Supplemental Table S1**), respectively, we first compared the genomic  
107 representation profiles of long (>500 bp) and short (<= 500bp) molecules, across 100-  
108 kilobase (kb) non-overlapping bins (Methods). Both genomic origins of long and short  
109 molecules displayed unevenly along the whole genome. We observed differential genomic  
110 representations between long and short molecules in some regions, with Chromosome 10  
111 representation shown as an example (**Fig. 1; Supplemental Fig. S1**). The differential  
112 signals in genomic representation were found in both SMRT (**Fig. 1A**) and ONT (**Fig. 1B**)  
113 sequencing data, with platform-specific differences. Specifically, we observed that long  
114 molecules tended to be overrepresented in regions that overlapped with light bands on  
115 Giemsa-stained chromosomes and underrepresented in dark-banded regions shown in  
116 ideograms (**Fig. 1**). The light and dark bands typically correspond to GC-rich euchromatic  
117 and AT-rich heterochromatic regions, respectively (Bickmore and Sumner 1989; BAC  
118 Resource Consortium et al. 2001). Thus, long molecules were preferentially derived from  
119 euchromatic regions compared to short molecules. This differential representation was  
120 particularly prominent on Chromosome 19 (**Supplemental Fig. S2**), which is well-known for  
121 high GC-content and considered as the most gene-rich chromosome (Grimwood et al. 2004).

122

123 As euchromatin generally tends to be gene-rich, we wondered if the gene density of a  
124 genomic region might positively correlate with the difference between long and short cfDNA  
125 molecules from that region. Indeed, the difference was significantly correlated to autosomal  
126 gene density (Pearson's  $r = 0.56$ ,  $P < 2.2 \times 10^{-16}$ ) and the gene-rich Chromosome 19

127 showed a correlation of 0.7 using non-pregnant controls from SMRT sequencing data. The  
128 correlation to Chromosome 19 gene density was 0.76 ( $P < 2.2 \times 10^{-16}$ ) for pregnant samples  
129 using ONT sequencing data. We further examined the potential relationship between  
130 genomic representation of cfDNA molecules and presence of DNA double-strand breaks.  
131 Previous studies have found that transcribed genes are hotspots for endogenous double-  
132 strand breaks (Crosetto et al. 2013; Lensing et al. 2016; Ballarino et al. 2022). A weak  
133 positive correlation (Pearson's  $r = 0.43$ ,  $P < 2.2 \times 10^{-16}$ ; **Supplemental Fig. S3**) between  
134 overrepresentation of long molecules and double-strand breaks detected from a  
135 lymphoblastoid cell line sample was observed (Methods).

136

### 137 **Association between long molecules and gene expression in human samples**

138 Motivated by the observation that long molecules appeared to be enriched in euchromatin  
139 that is associated with active transcription, we assessed the potential relationship between  
140 abundance of such molecules and transcriptional activity. First, we analyzed the abundance  
141 of cfDNA molecules originating from unexpressed and housekeeping genes (Methods). Long  
142 cfDNA molecules showed higher abundance over housekeeping genes and lower  
143 abundance over unexpressed genes in non-pregnant controls and pregnant subjects (**Fig.**  
144 **2A, B**). We further examined the abundance of cfDNA molecules originating from gene  
145 bodies of differentially expressed genes. Autosomal protein-coding genes were ranked  
146 based on their median expression across tissues and grouped into five sets, from a lowly  
147 expressed EXP1 to a highly expressed EXP5 set. In line with the enrichment in  
148 housekeeping genes, the analysis of non-pregnant controls revealed a stepwise increase of  
149 long cfDNA molecules abundance as expression levels increased from EXP1 to EXP5 ( $P =$   
150  $0.01374$ , Mann-Kendall trend test; Pearson's  $r = 0.78$ ), while no such trend was found in  
151 short molecules ( $P = 0.4032$ , Mann-Kendall trend test; Pearson's  $r = -0.28$ ; **Supplemental**  
152 **Fig. S4A**). We confirmed the observation that long molecule abundance was positively  
153 associated with gene expression from SMRT in pregnancies (**Supplemental Fig. S4B**) and

154 ONT sequencing data (**Supplemental Fig. S4C, D**). Analyzing the relationship at a higher  
155 resolution, we showed that the median abundances of long molecules were positively  
156 correlated (Pearson's  $r = 0.71$ ,  $P < 2.2 \times 10^{-16}$ ) with median transformed expression scores  
157 across tissues based on 15,556 expressed genes. In contrast, such positive correlation was  
158 not observed in short molecules (Pearson's  $r = 0.14$ ,  $P = 0.003$ ; **Fig. 2C**). The conclusion  
159 was further confirmed by the ONT data (**Fig. 2D**). Moreover, investigating overall methylation  
160 levels revealed lower methylation in long molecules compared with short molecules (**Fig. 2E,**  
161 **F**).

162

163 As the abundance of long cfDNA molecules in human plasma demonstrated a relationship  
164 with gene expression, we explored whether such signature could be used to distinguish  
165 patients with hepatocellular carcinoma (HCC) from subjects without HCC. Thus, we collected  
166 SMRT sequencing data using plasma samples from 20 healthy individuals, 19 hepatitis B  
167 virus (HBV) carriers and 48 patients with HCC (**Supplemental Table S2**). Median values of  
168 1,021,412, 286,492, and 712,223 high-quality circular consensus sequencing reads were  
169 obtained for healthy individuals, HBV carriers and patients with HCC, respectively. Long  
170 cfDNA molecules accounted for a median of 21.3%, 19.7% and 23.7% of total molecules in  
171 healthy subjects, HBV carriers and patients with HCC, respectively. As the number of  
172 molecules was relatively low, we attempted to examine a group of genes. The top 5000  
173 expressed genes in HCC tumor from The Cancer Genome Atlas dataset (Methods) were  
174 identified. We compared the abundance of long and short molecules over these HCC-  
175 associated genes. Long molecules showed significantly higher ( $P = 5.814 \times 10^{-5}$ , Kruskal-  
176 Wallis test) abundance in patients with HCC compared to non-HCC groups, whereas short  
177 molecules did not exhibit significant differences ( $P = 0.06504$ , Kruskal-Wallis test) in this  
178 dataset (**Fig. 3A**). A receiver operating characteristic (ROC) curve based on the abundance  
179 of long molecules yielded an area under the curve (AUC) of 0.76 to distinguish 39 non-HCC  
180 individuals and 48 patients with HCC. When the number of total molecules increased to 1

181 million for each sample, the AUC improved to 0.9 to distinguish 12 non-HCC individuals and  
182 18 patients with HCC (**Fig. 3B**).

183

184 In addition to the abundance of molecules, we assessed molecule ends in specific regions.  
185 End frequencies of cfDNA at transcription start sites (TSSs) and regulatory regions that  
186 include DHSs and CTCF binding sites (Methods) were examined. As the expression level  
187 rose, the normalized end frequencies of long molecules gradually increased in the vicinity of  
188 TSSs (e.g. 50-bp upstream and downstream of TSSs), rising from 1.01-fold in EXP1 to 1.7-  
189 fold in EXP5. Conversely, for short molecules, the normalized end frequencies had a  
190 decreasing tendency (**Fig. 4A; Supplemental Fig. S4E**). Long molecules pooled from the  
191 non-pregnant controls displayed higher end frequencies near the peak of tissue-invariant  
192 DHSs (**Fig. 4B**). We observed that the end frequencies of long cfDNA molecules were two  
193 times that of short molecules in DHSs (Methods). Analyzing end frequencies near tissue-  
194 specific DHSs revealed effects of tissue types on long cfDNA fragmentations. Specifically, at  
195 liver-specific and myeloid-specific DHSs, long cfDNA in non-pregnant controls tended to  
196 have stronger end preferences (1.4-fold) comparing to short cfDNA (**Supplemental Fig.**  
197 **S4F**). As a confirmation, the same pattern around tissue-invariant DHSs was observed in  
198 plasma of pregnant women from ONT sequencing (**Supplemental Fig. S4G**). At CTCF  
199 binding sites, long molecules exhibited higher cfDNA end frequencies as well, in both SMRT  
200 and ONT sequencing data (**Fig. 4C; Supplemental Fig. S4H**). These results have further  
201 highlighted that long molecules were preferentially derived from open chromatin regions and  
202 their abundance was correlated with transcriptional activities.

203

#### 204 **Differential cleavage profiles between long and short molecules surrounding CpGs**

205 Studies have demonstrated that the cfDNA cleavage preferences occurred at methylated  
206 CpG sites (Han et al. 2021; Zhou et al. 2022). Beyond the fragmentation around regulatory  
207 regions, we wondered whether fragmentation of long molecules at CpGs would be different  
208 from those of short molecules. We analyzed the cleavage patterns at commonly methylated

209 and unmethylated CpGs in pooled healthy data (Methods). The cleavage of short molecules  
210 recapitulated earlier findings that nuclease-mediated preferential cleavage occurs at  
211 methylated CpGs in healthy subjects analyzed by Illumina sequencing (Zhou et al. 2022).  
212 Specifically, the cfDNA cleavage was relatively enriched at the cytosine of methylated CpGs  
213 (position 0), followed by a rapid decrease at the nucleotide 1 position immediately preceding  
214 methylated CpGs. We found that the cleavage of long molecules did not show such a pattern  
215 at methylated cytosines but exhibited the preferred cutting at positions several nucleotides  
216 away from the methylated CpGs (e.g. positions -4, -2 and 4) (**Supplemental Fig. S5A**).

217

218 To explore whether the cleavage profile around CpGs for long molecules would be  
219 associated with diseases, we evaluated cleavage patterns using the HCC cohort data  
220 mentioned above. Considering the relatively low number of molecules, we evaluated  
221 cleavage patterns related to all autosomal CpGs. As the human genome is highly methylated,  
222 with 70%-80% of CpGs being methylated (Ziller et al. 2013), the cleavage profiles of short  
223 molecules showed a cutting preference at position 0 cytosine (**Fig. 5A**). Due to the low  
224 number of sequenced molecules, long molecules showed higher variability in cleavage  
225 profiles. The cleavages of long molecules, however, consistently showed significant higher  
226 proportions ( $P = 2.6 \times 10^{-6}$ , paired Wilcoxon test,  $P$  values adjusted by Benjamini-Hochberg  
227 method) than the adjacent downstream nucleotide at positions -4, -2, 1 and 4 in healthy  
228 subjects and HBV carriers (**Fig. 5B; Supplemental Fig. S5B**). The ratio of 5' CGN to NCG  
229 end motifs (CGN/NCG motif ratio) has been shown to inform methylation level and can be  
230 used to distinguish patients with HCC from non-HCC individuals using Illumina sequencing  
231 (Zhou et al. 2022). In this SMRT sequencing dataset, patients with HCC showed overall  
232 lower CGN/NCG motif ratios ( $P = 2.081 \times 10^{-7}$ , Kruskal-Wallis test) for short molecules. On  
233 the contrary, for long molecules, CGN/NCG motif ratios of patients with HCC showed no  
234 significant difference ( $P = 0.05137$ , Kruskal-Wallis test) among the three groups (**Fig. 5C**).

235

236 We further attempted to use the ratio between the molecules ending at positions -4, -2, 1  
237 and 4 and those ending at position -1. Using this cleavage ratio, both long and short  
238 molecules displayed significant differences ( $P = 3.893 \times 10^{-8}$  and  $P = 8.036 \times 10^{-8}$ , Kruskal-  
239 Wallis test) between patients with HCC and non-HCC subjects (**Fig. 5D**). As a result, long  
240 molecules alone achieved an AUC of 0.85 to distinguish 39 non-HCC individuals and 48  
241 patients with HCC. Combining the ratio of long and short molecules, the AUC increased to  
242 0.87 (**Fig. 5E**).

243

#### 244 **Investigation of long molecules generation using knock-out mice**

245 The amount of cfDNA molecules with 5' ending in an A or G was reported to increase as the  
246 size of cfDNA molecules increased (Yu et al. 2021) and the nuclease DNA fragmentation  
247 factor subunit beta (DFFB) exhibited a preference to cut 5' to an A or G nucleotide (Han et al.  
248 2020). We hypothesized that long molecules were predominantly generated by DFFB. To  
249 gain potential mechanistic insights on generating long cfDNA molecules in human plasma,  
250 we studied mouse models in which different nuclease genes had been knocked out. Mice of  
251 the C57BL/6 background were used. We sequenced 4 wild-type (WT), 5 *Dffb*<sup>-/-</sup>, 5 *Dnase1*<sup>-/-</sup>,  
252 5 *Dnase1/3*<sup>-/-</sup>, and 5 *Dnase1*<sup>-/-</sup> and *Dnase1/3*<sup>-/-</sup> double knockout mice samples, with each  
253 sample pooled from 3-5 mice subjects, using SMRT sequencing (**Supplemental Table S3**).  
254 We first analyzed the size profile of plasma cfDNA from these mice by pooling sequenced  
255 data of the same genotype. Compared with wild-type mice, the frequencies of long cfDNA  
256 molecules above 1 kb in *Dffb*<sup>-/-</sup> mice were consistently lower, whereas the corresponding  
257 frequencies in  
258 *Dnase1*<sup>-/-</sup> mice were higher, and such an increasing pattern was further enhanced in  
259 *Dnase1/3*<sup>-/-</sup> mice and mice with double deletion of *Dnase1* and *Dnase1/3* (*Dnase1*<sup>-/-</sup> &  
260 *Dnase1/3*<sup>-/-</sup>) (**Fig. 6A**). For instance, the percentages of cfDNA molecules with a size of >  
261 500 bp were 21.7%, 16.8%, 29.3%, 37.2%, and 45% in wild-type, *Dffb*<sup>-/-</sup>, *Dnase1*<sup>-/-</sup>,  
262 *Dnase1/3*<sup>-/-</sup>, and *Dnase1*<sup>-/-</sup> & *Dnase1/3*<sup>-/-</sup> mice, respectively. The percentages of long  
263 molecules were significantly different ( $P = 0.002053$ , Kruskal-Wallis test) between different

264 types of mice, with *Dnase1l3*<sup>-/-</sup>, and *Dnase1*<sup>-/-</sup> & *Dnase1l3*<sup>-/-</sup> mice higher than wild-type mice  
265 (**Supplemental Fig. S6A**). Meanwhile, the frequencies of short molecules with a size < 250  
266 bp decreased in *Dnase1l3*<sup>-/-</sup> and *Dnase1*<sup>-/-</sup> & *Dnase1l3*<sup>-/-</sup> mice, with a relatively more  
267 pronounced reduction in mononucleosomal size (**Fig. 6B**). The data suggested that DFFB  
268 was at least in part responsible for the generation of long cfDNA molecules. The removal of  
269 DNASE1L3 and DNASE1 would lead to an enhanced effect in generating long cfDNA  
270 molecules, as DFFB might, in the absence of DNASE1L3 and DNASE1, act as a dominant  
271 nuclease in contributing cfDNA into circulation. The deletion of *Dffb* resulted in the reduction  
272 of long molecules originating from regions around TSSs which normally represent open  
273 chromatin states (**Fig. 6C**), further providing evidence that DFFB might be responsible for  
274 the patterns related to long cfDNA molecules. The increase of long cfDNA molecules in  
275 regions nearby TSSs was elevated in mice where DFFB functioned prominently by knocking  
276 out *Dnase1* and *Dnase1l3*. These patterns could be reproduced in open chromatin regions  
277 comprising DHSs and CTCF binding sites (**Supplemental Fig. S6B, C**). Moreover, the  
278 highly expressed genes tended to harbor more long molecules in *Dffb*-competent mice.  
279 However, such enrichment of long molecules was absent in *Dffb*-deficient mice (**Fig. 6D**).  
280

281 **DISCUSSION**

282 This study explores previously unknown properties of long cfDNA molecules in human  
283 plasma. First, long cfDNA molecules have been found to be preferentially originated from  
284 transcriptionally active regions of the genome, with overall lower methylation and the end  
285 frequencies of long cfDNA molecules being more enriched in transcriptional start sites and  
286 open chromatin regions. Second, long molecules exhibit a distinctive cleavage profile  
287 surrounding CpG sites. The cleavage ratio that makes use of multiple positions surrounding  
288 CpG sites can be used to distinguish non-HCC individuals from patients with HCC. In  
289 contrast to the previous report that the generation of short cfDNA molecules might be largely  
290 attributed to DNASE1L3 (Serpas et al. 2019; Chan et al. 2020), we demonstrate that the  
291 characteristic distribution of long cfDNA molecules in the genome might at least in part be  
292 attributed to the DFFB enzyme activity using genetically modified mice. This observation  
293 suggests that different nucleases involved in the fragmentation of long and short cfDNA  
294 molecules in human plasma.

295

296 By using knockout mice, a plasma DNA fragmentation model based on the end motifs of  
297 cfDNA molecules has been proposed, in which DFFB, DNASE1L3 and possibly other  
298 nucleases first digest DNA intracellularly, and then extracellular DNASE1L3 and DNASE1  
299 further degrade the plasma DNA (Han et al. 2020). DNASE1L3, DNASE1 and DFFB  
300 preferentially cut at C-end, T-end and A-end nucleotides of cfDNA fragments, respectively  
301 (Han et al. 2020). The long molecules uncovered via single-molecule sequencing platforms  
302 are likely to be, at least in part, the end product of DFFB other than DNASE1L3, which is  
303 supported by the observation of enrichment of A-end fragments in long cfDNA molecules (Yu  
304 et al. 2021). Of note, following knocking out of the *Dnase1l3* and/or *Dnase1* genes in mice,  
305 DFFB appears to take on the leading role in cfDNA fragmentation. Additionally, no clear  
306 correlation was found between total cfDNA concentration and proportion of long molecules.  
307 The higher proportions of long molecules originated from accessible chromatin regions in  
308 *Dnase1l3*<sup>-/-</sup> and *Dnase1*<sup>-/-</sup> & *Dnase1l3*<sup>-/-</sup> double knockout mice were reminiscent of the

309 characteristic of long molecules in human plasma. The long cfDNA molecules in human  
310 plasma are therefore likely to be the digestion products of DFFB, which have escaped  
311 further or secondary digestion by DNASE1L3. The exact biological mechanisms for the  
312 escape of cleavage for long cfDNA remain elusive. *Dnase1l3*<sup>-/-</sup> mice have been shown to  
313 develop features of systemic lupus erythematosus on both the 129 and B6 genetic  
314 backgrounds (Sisirak et al. 2016). DNASE1L3 deficiencies in human are associated with  
315 systemic lupus erythematosus (Al-Mayouf et al. 2011). Long cfDNA that are generated by  
316 DFFB might have potential immunomodulatory effects to the innate immune responses.

317

318 Prior studies systematically profiled plasma DNA end motifs using mice of different  
319 genotypes and found altered end motif frequencies in nuclease-deficient mice (Serpas et al.  
320 2019; Chan et al. 2020). By linking the end motifs and nuclease-deficient genotypes, a non-  
321 negative matrix factorization algorithm-based approach was developed to deconvolute  
322 contributions of nucleases. The work has demonstrated the feasibility of inferring nuclease  
323 activities in liquid biopsy and that DNASE1L3 is a leading contributor to fragmentation in the  
324 plasma cfDNA (Zhou et al. 2023). One reason for previous studies focusing on the  
325 DNASE1L3 activity is largely because that the Illumina sequencing platform is designed for  
326 short read sequencing, with an upper limit of sequenceable fragments being at  
327 approximately 600 bp. Long cfDNA molecules investigated in this study are more likely to  
328 reflect the DFFB activity. The genome-wide representational analysis of long cfDNA  
329 molecules provides another means to dissect the DFFB activity and link the enzyme activity  
330 to gene transcriptional activity. Further investigations to elucidate cfDNA fragmentation are  
331 warranted. A more comprehensive picture of various nuclease activities could be uncovered  
332 when examining the full spectrum of cfDNA molecules in different nuclease-knockout mouse  
333 models.

334

335 We have uncovered genomic representation difference between long and short cfDNA  
336 molecules in human plasma. The enrichment of long molecules in transcriptionally active

337 regions may have important implications. It has been suggested that DFFB first cleaves at  
338 nuclear scaffold attachment sites (Nagata et al. 2003) to unfold the chromatin. DFFB induces  
339 DNA breaks as a signal for cell fate (Larsen and Sørensen 2017). Though direct evidence of  
340 locations of DFFB-induced DNA breaks is lacking, we observed a weak positive correlation  
341 between overrepresentation of long cfDNA molecules and DNA double-strand breaks  
342 detected in a lymphoblastoid cell line sample. Future examinations at higher resolutions to  
343 address the relationship between transcription-associated DNA double-strand breaks and  
344 long DNA generation will be needed. Whether DFFB and other nucleases preferentially  
345 cleave at those damaged DNA breakpoints and thereafter induce downstream DNA  
346 fragmentation remain to be investigated. Additionally, the nuclear morphology at the time of  
347 cell death might be another dimension that affects the fragmentation and hence results in  
348 overrepresentation of long molecules in euchromatin (Toné et al. 2007). Beyond the  
349 mechanism, the enrichment of long molecules in genic regions highlights potential diagnostic  
350 value. Retaining characteristics of chromatin structure, non-random fragmentation of cfDNA  
351 reflects transcriptomic and epigenomic status of dying cells (Ivanov et al. 2015; Ulz et al.  
352 2016; Snyder et al. 2016; Lo et al. 2021). A number of studies showed that coverage profiles  
353 near TSSs correlate with gene expression and that accessibilities near the transcription  
354 factor binding sites reflect the transcription factor activity (Ulz et al. 2016, 2019). Notably, the  
355 normalized end frequencies at TSSs in long cfDNA molecules were higher than those of  
356 short cfDNA molecules, suggesting a potential higher utility for long cfDNA molecules for  
357 inferring gene expression when compared with their shorter cfDNA counterparts. Moreover,  
358 the abundance of long molecules demonstrated discriminating power in cancer detection.

359

360 Compared with sequencing-by-synthesis, the current throughput of single-molecule  
361 sequencing is relatively low, which presents challenges to analyses that require a larger  
362 number of molecules. The low number of long cfDNA molecules sequenced creates  
363 obstacles for analyzing the association with gene expression at a higher resolution,  
364 especially in the context of measuring expression of a single gene. Nevertheless, with

365 technological advances, such as with the recent launch of a higher throughput version of the  
366 SMRT-based sequencer, and further optimized protocols, the throughput is likely to improve,  
367 potentially facilitating the elucidation of further biological insights and the development of  
368 clinical tools. Alternatively, enriching the long molecules of interest may help to adequately  
369 harness the advantages for diagnostic purposes. The performance of cancer detection using  
370 long molecule abundance and fragmentation patterns can be further enhanced. In addition,  
371 a larger cohort that allows validation of these biomarkers is desired. Platform-specific  
372 differences were observed from SMRT and ONT sequencing data. Future investigations into  
373 the use of a specific platform to maximize the use of long cfDNA are warranted.

374

375 In summary, this work expands the current knowledge of fragmentation properties of long  
376 cfDNA molecules. The characteristic distribution of long cfDNA molecules along the human  
377 genome is correlated to transcriptional activities and disease states, potentially opening up  
378 possibilities for biomarker discovery and new clinical applications.

## 379 **METHODS**

### 380 **Human sample collection**

381 The study was approved by the Joint Chinese University of Hong Kong – Hospital Authority  
382 New Territories East Cluster Clinical Research Ethics Committee. PacBio SMRT sequencing  
383 data of plasma cfDNA, including healthy individuals (N=15), pregnancies of different  
384 trimesters (N=28), HBV carriers (N=13) and patients with HCC (N=13), were obtained from  
385 previous studies (Yu et al. 2021; Choy et al. 2022). Additionally, we collected blood samples  
386 from healthy individuals (N=5), HBV carriers (N=6) and patients with HCC (N=35) from the  
387 Prince of Wales Hospital, Hong Kong, with written informed consent. Nanopore sequencing  
388 data of plasma cfDNA from pregnancies of different trimesters (N=31) were obtained from  
389 the prior study (Yu et al. 2023b). Additionally, plasma cfDNA from non-pregnant healthy  
390 individuals (N=5) were collected and subjected to ONT sequencing. The blood samples were  
391 processed by a double centrifugation protocol (1,600 x g for 10 min at 4°C, followed by  
392 further centrifugation of the plasma at 16,000 x g for 10 min at 4°C) as previously described  
393 (Chiu et al. 2001).

394

### 395 **Murine Models**

396 Mice with a CRISPR-Cas9-targeted deletion of exon 5 in *Dnase1/3* of the C57BL/6 (B6)  
397 background were generated by The Jackson Laboratory. Mice carrying a targeted allele  
398 mutation of *Dnase1* (*Dnase1<sup>tm1.1(KOMP)Vlcg</sup>*) and mice carrying a targeted allele of *Dffb*  
399 (*Dffb<sup>C57BL/6N-Dffbem1Wtsj</sup>*), both on the B6 background were obtained from the Knockout Mouse  
400 Project Repository of the University of California, Davis. These mice were obtained under a  
401 third-party distribution agreement and are nontransferable. *Dnase1/3<sup>-/-</sup>* mice were cross-bred  
402 with *Dnase1<sup>-/-</sup>* mice to produce *Dnase1<sup>-/-</sup>* & *Dnase1/3<sup>-/-</sup>* mice in the Laboratory Animal  
403 Services Centre of The Chinese University of Hong Kong (CUHK). Mice including wild-type  
404 B6 were maintained in the same facility. All experimental procedures were approved by the  
405 Animal Experimentation Ethics Committee of CUHK and performed in compliance with the

406 Guide for the Care and Use of Laboratory Animals (8<sup>th</sup> edition) established by the National  
407 Institutes of Health.

408

#### 409 **Murine sample collection**

410 Mice were sacrificed and exsanguinated by cardiac puncture. Blood was transferred into  
411 EDTA-containing collection tubes. The blood samples were processed by a double  
412 centrifugation protocol (1,600 x *g* for 10 min at 4°C, followed by further centrifugation of the  
413 plasma at 16,000 x *g* for 10 min at 4°C) as previously described (Chiu et al. 2001). The  
414 resulting plasma was collected.

415

#### 416 **Library preparation and sequencing**

417 Plasma DNA was extracted from 1.2 to 4 mL of human plasma and 1.1 to 1.4 mL of pooled  
418 mice plasma with the QIAamp Circulating Nucleic Acid Kit (Qiagen) according to the  
419 manufacturer's instructions. SMRTbell Express Template Prep Kit 2.0 (PacBio) was used for  
420 the library preparation of plasma DNA. Briefly, DNA molecules were ligated with hairpin  
421 adaptors to form a circularized template. Sequencing Primer v4 was annealed to the  
422 sequencing template, followed by binding of polymerase to templates using a Sequel II  
423 Binding Kit 2.1 and Internal Control Kit 1.0 (PacBio). SMRT Cell 8M was used for  
424 sequencing with a Sequel II Sequencing 2.0 Kit (PacBio), and sequencing movies were  
425 collected for 30 hours. Nanopore library preparation and sequencing of pregnant samples  
426 are as described before (Yu et al. 2023b). For healthy samples, nanopore libraries were  
427 prepared using Native Barcoding Kit (SQK-NBD114.96, ONT) according to the  
428 manufacturer's instructions except that a bead-to-sample ratio of 3 was used in all clean-up  
429 steps using AMPure XP beads with prolonged incubation time for DNA repair, end-prep,  
430 barcode and adapter ligation. Short Fragment Buffer, which retained DNA fragments of all  
431 sizes, was used in adapter ligation. Each library was loaded onto a PromethION Flow Cell  
432 R10.4.1 and sequenced on a PromethION device for 72 hours.

433

**434 Sequence alignment**

435 Circular consensus sequencing (CSS) reads that were generated from at least 3 PacBio  
436 subreads were kept. CCS reads from human and mice data were aligned to hg19 and mm10  
437 respectively using blasr (Chaisson and Tesler 2012). Nanopore-generated sequences were  
438 aligned to hg19 with minimap2 (Li 2018) as described in (Yu et al. 2023b). The size  
439 distributions of the raw sequenced molecules before and after alignment were examined.  
440 SMRT sequencing data showed highly overlapped size distributions, while ONT sequencing  
441 data showed deviations comparing before and after alignment (**Supplemental Fig. S7**).  
442 Chimeric reads were filtered out prior to downstream analyses. For human sequence data,  
443 re-analysis of the data using the GRCh38 human reference genome would not affect the  
444 results significantly, as the major difference between these two versions of the human  
445 reference genome is the sequence representation for centromeres and highly repetitive  
446 regions. Similarly, for mouse sequence data, re-analysis of the data using the GRCm39  
447 mouse reference genome would not affect the results significantly because only uniquely  
448 aligned reads were included in downstream analyses.

449

**450 Genomic representation analysis**

451 The human genome was partitioned into non-overlapping 100 kb bins. Sequences that  
452 mapped to blacklist regions (Amemiya et al. 2019) were removed. Long and short molecules  
453 were counted in each bin. The bin counts were smoothed by a 1 Mb moving average window  
454 and normalized by the median bin count of autosomes. Cytobands information for  
455 chromosome ideograms were obtained from UCSC Genome Browser  
456 (<http://genome.ucsc.edu/>). Gene densities were estimated by number of genes in a 100-kb  
457 window. The correlations between normalized bin counts and gene densities were calculated  
458 using Pearson's method and smoothed data applying 1Mb-window were used.

459

460 To examine the correlation between genomic representation of cfDNA molecules and  
461 presence of DNA double-strand breaks, we obtained double-strand breaks detected from a  
462 lymphoblastoid cell line sample (Bouwman et al. 2020). DNA double-strand breaks were  
463 counted in each 100-kb bin and normalized by the median bin count.

464

465 To calculate the proportions of long molecules originated from TSSs, DHSs and CTCF  
466 binding sites in mice, we counted long and short molecules overlapped with these regions.  
467 2000 bp upstream and downstream of transcription start sites or peak points of DHSs and  
468 CTCF binding sites were regarded as regions of interest. At least 50% size of a molecule  
469 needs to be overlapped with these regions of interest to be counted. The proportions of long  
470 molecules were computed as the number of long molecules divided by the total number of  
471 molecules overlapped with the regions of interest.

472

### 473 **Genic abundance analysis**

474 To quantify molecule abundance over genes, we adopted similar quantification of transcripts  
475 using reads per kilobase per million mapped reads. We counted molecules overlapped with  
476 gene bodies and normalized by the total number of molecules and the total length of the  
477 genes. For both long and short molecules, we required at least 50% size of a molecule to be  
478 overlapped with a gene to be counted. To calculate the Pearson's correlation between  
479 molecule abundance and gene expression, we used the median molecule abundance of  
480 samples and correlated with the median gene expression of each gene set. The Mann-  
481 Kendall trend test was performed using the median molecule abundance of samples.

482

483 HCC cancer-associated genes were identified using curated HCC expression data from The  
484 Cancer Genome Atlas (Cancer Genome Atlas Research Network et al. 2013). Expressions  
485 were averaged across 356 HCC tumor tissue samples and ranked by the averaged  
486 expression. Top 5000 expressed genes were selected and used as HCC cancer-associated

487 genes. The abundance of long and short molecules over these 5000 genes were calculated.  
488 AUCs were calculated based on the long molecule index.

489

490 Single cell gene expression data were downloaded from The Human Protein Atlas  
491 (<https://www.proteinatlas.org/>) and The Tabular Muris project (Tabula Muris Consortium et al.  
492 2018) for human and mouse, respectively. In human analysis, housekeeping genes were  
493 retrieved from a previous study (Eisenberg and Levanon 2013) and 3510 autosomal  
494 housekeeping genes were used. 2154 genes with a mean of less than 1 normalized  
495 transcript per million across 54 tissues were defined as unexpressed genes. Autosomal  
496 protein coding genes were stratified by median expression across 79 cell types. Gene  
497 groups of EXP1-5 correspond to genes with median normalized transcript per million in the  
498 range of [0, 0.001), [0.001, 1), [1, 8), [8, 30), [30, 11800), with each containing 4836, 3163,  
499 4187, 4816, and 3088 genes, respectively. To correlate the abundance of molecules with  
500 gene expression at relatively high resolution in human samples, median expression level of  
501 a gene across cell types were log-transformed and scaled to a score ranging from 0 to 1000.  
502 Median molecule abundance of genes at each scaled expression level was used. Genes  
503 with expression scores of lower and upper quantiles (outliers) were discarded from the  
504 correlation analysis.

505

506 Murine genes were stratified by median expression across 20 tissue types. The low, medium  
507 and high expression gene groups in mice correspond to autosomal protein coding genes  
508 with median normalized transcript count in the range of [0, 50), [50,180) and [180,450), with  
509 each containing 10987, 3883 and 5195 genes, respectively. The percentage of change in  
510 long molecule abundance between the knockout and wild-type mice was computed using the  
511 difference of long molecule abundance between knockout and wild-type mice, and divided  
512 by the abundance of long molecules in the wild-type mice.

513

514 **End frequency analysis**

515 End frequency analysis was performed surrounding the defined region of interests. The  
516 region of interests, including CpGs, DHSs, CTCF binding sites and TSSs, were extracted as  
517 follows.

518

519 Transcription start sites were retrieved using NCBI RefSeq data  
520 (<https://www.ncbi.nlm.nih.gov/refseq/>). DHSs and CTCF binding sites were identified from  
521 publicly available data. Tissue-invariant DHSs from the study (Meuleman et al. 2020) were  
522 identified and genomic coordinates were converted from hg38 to hg19 using the UCSC  
523 liftOver tool. 44,997 tissue-invariant DHSs were used in the analysis. 157,556 CTCF binding  
524 sites from the Encyclopedia of DNA Elements (ENCODE) uniformly processed transcription  
525 factor binding site clusters in human were used for analysis (Dunham et al. 2012). 107,227  
526 DHSs (ENCFF855RCO) and 62,461 CTCF binding sites (ENCFF883UPM) from CH12.LX  
527 mouse cell line available via ENCODE were used in mice analysis.

528

529 The peak or mid-point of the region of interests was regarded as position 0. The end  
530 frequencies within +/- 2000 bp were aggregated and normalized with median count of the  
531 region. For regulatory regions (DHS and CTCF) and TSS, we applied a 10 bp and 20 bp  
532 window to smooth the count. Fold change of end frequencies around DHSs was measured  
533 as follows. End frequencies within a DHS region (a median size of 211 bp) were aggregated.  
534 In parallel, end frequencies of the equivalent region that are 3000 bp away upstream and  
535 downstream were aggregated and used as a normalization factor. Fold change of end  
536 frequencies around TSSs was measured using a ratio between long and short in human.  
537 Normalized end frequencies from position -50 to 50 were aggregated to compute the ratio.

538

### 539 **Cleavage profiles surrounding CpGs**

540 For human data analysis, commonly methylated and unmethylated CpG sites were defined  
541 as methylation level > 70% and < 30% respectively in human buffy coat, liver and placental  
542 tissues. Whole genome bisulfite sequencing subjected to the Illumina platform from the three

543 tissues available from a previous study (Sun et al. 2015) were used to identify 2,422,965  
544 commonly methylated and 558,267 commonly unmethylated CpGs. For analysis of all CpGs,  
545 26,752,699 autosomal CpGs from the human genome hg19 were obtained. The cleavage  
546 proportion to measure the cutting frequency were measured as the number of molecules  
547 ends at a particular site divided by sequencing depth at the site, as described in (Zhou et al.  
548 2022). The CGN/NCG motif ratio was calculated using the number of molecules ends at  
549 position 0 divided by the number of molecules ends at position -1. The cleavage ratio was  
550 calculated using the aggregated number of molecules ends at position -4, -2, 1 and 4 divided  
551 by the number of molecules ends at position -1.

552

### 553 **Statistical analysis**

554 The Mann-Kendall trend test was used for testing monotonic trend in molecule abundance  
555 over gene expression groups. The alternative hypothesis that monotonic increasing trend is  
556 present was assumed. The Pearson's correlation test was used to measure relations  
557 between the abundance of molecules and gene expression. The Wilcoxon rank-sum test  
558 was used to compare two groups at a significance level of 0.05. The Kruskal-Wallis test or  
559 Friedman test in case of dependent groups was used to compare three or more groups at a  
560 significance level of 0.05. Post hoc pairwise Wilcoxon tests were performed with Benjamini-  
561 Hochberg adjustment to yield pairwise *P* values.

562

### 563 **DATA ACCESS**

564 The sequencing data generated in this work have been submitted to the European Genome-  
565 phenome Archive (EGA; <https://web2.ega-archive.org/>) under accession number  
566 EGAS00001005515.

567

### 568 **COMPETING INTEREST**

569 K.C.A.C. and Y.M.D.L. hold equities in DRA, Take2, Grail/Illumina, and Insighta. P.J. and  
570 W.K.J.L hold equities in Grail/Illumina. P.J. is a consultant to KingMed Future. W.P. is a  
571 consultant to Take2. K.C.A.C. and P.J. are Directors of DRA, Take2, KingMed Future and  
572 Insighta. W.K.J.L is a Director of DRA. S.C.Y.Y. received financial support from Oxford  
573 Nanopore for attending meetings. H.C., P.J., K.C.A.C. and Y.M.D.L. filed a US patent  
574 application (No. 63/544,014), entitled “Genomic origin, fragmentomics, and transcriptional  
575 correlation of long cell-free DNA” on October 13, 2023, based on the data in this study.  
576 Patent royalties are received from Grail, Illumina, DRA, Take2, and Xcelom.

577

#### 578 **ACKNOWLEDGEMENTS**

579 The work was supported by the Innovation and Technology Commission (InnoHK Initiative).  
580 Y.M.D.L. received an endowed chair from the Li Ka Shing Foundation. We would like to  
581 thank Ms. Angel Lai, Mr. Chris Kum, Mr. Saravanan Ramakrishnan, Mr. Xingfu Qin and Mr.  
582 Danny Wong for their technical assistance.

583 **FIGURE LEGENDS**

584

585 **Figure 1. Distribution of long and short molecules from human plasma DNA.** (A)  
586 Comparison of genomic representation on Chromosome 10 between long and short DNA  
587 molecules in 15 non-pregnant controls using SMRT sequencing. Overrepresentation and  
588 underrepresentation of long cfDNA molecules with respect to short molecules are indicated  
589 in blue and red, respectively. The genomic representation was determined based on 100-kb  
590 bins, and was further smoothed by a 1-Mb moving average sliding window. The horizontal  
591 solid lines indicate normalized median differences between long and short molecules. The  
592 dashed rectangular boxes indicate one euchromatic (i) and one heterochromatic (ii) region.  
593 The track in between overrepresentation and underrepresentation of long molecules shows  
594 the chromosome ideogram. The ideogram band colors correspond to cytogenetic bands in  
595 UCSC Genome Browser. Darker bands are AT-rich and lighter bands are GC-rich.  
596 Centromeric regions are indicated in dark green. The bottom track displays gene densities  
597 estimated by number of genes in 100-kb windows. (B) Comparison of genomic  
598 representation on Chromosome 10 between long and short DNA molecules in 31 pregnant  
599 samples using ONT sequencing.

600

601 **Figure 2. The abundance of long cfDNA molecules exhibits a positive correlation with**  
602 **transcriptional activity.** (A, B) The abundance of long and short molecules on gene bodies  
603 of unexpressed and housekeeping genes for non-pregnant controls and pregnant subjects.  
604 (A) Data from SMRT sequencing. (B) Data from ONT sequencing. (C, D) The correlation  
605 between gene expression and molecules abundance. Median expression level of a gene  
606 across tissues was log-transformed and scaled. The median molecule abundance was  
607 derived from scaled expression levels. (C) Pooled data of 15 non-pregnant controls from  
608 SMRT sequencing. (D) Pooled data of 31 pregnant samples from ONT sequencing. *P* values  
609 by Pearson's correlation test (C, D). (E, F) Comparisons of DNA methylation between long  
610 and short molecules. Data points from one sample are connected with a black line. (E) Data

611 from SMRT sequencing. (F) Data from ONT sequencing. *P* values by Wilcoxon rank-sum  
612 test (E, F).

613

614 **Figure 3. The abundance of long cfDNA molecules for HCC detection.** (A) Comparison  
615 of the abundance of SMRT sequencing molecules among healthy individuals, HBV carriers  
616 and patients with HCC. The abundance of long and short molecules was measured using  
617 the top 5000 expressed genes in HCC tumor tissues. The Kruskal-Wallis test *P* value for  
618 differences among groups. Post hoc pairwise Wilcoxon rank sum test *P* values with  
619 Benjamini-Hochberg adjustment are shown above horizontal lines. (B) ROCs of long  
620 molecule abundance measured in (A) for distinguishing individuals without HCC, including  
621 healthy subjects and HBV carriers, and with HCC. Multiple thresholds, including 0.1 million  
622 (Total > 0.1M), 0.3 million (Total > 0.3M), 0.5 million (Total > 0.5M) and 1 million (Total > 1M)  
623 molecules from a sample, were used to include samples for constructing ROCs.

624

625 **Figure 4. Normalized end frequencies of SMRT sequencing data.** (A) Normalized end  
626 frequencies of SMRT long and short molecules pooled from plasma of healthy individuals at  
627 TSSs of expression-stratified gene groups EXP1 to EXP5, corresponding to low to high  
628 expression. Transcription start positions are denoted as position 0. All transcription start  
629 sites were strand-adjusted so that positive positions are in the direction of transcription. (B, C)  
630 Normalized end frequencies of SMRT long and short molecules pooled from plasma of  
631 healthy individuals at DHSs (B) and CTCF binding sites (C). DHSs or CTCF binding sites  
632 peaks are denoted as position 0; downstream and upstream 2000 bp are shown.

633

634 **Figure 5. Cleavage profiles of long and short molecules surrounding CpGs.** (A, B)  
635 Cleavage profiles surrounding all autosomal CpGs for short (A) and long (B) cfDNA  
636 molecules. Each line represents one sample. A cleavage window of 11 bases is shown.  
637 Position 0 and 1 indicate cytosine and guanine, respectively. (C) Box plot of CGN/NCG motif  
638 ratios for long and short molecules. (D) Box plot of cleavage ratios between aggregating

639 position -4, -2, 1 and 4 and position -1. (E) AUCs for distinguishing patients with HCC from  
640 non-HCC subjects using cleavage ratios in (D). *P* values of differences among groups by  
641 Kruskal-Wallis tests. Post hoc pairwise *P* values by Wilcoxon rank sum tests with Benjamini-  
642 Hochberg adjustment shown above horizontal lines (C, D).

643

644 **Figure 6. Nuclease-mediated fragmentation in knockout mice.** (A) Pooled molecule size  
645 distributions of wild-type and nuclease-deficient mice from SMRT sequencing. Visualization  
646 of size in the range of 0 to 5000 bp and  $\log_{10}$ -transformed frequencies were used. (B) Zoom-  
647 in plot of (A) showing size in the range of 0 to 250 bp on the linear scale. (C) Boxplot  
648 showing the proportions of long molecules originated from 2000 bp upstream and  
649 downstream of transcription start sites. (D) Percent of changes in pooled long molecules  
650 (>500 bp) abundance relative to wild-type mice on low, medium and high expression gene  
651 groups.

652 **REFERENCES**

- 653 Al-Mayouf SM, Sunker A, Abdwani R, Abrawi SA, Almurshedi F, Alhashmi N, Al Sonbul A,  
654 Sewairi W, Qari A, Abdallah E, et al. 2011. Loss-of-function variant in DNASE1L3  
655 causes a familial form of systemic lupus erythematosus. *Nat Genet* **43**: 1186–1188.
- 656 Amemiya HM, Kundaje A, Boyle AP. 2019. The ENCODE blacklist: identification of  
657 problematic regions of the genome. *Sci Rep* **9**: 9354.
- 658 BAC Resource Consortium T, Cheung VG, Nowak N, Jang W, Kirsch IR, Zhao S, Chen X-N,  
659 Furey TS, Kim U-J, Kuo W-L, et al. 2001. Integration of cytogenetic landmarks into  
660 the draft sequence of the human genome. *Nature* **409**: 953–958.
- 661 Ballarino R, Bouwman BAM, Agostini F, Harbers L, Diekmann C, Wernersson E, Bienko M,  
662 Crosetto N. 2022. An atlas of endogenous DNA double-strand breaks arising during  
663 human neural cell fate determination. *Sci Data* **9**: 400.
- 664 Bickmore WA, Sumner AT. 1989. Mammalian chromosome banding — an expression of  
665 genome organization. *Trends Genet* **5**: 144–148.
- 666 Bouwman BAM, Agostini F, Garnerone S, Petrosino G, Gothe HJ, Sayols S, Moor AE,  
667 Itzkovitz S, Bienko M, Roukos V, et al. 2020. Genome-wide detection of DNA double-  
668 strand breaks by in-suspension BLISS. *Nat Protoc* **15**: 3894–3941.
- 669 Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, Shaw KRM,  
670 Ozenberger BA, Ellrott K, Shmulevich I, Sander C, Stuart JM. 2013. The Cancer  
671 Genome Atlas Pan-Cancer analysis project. *Nat Genet* **45**: 1113–1120.
- 672 Chaisson MJ, Tesler G. 2012. Mapping single molecule sequencing reads using basic local  
673 alignment with successive refinement (BLASR): application and theory. *BMC*  
674 *Bioinformatics* **13**: 238.
- 675 Chan KCA, Zhang J, Hui ABY, Wong N, Lau TK, Leung TN, Lo K-W, Huang DWS, Lo YMD.  
676 2004. Size distributions of maternal and fetal DNA in maternal plasma. *Clin Chem* **50**:  
677 88–92.
- 678 Chan RWY, Serpas L, Ni M, Volpi S, Hiraki LT, Tam L-S, Rashidfarrokhi A, Wong PCH, Tam  
679 LHP, Wang Y, et al. 2020. Plasma DNA profile associated with DNASE1L3 gene  
680 mutations: clinical observations, relationships to nuclease substrate preference, and  
681 in vivo correction. *Am J Hum Genet* **107**: 882–894.
- 682 Chiu RWK, Poon LLM, Lau TK, Leung TN, Wong EMC, Lo YMD. 2001. Effects of blood-  
683 processing protocols on fetal and total DNA quantification in maternal plasma. *Clin*  
684 *Chem* **47**: 1607–1613.
- 685 Choy LYL, Peng W, Jiang P, Cheng SH, Yu SCY, Shang H, Olivia Tse OY, Wong J, Wong  
686 VWS, Wong GLH, et al. 2022. Single-molecule sequencing enables long cell-free  
687 DNA detection and direct methylation analysis for cancer patients. *Clin Chem* **68**:  
688 1151–1163.
- 689 Crosetto N, Mitra A, Silva MJ, Bienko M, Dojer N, Wang Q, Karaca E, Chiarle R, Skrzypczak  
690 M, Ginalski K, et al. 2013. Nucleotide-resolution DNA double-strand break mapping  
691 by next-generation sequencing. *Nat Methods* **10**: 361–365.

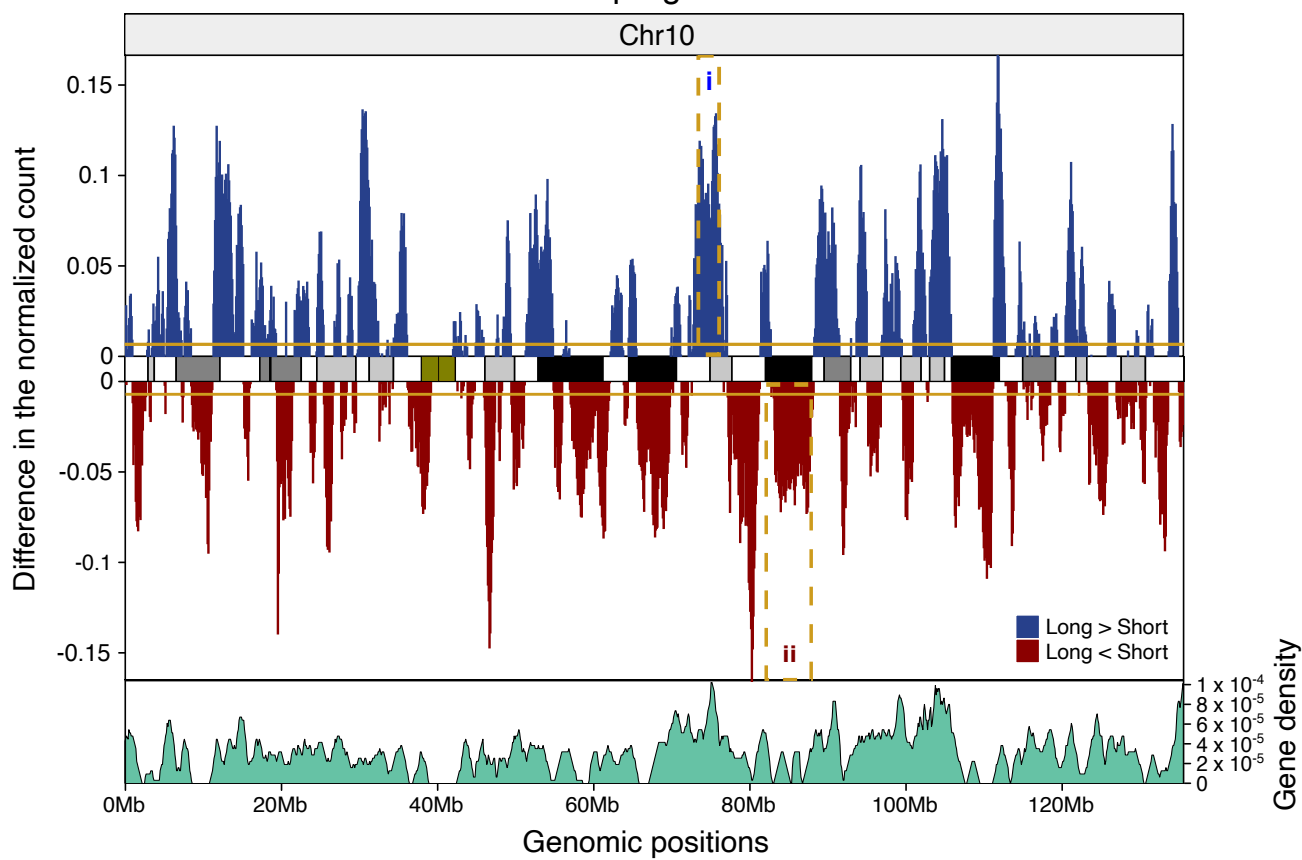
- 692 Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis CA, Doyle F, Epstein CB, Frietze S,  
693 Harrow J, Kaul R, et al. 2012. An integrated encyclopedia of DNA elements in the  
694 human genome. *Nature* **489**: 57–74.
- 695 Eisenberg E, Levanon EY. 2013. Human housekeeping genes, revisited. *Trends Genet* **29**:  
696 569–574.
- 697 Grimwood J, Gordon LA, Olsen A, Terry A, Schmutz J, Lamerdin J, Hellsten U, Goodstein D,  
698 Couronne O, Tran-Gyamfi M, et al. 2004. The DNA sequence and biology of human  
699 chromosome 19. *Nature* **428**: 529–535.
- 700 Han DSC, Ni M, Chan RWY, Chan VWH, Lui KO, Chiu RWK, Lo YMD. 2020. The Biology of  
701 cell-free DNA fragmentation and the roles of DNASE1, DNASE1L3, and DFFB. *Am J*  
702 *Hum Genet* **106**: 202–214.
- 703 Han DSC, Ni M, Chan RWY, Wong DKL, Hiraki LT, Volpi S, Jiang P, Lui KO, Chan KCA, Chiu  
704 RWK, et al. 2021. Nuclease deficiencies alter plasma cell-free DNA methylation  
705 profiles. *Genome Res* **31**: 2008–2021.
- 706 Ivanov M, Baranova A, Butler T, Spellman P, Mileyko V. 2015. Non-random fragmentation  
707 patterns in circulating cell-free DNA reflect epigenetic regulation. *BMC Genomics* **16**:  
708 S1.
- 709 Jiang P, Chan CWM, Chan KCA, Cheng SH, Wong J, Wong VW-S, Wong GLH, Chan SL,  
710 Mok TSK, Chan HLY, et al. 2015. Lengthening and shortening of plasma DNA in  
711 hepatocellular carcinoma patients. *Proc Natl Acad Sci U S A* **112**: E1317–E1325.
- 712 Katsman E, Orlanski S, Martignano F, Fox-Fisher I, Shemer R, Dor Y, Zick A, Eden A, Petrini  
713 I, Conticello SG, et al. 2022. Detecting cell-of-origin and cancer-specific methylation  
714 features of cell-free DNA from nanopore sequencing. *Genome Biol* **23**: 158.
- 715 Larsen BD, Sørensen CS. 2017. The caspase-activated DNase: apoptosis and beyond. *The*  
716 *FEBS Journal* **284**: 1160–1170.
- 717 Lau BT, Almeda A, Schauer M, McNamara M, Bai X, Meng Q, Partha M, Grimes SM, Lee H,  
718 Heestand GM, et al. 2023. Single-molecule methylation profiles of cell-free DNA in  
719 cancer with nanopore sequencing. *Genome Med* **15**: 33.
- 720 Lensing SV, Marsico G, Hänsel-Hertsch R, Lam EY, Tannahill D, Balasubramanian S. 2016.  
721 DSBCapture: in situ capture and sequencing of DNA breaks. *Nat Methods* **13**: 855–  
722 857.
- 723 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**:  
724 3094–3100.
- 725 Lo YMD, Chan KCA, Sun H, Chen EZ, Jiang P, Lun FMF, Zheng YW, Leung TY, Lau TK,  
726 Cantor CR, et al. 2010. Maternal plasma DNA sequencing reveals the genome-wide  
727 genetic and mutational profile of the fetus. *Sci Transl Med* **2**: 61ra91-61ra91.
- 728 Lo YMD, Han DSC, Jiang P, Chiu RWK. 2021. Epigenetics, fragmentomics, and topology of  
729 cell-free DNA in liquid biopsies. *Science* **372**, eaaw3616.
- 730 Meuleman W, Muratov A, Rynes E, Halow J, Lee K, Bates D, Diegel M, Dunn D, Neri F,  
731 Teodosiadis A, et al. 2020. Index and biological spectrum of human DNase I  
732 hypersensitive sites. *Nature* **584**: 244–251.

- 733 Mouliere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, Mair R,  
734 Goranova T, Marass F, Heider K, et al. 2018. Enhanced detection of circulating tumor  
735 DNA by fragment size analysis. *Sci Transl Med* **10**, eaat4921.
- 736 Nagata S, Nagase H, Kawane K, Mukae N, Fukuyama H. 2003. Degradation of  
737 chromosomal DNA during apoptosis. *Cell Death Differ* **10**: 108–116.
- 738 Serpas L, Chan RWY, Jiang P, Ni M, Sun K, Rashidfarrokhi A, Soni C, Sisirak V, Lee W-S,  
739 Cheng SH, et al. 2019. Dnase1l3 deletion causes aberrations in length and end-motif  
740 frequencies in plasma DNA. *Proc Natl Acad Sci U S A* **116**: 641–649.
- 741 Sisirak V, Sally B, D'Agati V, Martinez-Ortiz W, Özçakar ZB, David J, Rashidfarrokhi A, Yeste  
742 A, Panea C, Chida AS, et al. 2016. Digestion of chromatin in apoptotic cell  
743 microparticles prevents autoimmunity. *Cell* **166**: 88–101.
- 744 Snyder MW, Kircher M, Hill AJ, Daza RM, Shendure J. 2016. Cell-free DNA comprises an in  
745 vivo nucleosome footprint that informs its tissues-of-origin. *Cell* **164**: 57–68.
- 746 Sun K, Jiang P, Chan KCA, Wong J, Cheng YKY, Liang RHS, Chan W, Ma ESK, Chan SL,  
747 Cheng SH, et al. 2015. Plasma DNA tissue mapping by genome-wide methylation  
748 sequencing for noninvasive prenatal, cancer, and transplantation assessments. *Proc*  
749 *Natl Acad Sci U S A* **112**: E5503–E5512.
- 750 Sun K, Jiang P, Cheng SH, Cheng THT, Wong J, Wong VWS, Ng SSM, Ma BBY, Leung TY,  
751 Chan SL, et al. 2019. Orientation-aware plasma cell-free DNA fragmentation analysis  
752 in open chromatin regions informs tissue of origin. *Genome Res* **29**: 418–427.
- 753 Tabula Muris Consortium, Overall coordination, Logistical coordination, Organ collection and  
754 processing, Library preparation and sequencing, Computational data analysis, Cell  
755 type annotation, Writing group, Supplemental text writing group, Principal  
756 investigators. 2018. Single-cell transcriptomics of 20 mouse organs creates a Tabula  
757 Muris. *Nature* **562**: 367–372.
- 758 Toné S, Sugimoto K, Tanda K, Suda T, Uehira K, Kanouchi H, Samejima K, Minatogawa Y,  
759 Earnshaw WC. 2007. Three distinct stages of apoptotic nuclear condensation  
760 revealed by time-lapse imaging, biochemical and electron microscopy analysis of  
761 cell-free apoptosis. *Exp Cell Res* **313**: 3635–3644.
- 762 Tse OYO, Jiang P, Cheng SH, Peng W, Shang H, Wong J, Chan SL, Poon LCY, Leung TY,  
763 Chan KCA, et al. 2021. Genome-wide detection of cytosine methylation by single  
764 molecule real-time sequencing. *Proc Natl Acad Sci U S A* **118**: e2019768118.
- 765 Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, Wölfler A, Zebisch A, Gerger A,  
766 Pristauz G, et al. 2019. Inference of transcription factor binding from cell-free DNA  
767 enables tumor subtype prediction and early detection. *Nature Commun* **10**: 4666.
- 768 Ulz P, Thallinger GG, Auer M, Graf R, Kashofer K, Jahn SW, Abete L, Pristauz G, Petru E,  
769 Geigl JB, et al. 2016. Inferring expressed genes by whole-genome sequencing of  
770 plasma DNA. *Nat Genet* **48**: 1273.
- 771 Yu SCY, Chan KCA, Zheng YWL, Jiang P, Liao GJW, Sun H, Akolekar R, Leung TY, Go ATJI,  
772 Vugt JMG van, et al. 2014. Size-based molecular diagnostics using plasma DNA for  
773 noninvasive prenatal testing. *Proc Natl Acad Sci U S A* **111**: 8583–8588.

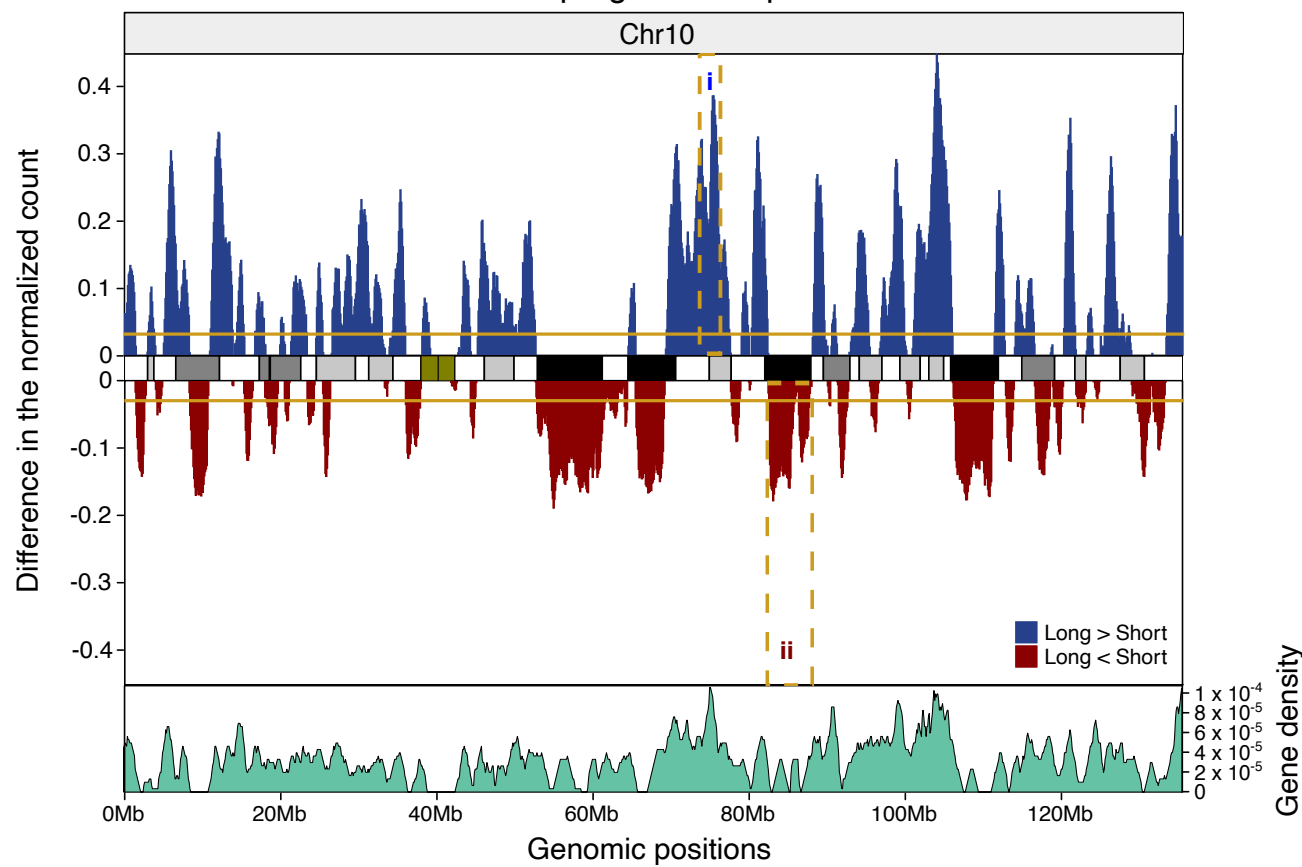
- 774 Yu SCY, Choy LYL, Lo YMD. 2023a. 'Longing' for the next generation of liquid biopsy: the  
775 diagnostic potential of long cell-free DNA in oncology and prenatal testing. *Mol Diagn*  
776 *Ther* **27**: 563-571.
- 777 Yu SCY, Deng J, Qiao R, Cheng SH, Peng W, Lau SL, Choy LYL, Leung TY, Wong J, Wong  
778 VW-S, et al. 2023b. Comparison of single molecule, real-time sequencing and  
779 nanopore sequencing for analysis of the size, end-motif, and tissue-of-origin of long  
780 cell-free DNA in plasma. *Clin Chem* **69**: 168–179.
- 781 Yu SCY, Jiang P, Peng W, Cheng SH, Cheung YTT, Tse OYO, Shang H, Poon LC, Leung TY,  
782 Chan KCA, et al. 2021. Single-molecule sequencing reveals a large population of  
783 long cell-free DNA molecules in maternal plasma. *Proc Natl Acad Sci U S A* **118**:  
784 e2114937118.
- 785 Zhou Q, Kang G, Jiang P, Qiao R, Lam WKJ, Yu SCY, Ma M-JL, Ji L, Cheng SH, Gai W, et al.  
786 2022. Epigenetic analysis of cell-free DNA by fragmentomic profiling. *Proc Natl Acad*  
787 *Sci U S A* **119**: e2209852119.
- 788 Zhou Z, Ma M-JL, Chan RWY, Lam WKJ, Peng W, Gai W, Hu X, Ding SC, Ji L, Zhou Q, et al.  
789 2023. Fragmentation landscape of cell-free DNA revealed by deconvolutional  
790 analysis of end motifs. *Proc Natl Acad Sci U S A* **120**: e2220982120.
- 791 Ziller MJ, Gu H, Müller F, Donaghey J, Tsai LT-Y, Kohlbacher O, De Jager PL, Rosen ED,  
792 Bennett DA, Bernstein BE, et al. 2013. Charting a dynamic DNA methylation  
793 landscape of the human genome. *Nature* **500**: 477–481.
- 794

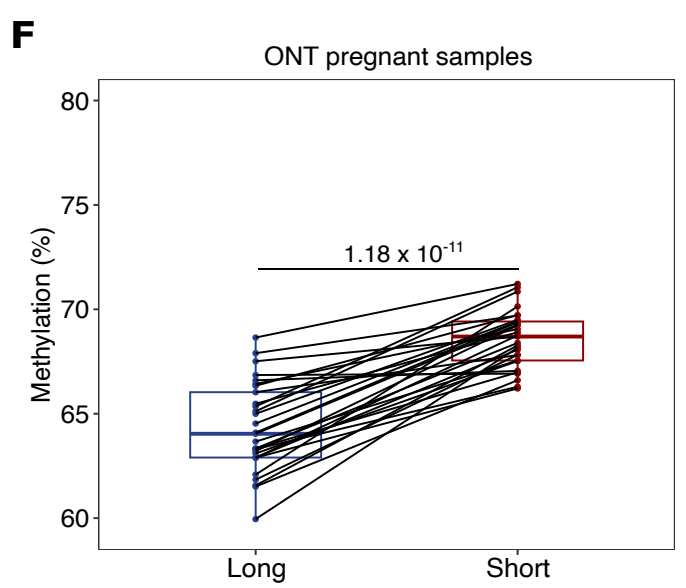
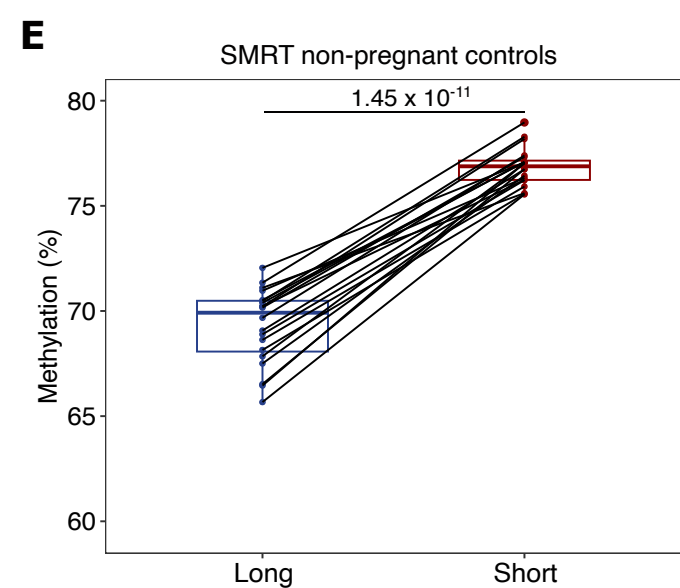
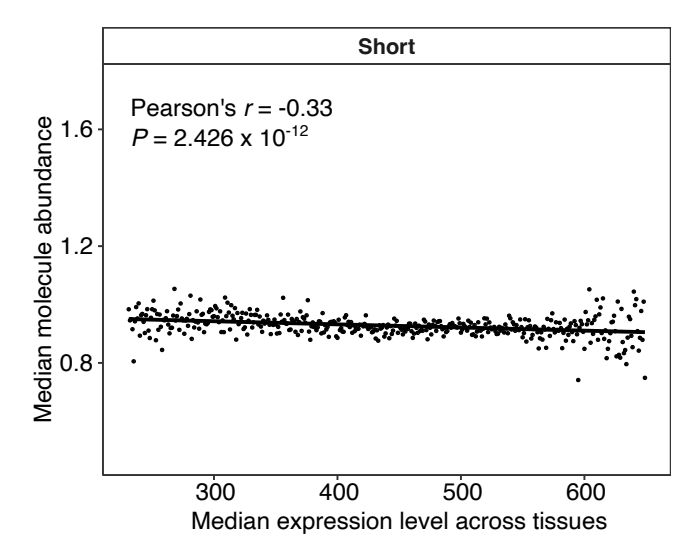
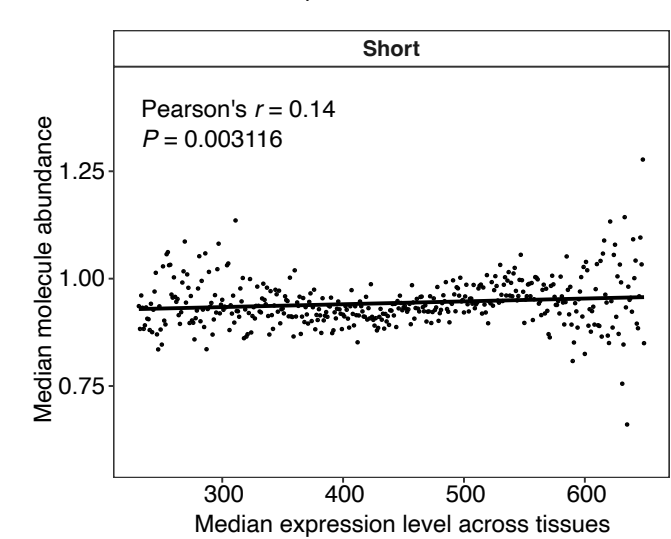
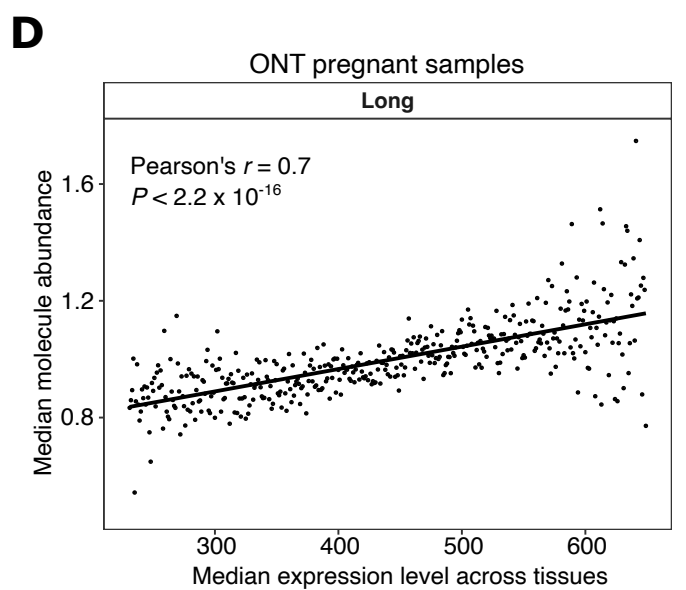
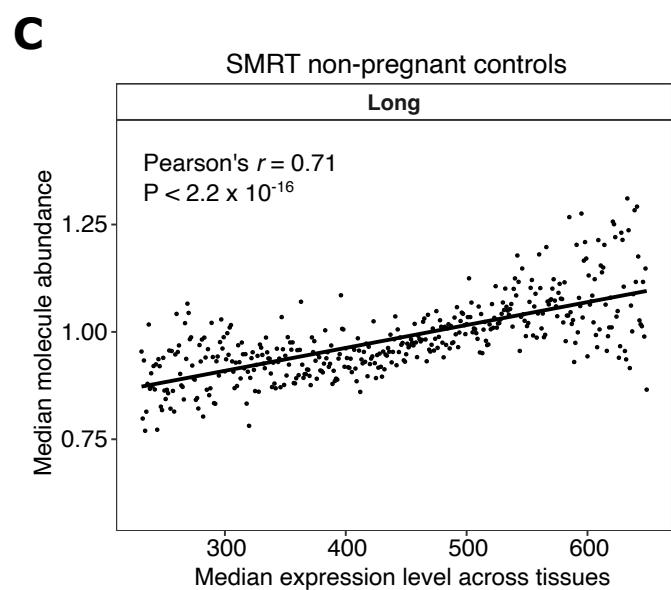
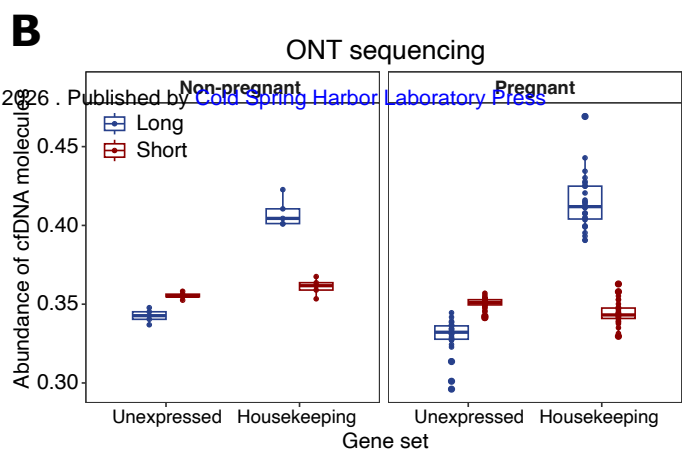
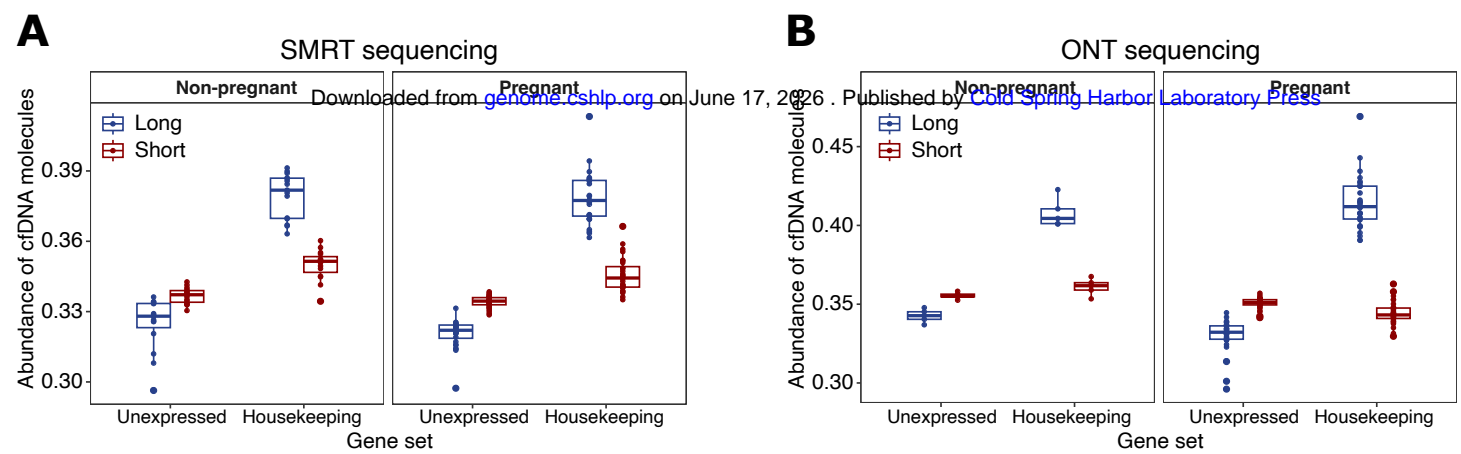
**A**

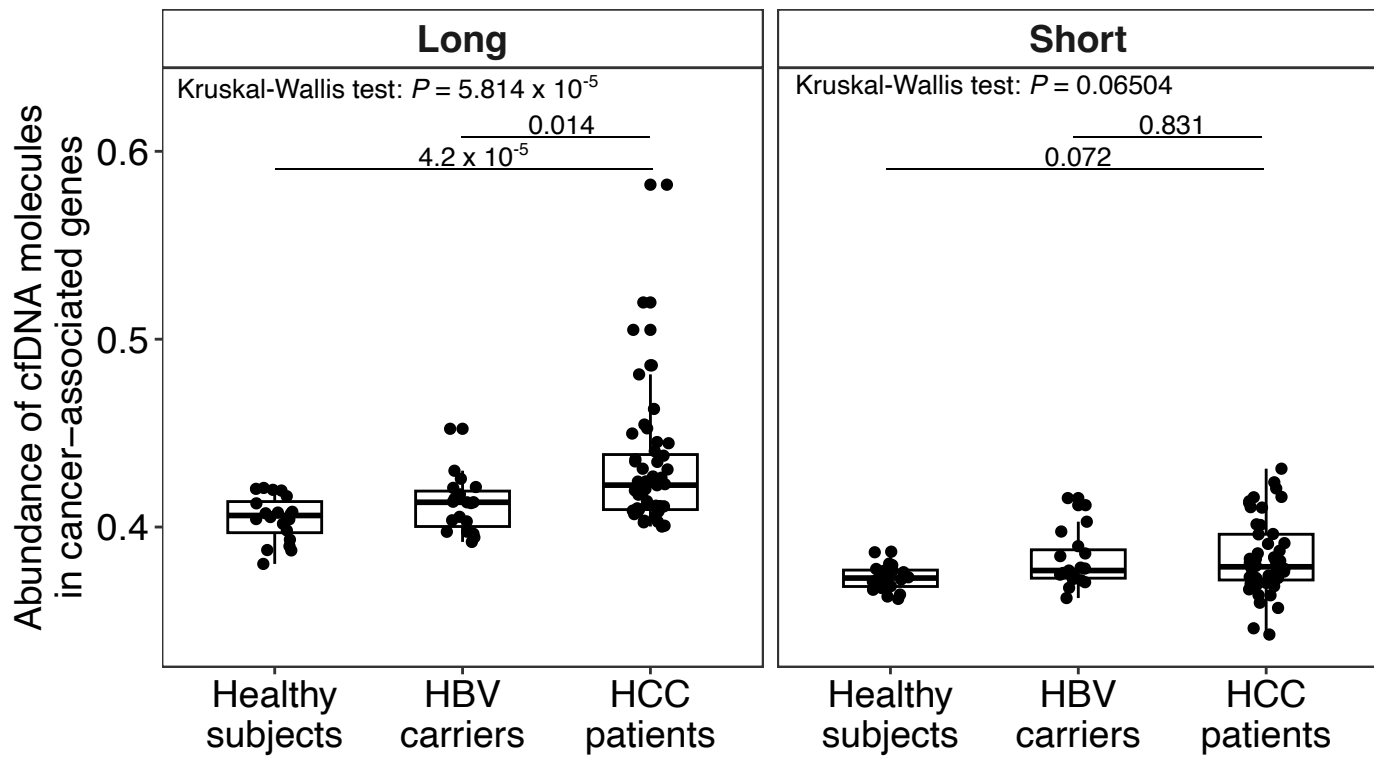
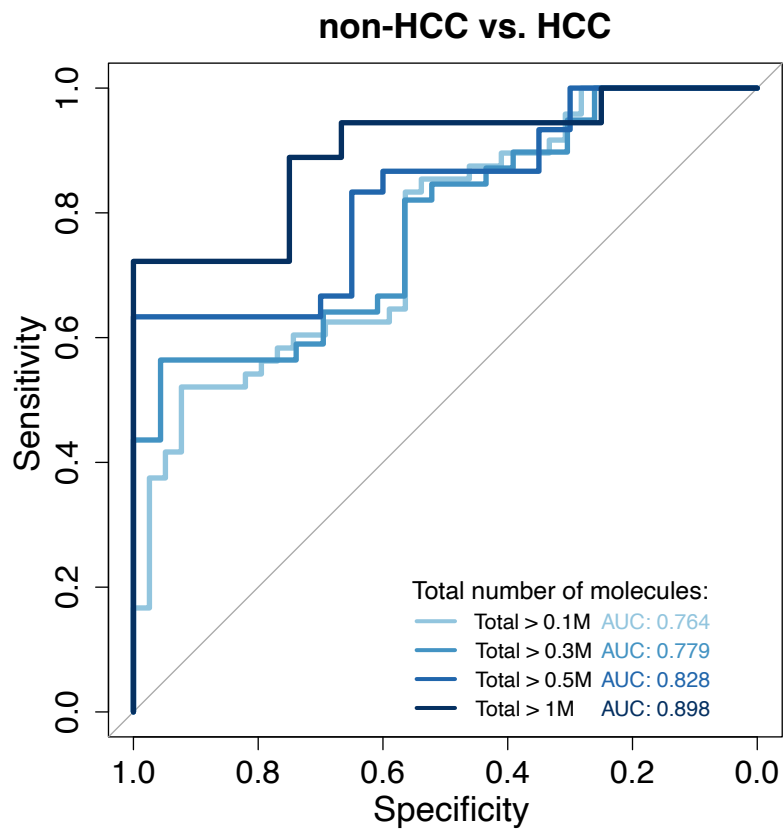
## SMRT non-pregnant controls

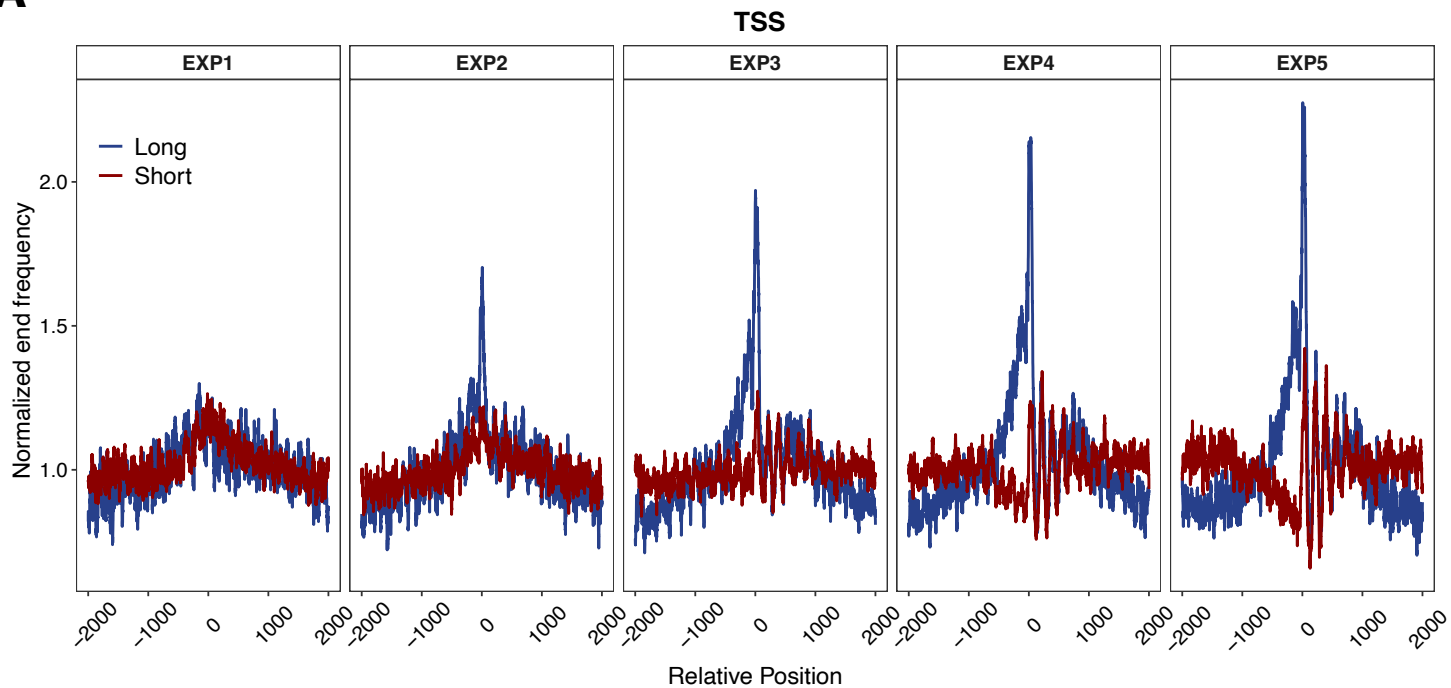
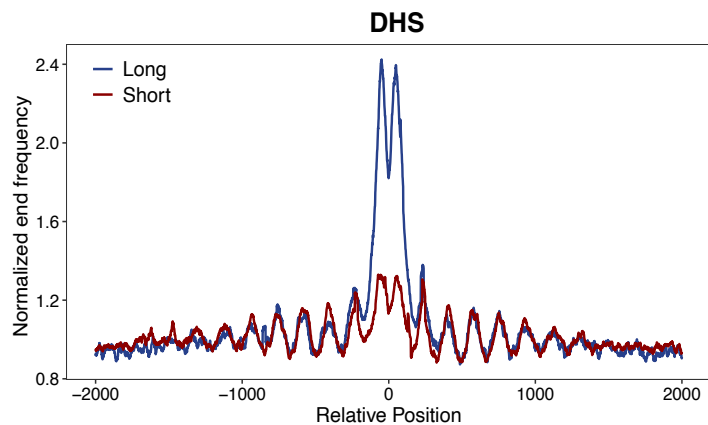
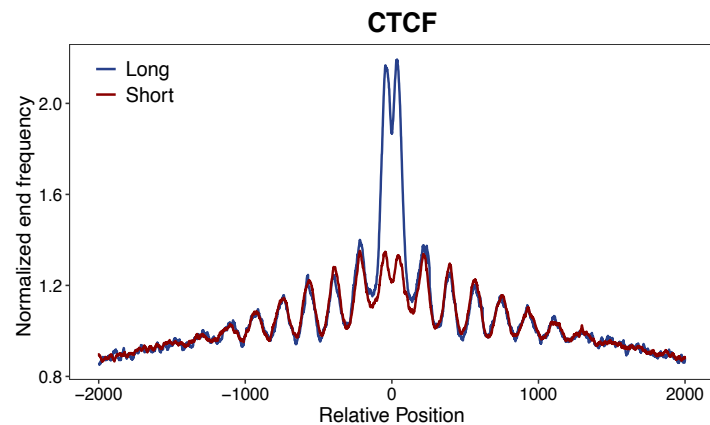
**B**

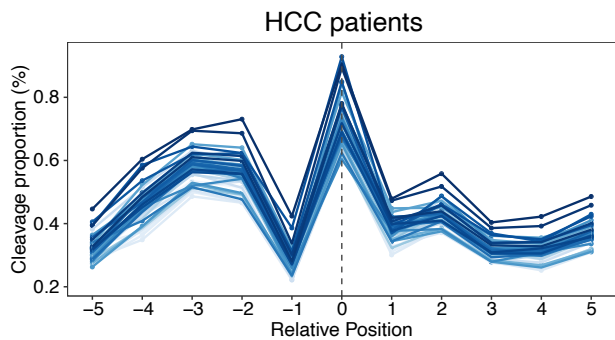
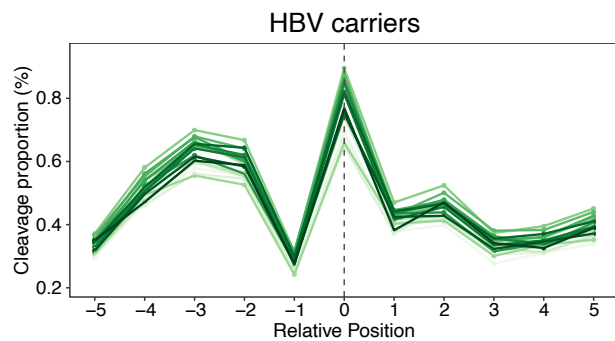
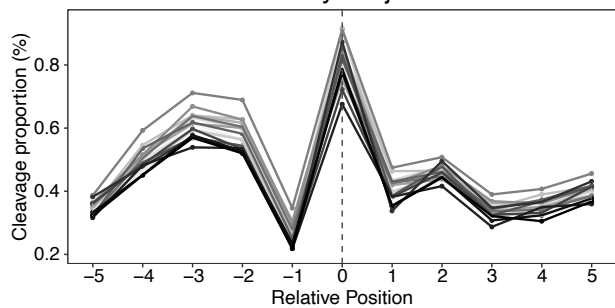
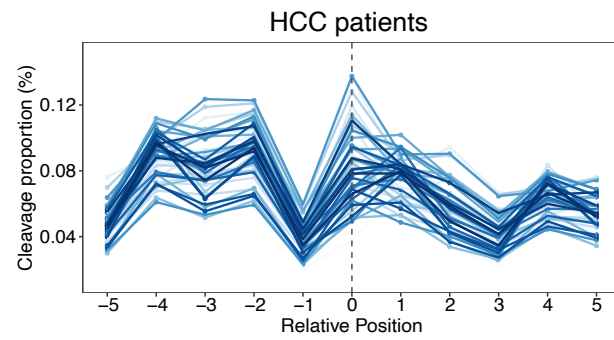
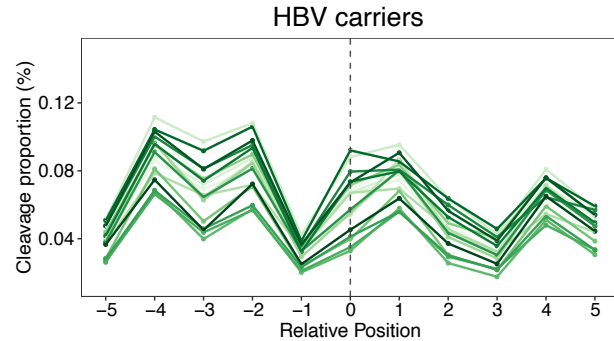
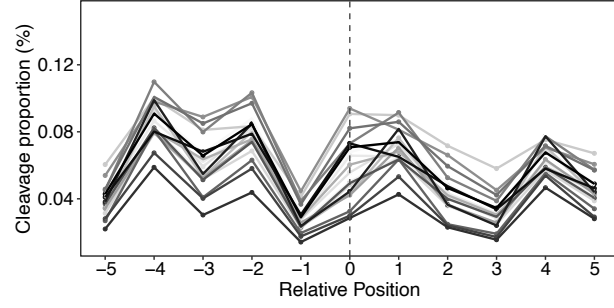
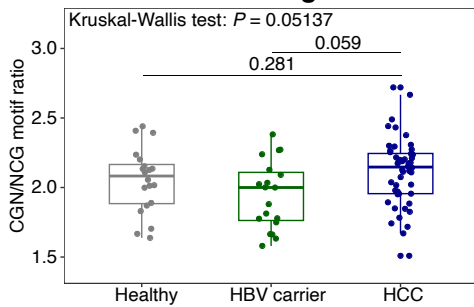
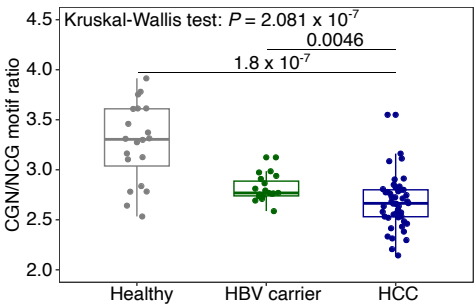
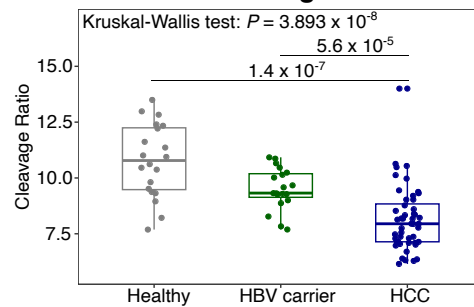
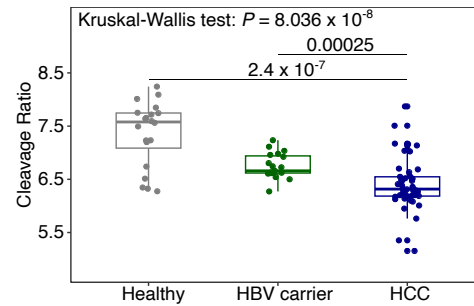
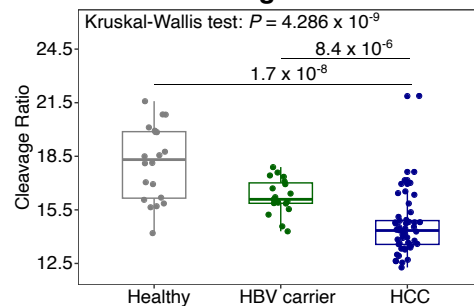
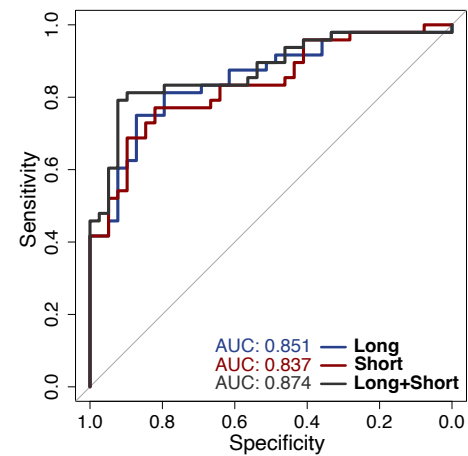
## ONT pregnant samples

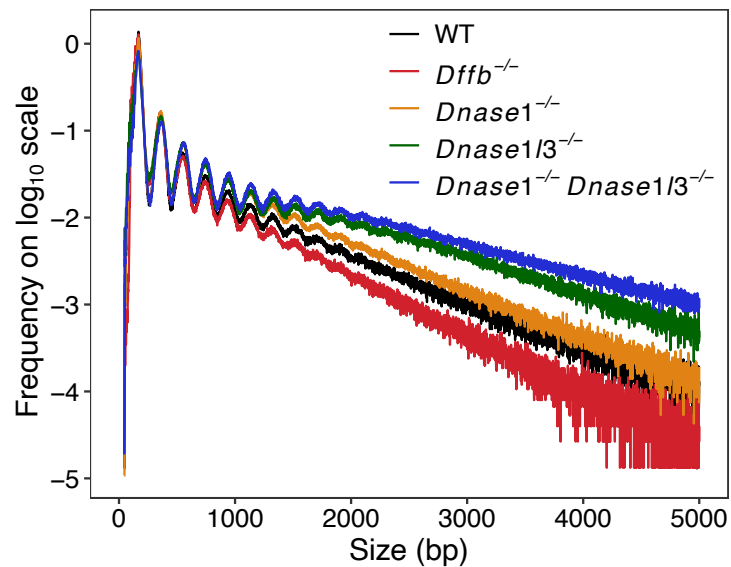
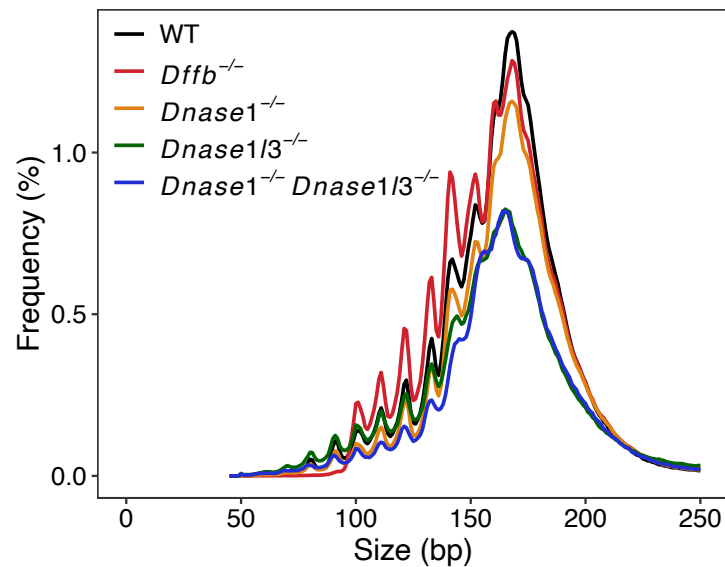
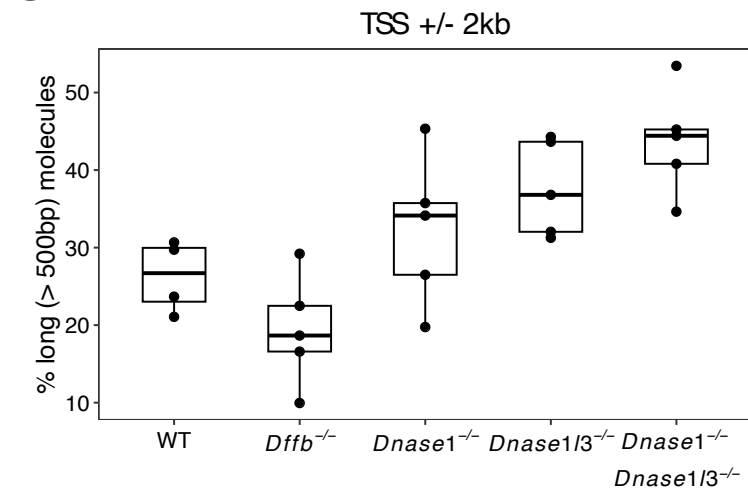




**A****B**

**A****B****C**

**A****Short molecules**Downloaded from [genome.cshlp.org](http://genome.cshlp.org) on June 17, 2026 . Published by [Cold Spring Harbor Laboratory Press](http://ColdSpringHarborLaboratoryPress)**B****Long molecules****C****Long****Short****D****Long****Short****Long+Short****E****non-HCC vs. HCC**

**A****B****C****D**