

Repeats Mimic Pathogen-Associated Patterns Across a Vast Evolutionary Landscape

Petr Šulc^{1,2,*†}, Andrea Di Gioacchino^{3,*†}, Alexander Solovyov^{4,*}, Sajid A. Marhon^{5,+}, Siyu Sun^{4,+},
Håvard T Lindholm^{5,+}, Raymond Chen⁵, Amir Hosseini^{5,6}, Hua Jiang⁷, Bao-Han Ly⁸, Parinaz
Mehdipour^{5,6}, Omar Abdel-Wahab^{9,10}, Nicolas Vabret¹¹, John LaCava^{7,8,†}, Daniel D. De
Carvalho^{5,10,#,†}, Rémi Monasson^{3,#}, Simona Cocco^{3,#,†}, Benjamin D. Greenbaum^{4,13,#,†}

¹ School of Molecular Sciences and Center for Molecular Design and Biomimetics, The Biodesign Institute, Arizona State University, Tempe, AZ 85281, USA

² Life and Medical Sciences (LIMES) Institute, University of Bonn, 53121 Bonn, Germany

³ Laboratoire de Physique de l'Ecole Normale Supérieure, PSL & CNRS UMR8063, Sorbonne Université, Université de Paris, Paris, France

⁴ Computational Oncology, Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

⁵ Princess Margaret Cancer Centre, University Health Network, Toronto, ON M5G 1L7, Canada

⁶ Ludwig Institute for Cancer Research, Nuffield Department of Medicine, University of Oxford, Oxford, OX3 7DQ, UK

⁷ Laboratory of Cellular and Structural Biology, The Rockefeller University, New York, NY 10065, USA

⁸ European Research Institute for the Biology of Ageing, University Medical Center Groningen, Groningen, The Netherlands

⁹ Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

¹⁰ Leukemia Service, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

¹¹ Tisch Cancer Institute, Icahn School of Medicine at Mount Sinai, New York, NY 10029 USA

¹² Department of Medical Biophysics, University of Toronto, Toronto, ON M5G 1L7, Canada

¹³ Physiology, Biophysics & Systems Biology, Weill Cornell Medicine, Weill Cornell Medical College, New York, NY 10065, USA

* Denotes Equal Contribution

+ Denotes Equal Contribution

Denotes Senior Author

† Correspondence: psulc@asu.edu, andrea.digioacchino@phys.ens.fr, j.p.lacava@umcg.nl, Daniel.DeCarvalho@uhnresearch.ca, simona.cocco@phys.ens.fr, greenbab@mskcc.org

ABSTRACT

An emerging hallmark across human diseases – such as cancer, autoimmune and neurodegenerative disorders – is the aberrant transcription of typically silenced repetitive elements. Once active, a subset of repeats may be capable of “viral mimicry”: the display of pathogen-associated molecular patterns (PAMPs) that can, in principle, bind pattern recognition receptors (PRRs) of the innate immune system and trigger inflammation. Yet how to quantify the landscape of viral mimicry and how it is shaped by natural selection remains a critical gap in our understanding of both genome evolution and the immunological basis of disease. We propose a theoretical framework to quantify selective forces on virus-like features as the entropic cost a sequence pays to hold a non-self PAMP and show our approach can predict classes of viral-mimicry within the human genome and across eukaryotes. We quantify the breadth and conservation of viral mimicry across multiple species for the first time and integrate selective forces into predictive evolutionary models. We show HSATII and intact LINE-1 (L1) are under selection to maintain CpG motifs, and specific Alu families likewise maintain the proximal presence of inverted copies to form double-stranded RNA (dsRNA). We validate our approach by predicting high CpG L1 ligands of L1 proteins and the innate receptor *ZCCHC3*, and dsRNA present both intracellularly and as MDA5 ligands. We conclude viral mimicry is a general evolutionary mechanism whereby genomes co-opt pathogen-associated features generated by prone repetitive sequences, likely offering an advantage as a quality control system against transcriptional dysregulation.

53 **MAIN TEXT**

54
55 It has recently become clear that repetitive elements, which represent most of the human genome
56 and often derive from integrated viruses and genome parasites, can function as “non-self” pathogen-
57 associated molecular patterns (PAMPs). Under aberrant conditions, such as cancer¹ and viral
58 infection²⁻⁴, repeats may be overexpressed, where the PAMPs they display can engage the innate
59 immune system⁵⁻¹⁰. Consistently, a growing body of literature has demonstrated aberrant expression
60 of immunostimulatory repeats across an array of human diseases, such as in aging¹¹ and
61 autoimmunity¹², implying viral mimicry may be a fundamental feature across inflammatory diseases.
62 The ability to quantify PAMPs capable of being sensed by pattern recognition receptors (PRRs) is
63 also of considerable theoretical interest¹³. Mathematical models of the evolution of human H1N1
64 influenza since the 1918 pandemic showed an attenuation of CpG motifs, leading to the prediction
65 such motifs can be targeted by PRRs^{14,15}. It was subsequently discovered the protein ZAP
66 (*ZC3HAV1*) is a PRR targeting CpG motifs, indicating inferences drawn from genome evolution can
67 predict PAMPs and PRR specificities relevant to the adaptation of emerging viruses^{16,17}, including
68 SARS-CoV-2^{18,19}. It has been more difficult to predict specificities associated with structure
69 formation²⁰, such as the recognition of double-stranded RNA PAMPs by MDA5 (*IFIH1*) or TLR3²¹, or
70 other more complex PAMPs, such as the creation of DNA:RNA hybrids. Importantly, viral mimicry
71 can be leveraged therapeutically: the expression of immunostimulatory repeats is inducible by
72 epigenetic drugs, leading to triggering of innate immune sensors and induction of an interferon
73 response⁵⁻⁹.

74
75 Several fundamental questions remain, such as which PRRs can be activated by which specific
76 repeats, if viral mimicry serves a functional role in the genome as an evolved checkpoint for loss of
77 epigenetic regulation or genome fidelity, and whether tumors and pathogens adapt to manipulate
78 mimicry to their own selective advantage^{22,23}. In one evolutionary scenario, repeats which contain
79 PAMPs in somatically silenced regions can offer a fitness advantage to cells due to their ability to
80 trigger PRRs under epigenetic stress, eliminating dysregulated cells and maintaining tissue
81 homeostasis^{22,23}. Such features could then be maintained by natural selection. Alternatively, in a
82 neutral scenario, it may be that high RNA concentration resulting from transcriptional dysregulation
83 can engage PAMPs non-specifically, and their sensing is a by-product of dysregulation rather than of
84 selection acting on specific features. Discriminating between such scenarios is critical to
85 understanding how non-self mimicry by the self-genome evolved, and how it can be leveraged for
86 emerging therapies and honed for existing ones. There is therefore a pressing need for computational
87 approaches to quantify the presence of viral mimics, define their immunological features and quantify
88 the evolutionary dynamics of their (putative) PAMP content. We utilize a novel approach from
89 statistical physics to quantify nucleic-acid motifs and double-stranded structures under selective
90 forces and use multiple assays to validate our predictions. In doing so we define specific categories
91 of repeat families that were likely retained by natural selection to display viral mimics.

92 Inference and evolutionary dynamics of pathogen-associated patterns in repeats

93
94 We utilize the framework of selective and entropic forces to infer the presence and evolutionary
95 dynamics of a PAMP or set of PAMPs in a genome¹⁴. Sequences incorporated into a genome, subject
96 to constraints such as local nucleic acid content, accumulate mutations during evolution to resemble,
97 on average, the self-genome while selective forces, such as those acting on an atypical “non-self”
98 pattern, oppose such a trend (**Fig. 1A**). The selective force is an intensive parameter that can be
99 interpreted as a measure of the depletion (negative) or excess (positive) of a feature in a genome
100 sequence beyond the degree it would be expected based on the nucleotide statistics within the
101 sequence. This framework is ideal for the study of viral mimicry and PAMP detection as selective
102 forces can be readily compared between groups of sequences independently of their length. We infer
103 selective forces on a given sequence for one or more patterns introducing them as parameters in a
104 Maximum Entropy model for genomic sequences²⁴ (Methods). Our inference algorithm uses exact
105 transfer matrix methods from statistical physics which, unlike earlier approaches¹⁵, are
106 computationally efficient (scaling with the length of the sequence) and facilitate the analysis of longer
107 sequences and large sequence datasets (Methods).

108
109 A repetitive element is primarily defined by the presence of multiple copies (inserts) of its sequence
110 in the genome. As additional repetitive copies accumulate, we measure the evolutionary dynamics of
111 PAMPs as they diverge from their original sequence. We use two approaches. The first uses
112 relaxation dynamics: a new repeat in a genome evolves until it reaches a genomic equilibrium value
113 determined by a balance of factors such as constraints on nucleic acid usage and forces on sequence
114 patterns. By analogy, selective forces drove the avian origin 1918 influenza virus towards a new
115 equilibrium in humans with lower CpG content, and it was subsequently found the PRR ZAP
116 consistently targets CpG with greater affinity in humans than birds²⁵. The second approach uses a
117 Kimura-based model as a proxy for the neutral evolution of a sequence with given PAMP content.
118 We implement this variant of the Kimura model numerically to provide a null model of repeat evolution
119 within a genome. As in the standard Kimura model, we use different mutation probabilities for
120 transitions (a purine mutating into a purine or a pyrimidine into a pyrimidine) and transversions (a
121 purine mutating into a pyrimidine or vice-versa), with the former being more probable than the latter
122 (see Methods). Additionally, we use different ratios of mutation rates corresponding to nucleotide
123 transitions and transversions in CpG and non-CpG context²⁶. We calculated the dinucleotide
124 distribution stationary value, obtained as the stationary vector of the stochastic matrix with entries
125 corresponding to probabilities of mutating from one dinucleotide to another dinucleotide (see Methods
126 and Table 1 therein).

127
128 To test our approach, we quantify the overall degree of motif usage similarity between families of
129 human infecting viruses and regions of the human genome. We infer Maximum Entropy models with
130 forces on all single, di-, and tri-nucleotide motifs for a set of human repeat families and compare them
131 to models inferred for families of viruses which infect humans (**Fig. 1B, C**). To quantify similarity of
132 motif usage in the two sets of families we use the symmetrized Kullback-Leibler divergence (details
133 about its computation are given in Methods) between the corresponding models. Primarily viral and
134 human genomes share similar overall motif usage, a form of mimicry that is likely a product of shared
135 constraints on nucleotide usage across organisms and viruses, with some minor variation between
136 viral families. Coding regions in the human genome show stronger overall similarity to human infecting
137 viruses, most of whose genomes are devoted to coding, than non-coding regions, although large
138 variation exists in the latter. For instance, consistent with the overall trend, HERVK repeats show the
139 strongest similarity with viruses among repeat regions. As a stark exception, we find far less motif
140 usage similarity between Alu repeats and HSATII than either to the rest of the human genome or to
141 human infecting viruses. Neither repeats encode known proteins and both are thought to have non-

142 viral origins²⁷, indicating such regions may be subject to different evolutionary pressures from the
143 other repeats considered here.

144

145 **Landscape of repeats with selective forces on CpG dinucleotides**

146

147 CpG dinucleotides in humans are PAMPs in DNA, recognized via TLR9²⁸, and, as has been seen
148 more recently, CpGs in RNA are engaged via ZAP¹⁷. We compare the evolution of individual
149 dinucleotide motifs (quantified by calculating the selective force, x_m , on a dinucleotide motif, m , as
150 defined in Methods) between the original consensus sequence, representing the sequence most
151 likely to resemble a founding ancestral repeat insertion, and its subsequent copies in the genome
152 (**Fig. 2**). We analyzed x_{CpG} , and all other x_m , for all repeat families annotated in the DFAM database²⁹,
153 finding outliers such as Alu repeats and HSATII, the latter consistent with previous results¹⁰. Typically,
154 CpG content in the human genome is highly underrepresented as CpG sites mutate at a much faster
155 rate than the rest of the genome^{26,30,31}. We plot the mean difference in x_{CpG} per repeat family versus
156 x_{CpG} for the consensus insert (**Fig. 2A**). Consistent with our assumptions, we see families where the
157 selective force on CpG dinucleotides for the progenitor insert was greater than -1.9 have decreased
158 their force to this value, while those less than -1.9 have increased their value. We therefore establish
159 a genome-wide equilibrium in line with equilibria observed for human adapted viruses such as
160 influenza^{15,18}. If a repeat is not subject to selection, one would expect its insertions to evolve according
161 to a Kimura model with respective mutation rates for transitions and transversions, an approach used
162 in sequence evolution models to explain lower CpG content in vertebrate genomes³²⁻³⁴. **Fig. 2B**
163 shows the relaxation of x_{CpG} as a function of the Kimura distance³⁵ used for each individual repeat
164 sequence, as a proxy for time since insertion. The Kimura distance is the expected number of
165 mutations accumulated in a given period of time by a sequence that evolves with a higher probability
166 of transitions over transversions. It represents the expected number of differences between two
167 sequences after a given period of time at fixed mutational rates. We use it as a proxy for time since
168 insertion for each individual repeat sequence, relative to other elements of the same family. Most
169 repeat families show relaxation to the mean genome force expected from the neutral model, further
170 implying HSATII may be specifically under selection to hold this PAMP. Moreover, HSATII is the most
171 represented repeat among those overlapping with high- x_{CpG} ($x_{\text{CpG}} > 0$) genomic regions in the human
172 genome (**Fig 2C**).

173

174 As L1 elements have the most copies in the genome, they are most amenable to our approach. Their
175 copies are estimated to constitute about 20% of human genome³⁶. Here we only consider full-length
176 inserts, as annotated in L1Base2, and contrast those designated as fully intact (denoted FLI), from
177 full-length sequences designated as non-intact (FLNI)³⁷. Fully functional L1 DNA sequences are
178 regulated by hyper-methylation at CpG sites in their promoter, to inhibit their transcription^{38,39}. Indeed,
179 we find FLI L1 have higher CpG content than FLNI (**Fig. 2D**), though most conserved CpGs are not
180 in the promoter region (**Supplementary Fig. S1**). We find that as a L1 genome insertion ceases to
181 contain an intact copy, its CpG content decays with the Kimura distance to the consensus, reaching
182 the genome mean in a predictable way according to the Kimura model for neutral genomic evolution
183 (**Fig. 2E**). The most recent inserts into the human genome therefore appear to not have equilibrated.
184 It is important to identify all such cases because families that have not saturated are candidates for
185 viral mimicry via PAMP display, such as when LINE-1 is overexpressed in tumors^{1,40-42}. For Alu
186 repeats we observe a pattern of CpG-content relaxation similar to LINE-1, but only when considering
187 together the major Alu subfamilies (AluY, AluS, AluJ). The younger AluY and, to a lesser extent, AluS
188 are not yet equilibrated and still possess PAMP-like high CpG content. (**Extended Data Fig. 1A**). For
189 HSATII, evolutionary dynamics of the force relaxation (**Fig. 2E**) corresponds to saturation at force
190 approximately equal to -0.4 , well above the equilibrium distribution computed from the Kimura model,
191 implying its ability to retain CpGs is maintained by selection. For most families the data points are

192 scarce and noisy, making a relaxation fit such as the one shown for HSATII and LINE-1 difficult.
193 **Supplementary Table 1** lists the full repeat atlas of CpG content, computed both for the consensus
194 repeat and as an average over the inserts in the genome. CpG-rich regions ($x_{\text{CpG}} > 0$) in the human
195 genome mostly concentrate in intergenic and, to a lesser extent, intronic regions (**Extended Data**
196 **Fig. 1B**), and are listed in **Supplementary Table 2**.

197
198 We reasoned that the force acting on CpGs in intact L1 species is enforced by the *in cis* binding of
199 L1 encoded proteins. To determine if this is the case, we analyzed the RNAs affiliated with both L1
200 ORF1p and ORF2p by RNA co-immunoprecipitation sequencing (RIP-seq). Notably ORF2p is the
201 reverse transcriptase of L1. We conducted α -ORF1p RIP-seq in N2102Ep human embryonal
202 carcinoma cell lines⁴³ and, for the first time, α -ORF2p RIP-seq. Transcripts enriched by co-IP were
203 determined by differential expression analysis of RIP versus total RNA and matched mock IP controls
204 (Methods). Intact L1s are exclusively recovered in ORF1p and ORF2p binding transcripts (**Fig. 2F**;
205 **Extended Data Fig. 1C**) as expected. Consistently, ORF2p (**Fig. 2F**) and ORF1p (**Extended Data**
206 **Fig. 1C**) enriched transcripts ($\text{Log2FC} > 3$, adj. p-val. < 0.05) have CpG forces consistent with the
207 high x_{CpG} observed in fully intact transcripts and compared to controls ($\text{Log2FC} < -3$, adj. p-val. $<$
208 0.05). To further examine if innate immune receptors co-IP with high x_{CpG} L1 RNA, we examined the
209 ligands of *ZCCHC3*, a protein recently described as a co-sensor of cGAS⁴⁴ that has been found to
210 interact with L1 ribonucleoprotein in an RNA-dependent fashion^{45,46}. We find a substantial enrichment
211 in high CpG L1 RNA associating *ZCCHC3*, compared to both controls and to non-intact L1 (**Fig. 2G**).
212 Our findings indicate high x_{CpG} L1 RNA is both more likely to associate with L1 proteins and with a
213 putative innate immune sensor of L1 RNA. We therefore conclude high x_{CpG} is associated with both
214 replication competent L1 and innate immune sensing of L1.

215 216 **Landscape and evolution of repeats with selective forces on double-stranded RNA formation**

217
218 We further extend our approach, for the first time, to the formation of anomalous secondary structure
219 by calculating the force on double-stranded RNA (dsRNA) formation in repeats (**Fig. 3, Extended**
220 **Data Fig. 2**). This value quantifies the tendency of an RNA transcript to form double-stranded
221 segments. It is generally accepted that Toll-Like Receptor 3 (TLR3) is activated by short (approx. 30
222 bp) endosomal dsRNA and Retinoic acid-Inducible Gene I (RIG-I, *DDX58*) by short (tens of bases)
223 cytoplasmic dsRNA accompanied by a triphosphosphate⁴⁰, while Melanoma Differentiation Associated
224 protein 5 (MDA5) recognizes longer cytoplasmic dsRNA associated with RNA virus replication⁴⁷. We
225 calculated the double stranded force, x_{ds} (Methods), for repetitive families as well as ncRNA and
226 mRNA sequences in the human genome, and randomly generated sequences (**Fig. 3A**). While the
227 mean value of x_{ds} computed for functional mRNA sequences and noncoding sequences is close to
228 zero and essentially the same as the value for random sequences, the consensus sequences of
229 repeats contain multiple families with long complementary segments contributing to an increased
230 average x_{ds} value (34 families out of 980 analyzed have $x_{\text{ds}} > 0.5$). While the general trend is to relax
231 x_{ds} towards zero (**Fig. 3B**), we observe outliers having a higher positive x_{ds} value, indicating a
232 possible reservoir of double-stranded segments being maintained by selection. Including are the DNA
233 transposons Tigger4a (**Extended Data Fig. 2B**), MER107 and MER6B (**Fig. 3B**), which could be
234 transcribed under aberrant conditions.

235
236 To locate possible sources of double-stranded segments originating from the same transcript, we
237 scan the entire genome (hg38 assembly), using a window of transcripts of length 3000bp, comparable
238 to typical lengths of long ncRNAs⁴⁸. We quantified the sequence complexity of such complementary
239 segments (based on Kolmogorov complexity, Methods), as shown in **Fig. 3C**. The segments close to
240 the low complexity limit typically contain a repeating motif of only a few nucleic acids (such as
241 poly(AT)) while the longest segments have higher complexity, i.e. the regions that can form long

242 dsRNA are not exclusively simple repeats (as summarized in an atlas of all families analyzed,
243 **Supplementary Table 2, 3**). We then characterized the distribution of forces in the genomic scan:
244 we observe two peaks, a major one close to 0 and a smaller around 0.5 (**Fig. 3D**). The mean length
245 of the longest complementary segments found in the dataset with $x_{ds} > 0.5$ is 40 base pairs. We found
246 that for the majority (88%) of such regions the complementary segments in the 3000 bases long
247 regions overlap with known repeats. Greater than 43% of identified complementary segments
248 correspond to AluS, where a copy has inserted in a positive orientation close to one in a negative
249 orientation (inverted-repeat Alus, IR-Alus). AluS is the most represented Alu family in the human
250 genome (accounting for more than 60% of the Alu inserts), and it also has the highest fraction of IR
251 inserts, 59% (39% for AluJ and 18% for AluY). In particular, we noticed 73% of the IR-Alus in the
252 human genome consist of AluS IRs, but if we filter for high- x_{ds} the share of AluS repeats forming IR-
253 Alus increases to 86% (**Fig. 3F**). Besides Alu subfamilies (which constitute about 50% of long
254 complementary segments that overlap with known inserts), we also identified complementary
255 fragments from the ORF2 open reading frame of LINE-1⁴⁹. We found previously unannotated non-
256 inverted repeats prone to forming long-double stranded RNA (full list in **Supplementary Table 3**).
257 We conclude that while IR-Alus form the major class of binders, other unannotated inverted repeats
258 are also prone to dsRNA formation. Likewise, previous work hypothesized that dsRNA formed from
259 introns is a checkpoint against intron retention^{50,51}. We observed most regions (55.4%) with $x_{ds} > 0.5$
260 were over-represented at intronic regions (**Extended Data Fig. 2G**).

261
262 We next validated our ability to predict double-stranded forming regions. We first examined two
263 published datasets of RNA forming long dsRNA MDA5 receptor ligands, as their transcription has
264 been implicated as a response to genome-wide DNA demethylation^{7,52}. We find that RNA transcripts
265 binding MDA5 under DNA demethylation agents (AZA) display a double-peaked force distribution with
266 a predominance of the large x_{ds} peak. Inverted repeats, and notably the AluS family (**Fig 3. D,F**), only
267 populate the high x_{ds} peak. AluS repeats account for 89% of enriched IR-Alus in the MDA5-binding
268 experiment, in agreement with our prediction based only on quantifying high- x_{ds} sequences in the
269 human genome. A consistent result was found in a second MDA5 ligand dataset (**Extended Data**
270 **Fig. 2C**⁵²). To further validate our ability to predict dsRNA forming transcripts, we generated a novel
271 dataset of sequenced ligands of the J2 monoclonal antibody, an antibody able to recognize dsRNA
272 of greater than 40 bp, nearly identical in length to the average length of anomalous regions predicted
273 when $x_{ds} > 0.5$, in a set of patient-derived colorectal cancer cell lines (Methods, **Extended Data Fig.**
274 **2E**). Consistent with our predictions, we show an enrichment of high x_{ds} regions in J2 antibody binding
275 transcripts, and with a similar profile as the previously published MDA5 ligands. In this case we found
276 AluS repeats constitute 84% of the enriched IR-Alus, once more in agreement with the value predicted
277 for high- x_{ds} sequences with our framework (**Fig. 3E**). These results, based solely on *in silico* analysis
278 of the human genome using our framework, are a striking quantification of the experimental
279 observation that IR-Alus, and especially AluS IR, are the major source of self-RNA that form MDA5
280 agonists⁷, providing strong validation of the predictive power of our evolutionary model and, in turn,
281 the hypothesis that evolution selected this feature as an epigenetic checkpoint^{22,23}. We further
282 analyzed a dataset of inhibitors of RNA splicing which induce intron retention⁵³. We examined RNA
283 sequencing data from SF3B inhibitors which cause the retention of introns in SF3B1 K700E mutant
284 cells. Consistent with our model, we found splicing agents which lead to intron retention over express
285 the high double-stranded force intronic repeats we predicted (**Extended Data Fig. 3, Supplementary**
286 **Table 5**), supporting the potential ability to manipulate this feature using a cancer therapeutic
287 targeting RNA splicing. Consistently, for inhibitors less associated with intron retention the effect was
288 either weakened or not present. We therefore show a clear ability to predict inverted repeat regions
289 associated dsRNA formation.

291 **Presence and evolution of PAMPs in genomes across evolutionary scales**

292

293 To further understand whether PAMPs are held by selection, we examine the presence of repeats
294 with high forces on CpG dinucleotides and double-stranded RNA across 20 genomes (**Extended**
295 **Data Figs. 4-6**). We calculate the presence of outliers for high x_{CpG} and high x_{ds} regions across all
296 species. For humans and mouse we show that the presence of such anomalous regions is not
297 primarily due to CpG islands or enhancer regions, based on the FANTOM database^{54,55} (Methods,
298 **Extended Data Fig. 4**). We find high x_{ds} regions occur across many species, even those which do
299 not have the Alu family of SINE elements, providing further evidence such regions are likely a
300 byproduct of the reverse transcription machinery across genomes rather than a function of Alus
301 specifically. To establish such regions in other organisms are not due to low complexity regions, we
302 plot the complexity of high x_{ds} regions for the zebrafish genome (**Fig. 4A**). We find many genomic
303 regions which are not low complexity and would be prone to dsRNA formation if transcribed, implying
304 such regions may be a source of PAMPs across species. To the best of our knowledge, this is the
305 first quantification of the presence of likely PAMP-forming repeat regions outside of primates. The full
306 list of regions with $x_{\text{ds}} > 0.5$ we discovered in the zebrafish genome is reported in **Supplementary**
307 **Table 4**.

308
309 For repeat families identified in humans, we compared selective forces across organisms
310 phylogenetically close to humans. We used the Hominoidea superfamily, whose most recent common
311 ancestor has been proposed to date back to about 16 million years ago⁵⁶. We performed the same
312 analysis as for the human genome, scanning the genomes of five small and great apes and
313 comparing sequences with high x_{CpG} and x_{ds} values. We first considered the high x_{CpG} windows
314 ($x_{\text{CpG}} > 0$) and computed the conservation of these sequences across organisms (as quantified by
315 the Overlap Index, Methods). We observed that the number of high x_{CpG} sequences conserved
316 between humans and other apes decreases exponentially with their evolutionary distance (**Fig. 4B**);
317 an expected result, given the high CpG mutation rate. We further observed that, although most high
318 x_{CpG} genomic windows do not overlap with any repeat, the vast majority of conserved x_{CpG} sequences
319 that do overlap with a repeat are associated with HSATII. HSATII can be found in primates after the
320 branching between the *Pongo* genus and the other great apes, allowing us to pinpoint the HSATII
321 insertion in the primate genomes between 13.8 and 8.9 million years ago. Remarkably, since its
322 insertion into the genome HSATII sequences are conserved to a much greater extent than other
323 sequences in the high-CpG pool (**Fig. 4C**), suggesting a selective pressure maintains PAMPs in
324 HSATII. When we next considered high- x_{ds} genomic windows, we found them to be much more
325 generally conserved than high x_{CpG} regions (**Fig. 4D**). We found these results striking, since it is not
326 expected by a null model of sequence evolution and implies a selective pressure to keep these
327 windows functionally intact. When focusing on sequences overlapping with repeats, we confirmed
328 inverted Alu repeats are highly conserved in time since their appearance in primate genomes more
329 than 16.3 million year ago (**Fig. 4E**). We therefore conclude Alus, and particularly the AluS family,
330 are likely to have selectively maintained the ability to form double-stranded RNA.

DISCUSSION

We quantify the landscape and evolutionary dynamics of viral mimicry both across and within genomes. We find, generally, that virus infecting humans mimic the motif usage statistics of human coding regions, indicating shared global constraints on motifs for both viruses and their hosts, consistent with our previous work^{14,15}. There are strong exceptions, such as Alu repeats and HSATII, under less constraint. We find the high-copy satellite RNA HSATII is likely under selection to maintain its pathogen-associated CpG dinucleotides across primates since its origin nearly 10 million years ago and potentially functional L1 inserts maintain atypically high CpG content compared to non-functional copies. We validate the latter with novel co-IPs of high x_{CpG} L1 RNA with both L1 ORF1 and ORF2 proteins, indicating such L1s are more likely to be functional, and with the innate immune sensor *ZCCHC3*⁴⁴. Furthermore, we incorporate structure prediction into our method for the first time, which we validate in both published datasets and new dsRNA detecting antibody assays. In humans, many, but not all, dsRNA forming repeats come from inverted Alus, indicating double-stranded RNA mimicry is likely due to the error prone reverse transcriptional process, rather than being a specific property of the Alus other than their known parasitism of L1. We show a high degree of conservation of double-stranded RNA-forming Alus across primates, indicating selection has maintained their ability to display PAMPs. We find nontrivial potential PAMP forming regions across many genomes which lack either Alus or HSATII, implying reservoirs of potential PAMP formation likely exists within repeats across many organisms, which may have been acted upon in distinct ways in different species. The combination of our analysis within and across species raises the question of whether formation of double-stranded RNA is a function for which aspects of the LINE reverse transcription machinery has been selected for. We generally support the hypothesis that repeats are selected to maintain “non-self” PAMPs, whose induction and subsequent innate sensing may act as sensors for loss of heterochromatin, avoidance of genome instability^{22,23}, or aberrant RNA processing^{22,50,51,53}.

While a species may have evolved to maintain PAMPs, it can be difficult to establish whether PRR signaling is the primary reason for why a PAMP evolved in the first place. Inverted Alus can be hotspots for RNA-editing, altering gene expression over evolutionary times scales, while simultaneously acting a PAMP for PRRs such as MDA5⁵⁷⁻⁵⁹. HSATII may have a DNA regulatory function as well, as its DNA sequences can sequester chromatin regulatory proteins and trigger epigenetic change⁶⁰. Yet in cancer, where Alus and HSATII are often overexpressed, the same features can be sensed as PAMPs^{10,61}. Moreover, such functions are not mutually exclusive. The high CpG presence in L1 may have evolved both to allow active L1 species to remain silenced and to serve as a danger signal when aberrant demethylation occurs. For multicellular organisms with a high degree of epigenetic regulation and chromosomal organization, a repeat species with a non-immune function may be co-opted when it offers an opportunity to maintain stimulatory features to release a danger signal when epigenetic control is lost, such as during the release of repeats after p53 mutations, where immunostimulatory repeats may offer a back-up for p53 functions such as senescence^{6,62}.

Our work has several implications for how to quantify self versus non-self discrimination by the innate immune system. While we focus on motif usage and the formation of long double-stranded RNA structures, our framework is generalizable to other, more complex patterns and machine-learning approaches. Mathematically, our work highlights how approaches from statistical physics, such as maximum entropy and transfer matrix calculations can be used in efficient genome wide calculations and comparisons. The selective forces are intrinsic quantities which can be compared from sequence to sequence. Therefore, they are ideal for evolutionary analysis of genome features, the complexity of which can be added in future models. For instance, Y RNAs, implicated in RIG-I sensing during RNA virus infection⁴, have a more complicated feature set which includes an RNA modification⁶³. The potential association of high CpG content with replication competent L1 may also serve as a marker

382 for STING-cGAS²⁰ activation during reverse transcription or sensing of the ribonucleoprotein complex
383 by TRIM5 α ⁶⁴, and has been implicated here as a co-sensor with *ZCCHC3*⁴⁴. Using such methods to
384 “decipher” noncoding genome regions and to assign them a function may allow such regions to be
385 further exploited therapeutically. The implication is that we can learn a “repeat code” of self-agonists
386 within our genome held by selection to stimulate receptors under specific circumstances. Such work
387 will be enabled by emerging sequencing technologies, such as telomere-to-telomere⁶⁵ sequencing,
388 and broad sequencing of receptor ligand pairs. In doing so, we may discover a new set of phenotypes
389 hiding in the non-coding genome.

390 **References**

- 391
- 392 1 Ting, D. T. *et al.* Aberrant Overexpression of Satellite Repeats in Pancreatic and Other
393 Epithelial Cancers. *Science* **331**, 593-596, (2011).
- 394 2 Chiang, J. J. *et al.* Viral unmasking of cellular 5S rRNA pseudogene transcripts induces
395 RIG-I-mediated immunity. *Nat Immunol* **19**, 53-62, (2018).
- 396 3 Nogalski, M. T. *et al.* A tumor-specific endogenous repetitive element is induced by
397 herpesviruses. *Nat Commun* **10**, 90 (2019).
- 398 4 Vabret, N. *et al.* RNAs are conserved endogenous RIG-I ligands across RNA virus infection
399 and are targeted by HIV-1. *iScience* **25**, 104599 (2022).
- 400 5 Chiappinelli, K. B. *et al.* Inhibiting DNA Methylation Causes an Interferon Response in
401 Cancer via dsRNA Including Endogenous Retroviruses. *Cell* **162**, 974-986, (2015).
- 402 6 Leonova, K. I. *et al.* p53 cooperates with DNA methylation and a suicidal interferon
403 response to maintain epigenetic silencing of repeats and noncoding RNAs. *P Natl Acad Sci*
404 *USA* **110**, E89-E98, (2013).
- 405 7 Mehdi-pour, P. *et al.* Epigenetic therapy induces transcription of inverted SINEs and ADAR1
406 dependency. *Nature* **588**, 169-173, (2020).
- 407 8 Roulois, D. *et al.* DNA-Demethylating Agents Target Colorectal Cancer Cells by Inducing
408 Viral Mimicry by Endogenous Transcripts. *Cell* **162**, 961-973, (2015).
- 409 9 Sheng, W. Q. *et al.* LSD1 Ablation Stimulates Anti-tumor Immunity and Enables Checkpoint
410 Blockade. *Cell* **174**, 549-563, (2018).
- 411 10 Tanne, A. *et al.* Distinguishing the immunostimulatory properties of noncoding RNAs
412 expressed in cancer cells. *P Natl Acad Sci USA* **112**, 15154-15159, (2015).
- 413 11 De Cecco, M. *et al.* L1 drives IFN in senescent cells and promotes age-associated
414 inflammation. *Nature* **566**, 73-78, (2019).
- 415 12 Rice, G. I. *et al.* Reverse-Transcriptase Inhibitors in the Aicardi-Goutieres Syndrome. *New*
416 *Engl J Med* **379**, 2275-2277, (2018).
- 417 13 Vabret, N., Bhardwaj, N. & Greenbaum, B. D. Sequence-Specific Sensing of Nucleic Acids.
418 *Trends Immunol* **38**, 53-65, (2017).
- 419 14 Greenbaum, B. D., Cocco, S., Levine, A. J. & Monasson, R. Quantitative theory of entropic
420 forces acting on constrained nucleotide sequences applied to viruses. *P Natl Acad Sci USA*
421 **111**, 5054-5059, (2014).
- 422 15 Greenbaum, B. D., Levine, A. J., Bhanot, G. & Rabadan, R. Patterns of evolution and host
423 gene mimicry in influenza and other RNA viruses. *Plos Pathog* **4**, e1000079 (2008).
- 424 16 Stern, A. *et al.* The Evolutionary Pathway to Virulence of an RNA Virus. *Cell* **169**, 35-46,
425 (2017).
- 426 17 Takata, M. A. *et al.* CG dinucleotide suppression enables antiviral defence targeting non-self
427 RNA. *Nature* **550**, 124-127, (2017).
- 428 18 Di Gioacchino, A. *et al.* The Heterogeneous Landscape and Early Evolution of Pathogen-
429 Associated CpG Dinucleotides in SARS-CoV-2. *Mol Biol Evol* **38**, 2428-2445, (2021).
- 430 19 Digard, P., Lee, H. M., Sharp, C., Grey, F. & Gaunt, E. Intra-genome variability in the
431 dinucleotide composition of SARS-CoV-2. *Virus Evol* **6**, veaa057, (2020)
- 432 20 Mankan, A. K. *et al.* Cytosolic RNA:DNA hybrids activate the cGAS-STING axis. *EMBO J*
433 **33**, 2937-2946, (2014)
- 434 21 Jensen, S. & Thomsen, A. R. Sensing of RNA Viruses: a Review of Innate Immune
435 Receptors Involved in Recognizing RNA Virus Invasion. *J Virol* **86**, 2900-2910, (2012).
- 436 22 Chen, R., Ishak, C. A. & De Carvalho, D. D. Endogenous Retroelements and the Viral
437 Mimicry Response in Cancer Therapy and Cellular Homeostasis. *Cancer Discov* **11**, 2707-
438 2725, (2021).
- 439 23 Ishak, C. A. & De Carvalho, D. D. Reactivation of Endogenous Retroelements in Cancer
440 Development and Therapy. *Annu Rev Canc Biol* **4**, 159-176, (2020).

441 24 Jaynes, E. T. Information Theory and Statistical Mechanics. *Phys Rev* **106**, 620-630, (1957).
442 25 Goncalves-Carneiro, D., Takata, M. A., Ong, H., Shilton, A. & Bieniasz, P. D. Origin and
443 evolution of the zinc finger antiviral protein. *Plos Pathog* **17**, e1009545 (2021).
444 26 Sved, J. & Bird, A. The Expected Equilibrium of the Cpg Dinucleotide in Vertebrate
445 Genomes under a Mutation Model. *P Natl Acad Sci USA* **87**, 4692-4696, (1990).
446 27 Ullu, E. & Tschudi, C. Alu Sequences Are Processed 7sl Rna Genes. *Nature* **312**, 171-172,
447 (1984).
448 28 Bauer, S. *et al.* Human TLR9 confers responsiveness to bacterial DNA via species-specific
449 CpG motif recognition. *Proc Natl Acad Sci U S A* **98**, 9237-9242, (2001)
450 29 Hubley, R. *et al.* The Dfam database of repetitive DNA families. *Nucleic Acids Res* **44**, D81-
451 D89, (2016).
452 30 Bird, A. P. DNA Methylation and the Frequency of Cpg in Animal DNA. *Nucleic Acids Res* **8**,
453 1499-1504, (1980).
454 31 Shen, J. C., Rideout, W. M. & Jones, P. A. The Rate of Hydrolytic Deamination of 5-
455 Methylcytosine in Double-Stranded DNA. *Nucleic Acids Res* **22**, 972-976, (1994).
456 32 Arndt, P. F. Reconstruction of ancestral nucleotide sequences and estimation of substitution
457 frequencies in a star phylogeny. *Gene* **390**, 75-83, (2007).
458 33 Baele, G., Van de Peer, Y. & Vansteelandt, S. Modelling the ancestral sequence distribution
459 and model frequencies in context-dependent models for primate non-coding sequences.
460 *Bmc Evol Biol* **10**, 244 (2010).
461 34 Berard, J. & Gueguen, L. Accurate Estimation of Substitution Rates with Neighbor-
462 Dependent Models in a Phylogenetic Context. *Syst Biol* **61**, 510-521, (2012).
463 35 Kimura, M. A Simple Method for Estimating Evolutionary Rates of Base Substitutions
464 through Comparative Studies of Nucleotide-Sequences. *J Mol Evol* **16**, 111-120, (1980).
465 36 Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860-
466 921, (2001).
467 37 Penzkofer, T. *et al.* L1Base 2: more retrotransposition-active LINE-1s, more mammalian
468 genomes. *Nucleic Acids Res* **45**, D68-D73, (2017).
469 38 Hata, K. & Sakaki, Y. Identification of critical CpG sites for repression of L1 transcription by
470 DNA methylation. *Gene* **189**, 227-234, (1997).
471 39 Sanchez-Luque, F. J. *et al.* LINE-1 Evasion of Epigenetic Repression in Humans. *Mol Cell*
472 **75**, 590-604 e512 (2019)
473 40 Pichlmair, A. *et al.* RIG-I-mediated antiviral responses to single-stranded RNA bearing 5'-
474 phosphates. *Science* **314**, 997-1001, (2006).
475 41 Solovyov, A. *et al.* Global Cancer Transcriptome Quantifies Repeat Element Polarization
476 between Immunotherapy Responsive and T Cell Suppressive Classes. *Cell Rep* **23**, 512-
477 521, (2018).
478 42 Tubio, J. M. C. *et al.* Extensive transduction of nonrepetitive DNA mediated by L1
479 retrotransposition in cancer genomes. *Science* **345**, 531-531, 1251343n(2014).
480 43 Garcia-Perez, J. L. *et al.* Epigenetic silencing of engineered L1 retrotransposition events in
481 human embryonic carcinoma cells. *Nature* **466**, 769-773, (2010).
482 44 Lian, H. *et al.* ZCCHC3 is a co-sensor of cGAS for dsDNA recognition in innate immune
483 response. *Nat Commun* **9**, ARTN 3349 (2018)
484 45 Taylor, M. S. *et al.* Dissection of affinity captured LINE-1 macromolecular complexes. *Elife*
485 **7**, e30094 (2018).
486 46 Taylor, M. S. *et al.* Affinity Proteomics Reveals Human Host Factors Implicated in Discrete
487 Stages of LINE-1 Retrotransposition. *Cell* **155**, 1034-1048, (2013).
488 48 Wu, B. *et al.* Structural Basis for dsRNA Recognition, Filament Formation, and Antiviral
489 Signal Activation by MDA5. *Cell* **152**, 276-289, (2013).
490 48 Novikova, I. V., Hennelly, S.P. & Sanbonmatsu, K.Y. Sizing up long non-coding RNAs: do
491 lncRNAs have secondary and tertiary structure? *Bioarchitecture* **2**, 189-199 (2012).

492 49 Ardeljan, D. *et al.* LINE-1 ORF2p expression is nearly imperceptible in human cancers.
493 *Mobile DNA-Uk* **11**, 1 (2019).

494 50 Arrowsmith C, W. Q., Nie D, Alawi W, et al. Lupien M. Altered RNA splicing initiates the viral
495 mimicry response from inverted SINEs following type I PRMT inhibition in Triple-Negative
496 Breast Cancer. *Research Square* PPR368631, (2021).

497 51 Bowling, E. A. *et al.* Spliceosome-targeted therapies trigger an antiviral immune response in
498 triple-negative breast cancer. *Cell* **184**, 384-403, (2021).

499 52 Ahmad, S. *et al.* Breaching Self-Tolerance to Alu Duplex RNA Underlies MDA5-Mediated
500 Inflammation. *Cell* **172**, 797-810, (2018).

501 53 Seiler, M. *et al.* H3B-8800, an orally available small-molecule splicing modulator, induces
502 lethality in spliceosome-mutant cancers. *Nat Med* **24**, 497-504, (2018).

503 54 Lizio, M. *et al.* Update of the FANTOM web resource: expansion to provide additional
504 transcriptome atlases. *Nucleic Acids Res* **47**, D752-D758, (2019).

505 55 Lizio, M. *et al.* Gateways to the FANTOM5 promoter level mammalian expression atlas.
506 *Genome Biol* **16**, 22 (2015).

507 56 Tinh, V. N. *et al.* Mitochondrial evidence for multiple radiations in the evolutionary history of
508 small apes. *Bmc Evol Biol* **10**, 74 (2010).

509 57 Capshew, C. R., Dusenbury, K. L. & Hundley, H. A. Inverted Alu dsRNA structures do not
510 affect localization but can alter translation efficiency of human mRNAs independent of RNA
511 editing. *Nucleic Acids Res* **40**, 8637-8645 (2012)

512 58 Daniel, C., Silberberg, G., Behm, M. & Ohman, M. Alu elements shape the primate
513 transcriptome by cis-regulation of RNA editing. *Genome Biol* **15**, R28, (2014)

514 59 Kim, D. D. *et al.* Widespread RNA editing of embedded alu elements in the human
515 transcriptome. *Genome Res* **14**, 1719-1725, (2004)

516 60 Hall, L. L. *et al.* Demethylated HSATII DNA and HSATII RNA Foci Sequester PRC1 and
517 MeCP2 into Cancer-Specific Nuclear Bodies. *Cell Rep* **18**, 2943-2956, (2017).

518 61 Porter, R. L. *et al.* Satellite repeat RNA expression in epithelial ovarian cancer associates
519 with a tumor-immunosuppressive phenotype. *J Clin Invest* **132**, 155931(2022)

520 62 Levine, A. J. & Greenbaum, B. The Maintenance of Epigenetic States by p53: The Guardian
521 of the Epigenome. *Oncotarget* **3**, 1503-1504 (2012).

522 63 Hornung, V. *et al.* 5'-Triphosphate RNA is the ligand for RIG-I. *Science* **314**, 994-997, (2006)

523 64 Volkmann, B. *et al.* Human TRIM5 α senses and restricts LINE-1 elements. *P Natl Acad Sci*
524 *USA* **117**, 17965017976 (2020).

525 65 Hoyt, S. J. *et al.* From telomere to telomere: The transcriptional and epigenetic state of human
526 repeat elements. *Science* **376**, 57, (2022).

527 **Data and Code Availability Statement:** Original data will be made public upon acceptance. Code
528 will likewise be deposited on GitHub.

529

530 **Acknowledgements:**

531 This research was funded in part through the NIH/NCI Cancer Center Support Grant P30 CA008748
532 (A.S., S.S., B.G.); NIH grants R01AI081848 (N.V., B.G.), R01CA240924 (A.S., B.G.), R01GM126170
533 (J.L.), R01AG078925 (J.L.), P50 254838-01 (O.A.-W.) and U01CA228963 (A.S., S.S., B.G.);
534 Fondation de la Recherche Médicale: ANR-Flash Covid, Project SARS-Cov-2immunRNAs (S.C.,
535 R.M.); the V Foundation for Cancer Research (A.S.); the Mark Foundation ASPIRE award (B.G.); the
536 Pershing Square Sohn Prize-Mark Foundation Fellowship (A.S., O.A.-W., N.V., B.G.); the Edward P.
537 Evans Foundation (O.A.-W.); the Canadian Institute of Health Research (CIHR), New Investigator
538 salary award 201512MSH360794-228629 (D.D.C.); Canada Research Chair (D.D.C.); CIHR
539 Foundation Grant FDN 148430 (D.D.C.); CIHR Project Grant PJT 165986 (D.D.C.); NSERC 489073
540 (D.D.C.); and the European Union's Horizon 2020 research and innovation programme under the
541 Marie Skłodowska-Curie grant agreement No 101026293 (A.D.G.). The authors would like to
542 acknowledge productive conversations with Jef Boeke, Kathleen Burns, Katherine Chiappinelli,
543 Arnold Levine, Phil Sharp, Martin Taylor, David T. Ting, and the De Carvalho, Cocco, Greenbaum
544 and Monasson laboratories; and thank Nicole Rusk for reading and editing the manuscript. We would
545 also like to acknowledge support from the National Center for Dynamic Interactome Research (NIH
546 P41GM10982); the Genome Technology Center at NYULH, a shared resource partially supported by
547 the Cancer Center Support Grant P30CA016087 at the Laura and Isaac Perlmutter Cancer Center;
548 and the UMCG Research Sequencing Facility and Utrecht Sequencing Facility (USEQ; USEQ is
549 subsidized by the University Medical Center Utrecht and The Netherlands X-omics Initiative [NWO
550 project 184.034.019]).

551

552 **Author Contributions:**

553 Conceptualization: P.S., R.M., S.C., B.G.; Research Plan: P.S., R.M., S.C., B.G.; Mathematical
554 Modeling: P.S., A.D.G.; R.M., S.C., B.G.; Double-stranded Force Calculation: P.S., A.D.G.; R.M.,
555 S.C.; Model Implementation: P.S., A.D.G.; Comparative Genomic Analysis: A.S.; Data Analysis: P.S.,
556 A.D.G., A.S., H.L., S.M., S.S.; L1 experimental design and execution: J.L., H.J., B.H.L.; Double-
557 stranded RNA experimental design: D.D.C.; Interpretation: P.S., N.V., J.L., O.A.-W., D.D.C., R.M.,
558 S.C., B.G.; Writing: P.S., A.D.G.; R.M., S.C., B.G.; Reviewing & Editing: P.S., A.D.G., J.L., O.A.-W.,
559 D.D.C., R.M., S.C., B.G..

560

561 **Declaration of Interests:**

562 O.A.-W. has performed consulting for Incyte, Prelude Therapeutics, AstraZeneca, Merck, Janssen,
563 Pfizer Boulder, and LoxoOncology/Eli Lilly and is on the Scientific Advisory Board of AIChem and
564 Harmonic Discovery Inc. B.G. has received honoraria for speaking engagements from Merck, Bristol
565 Meyers Squibb, and Chugai Pharmaceuticals; has received research funding from Bristol Meyers
566 Squibb, Merck, and ROME Therapeutics; and has been a compensated consultant for Darwin Health,
567 Merck, PMV Pharma, Shennon Biotechnologies, and Rome Therapeutics of which he is a co-founder.
568 A.S. has been a compensated consultant for Rome Therapeutics. D.D.C. received research funding
569 from Pfizer and Nektar therapeutics; is a shareholder, co-founder and CSO of Adela (former DNAMx).
570 J.L. received research funding from: ROME Therapeutics, Ribon Therapeutics, and Refeyn; he
571 received compensation from Transposon Therapeutics, ROME Therapeutics and Oncolinea.

572

573 **Correspondence:** psulc@asu.edu, j.p.lacava@umcg.nl, Daniel.DeCarvahlo@uhnresearch.ca,
574 simona.cocco@phys.ens.fr, greenbab@mskcc.org

575 **FIGURE LEGENDS**

576
577 **Figure 1 | Competition between selective and entropic forces define presence of pathogen**
578 **associated patterns in the genome. A**, Representation of selective versus entropic forces on a
579 PAMP. Random sequences reproducing the nucleotide frequencies generically have numbers of
580 occurrences of a PAMP significantly different from what is observed in an actual genomic sequence.
581 Deviations imply constraints to enrich or avoid a PAMP, which we characterized by a, respectively,
582 positive or negative selective force. In our Maximum Entropy framework, the selective force
583 counterbalances the entropic force resulting from the loss of diversity (entropy) in sequences having
584 statistically abnormal PAMP numbers. **B**, Comparison of all nucleotide biases between repeat and
585 viral families. In this histogram each repeat family is compared with reference viral genomes of viral
586 families by computing the symmetrized Kullback-Leibler divergence of the probability distributions
587 associated to models inferred from viral and repeat sequences. Repeats with particularly high
588 divergence values are indicated by arrows. The value "Human genome (all)" is obtained with a model
589 trained on sequences randomly sampled from the human genome (hg38), while for "Human genome
590 (CDS)" we only consider coding sequences. **C**, Detailed comparison of nucleotide biases between
591 selected repeat and viral families. Each point is the symmetrized Kullback-Leibler divergence
592 between a repeat family (x-axis) and a specific viral family (as indicated by the color). Bars represent,
593 for each repeat family under consideration, the average value over the viral families.

594
595 **Figure 2 | Forces on CpG dinucleotides in the human genome. A**, Change in x_{CpG} computed on
596 all inserts annotated in hg38 for each repeat family, versus x_{CpG} of the consensus repeat reported in
597 the DFAM database. Alus and HSATII are highlighted as exceptions to the general trend. **B**, The
598 mean x_{CpG} of all inserts in a repeat family as a function of the Kimura distance from the consensus
599 sequence for each family. **C**, Annotation of high- x_{CpG} ($x_{CpG} > 0$) sequences in the human genome
600 according to their overlap with annotated repeats in the DFAM database. The + or - sign after the
601 repeat name indicates the sense in which the repeat is annotated in the database. "Unannotated"
602 sequences do not overlap with any repeat in the database. **D**, Scatter plot of x_{CpG} and x_{UpA} for LINE-
603 1 functional (blue) and non-functional (green) elements in the human genome. The white ellipse
604 corresponds to one standard deviation distance from the mean for x_{CpG} and x_{UpA} forces on FLI and
605 FLnI LINE-1 inserts respectively. **E**, x_{CpG} for FLnI inserts of LINE-1 and HSAT-II in human genome
606 as a function of average distance from the intact FLI sequences (for LINE-1) or the distance from the
607 consensus sequence (for HSAT-II). The force relaxation evolutionary model fit is shown for both
608 sequence families together with a Kimura (null) model fit. **F**, Distribution of x_{CpG} of L1 ORF2p binding
609 L1 transcripts in embryonal carcinoma cell line (N2102Ep). Functional intact LINEs are colored in
610 blue (BH corrected p-value labeled for t-test, **** denotes adjusted p-value < 0.01). ORF2p enriched
611 and depleted transcripts are selected by differential expression analysis between ORF2p-IP versus
612 Mock/total with $|\log_2FC|$ greater than 3 and adjusted p-value < 0.05 for Fisher Exact test on proportion
613 of x_{CpG} high versus x_{CpG} low of ORF2p enriched and depleted transcripts. **G**, x_{CpG} on *ZCCHC3*
614 binding LINE transcripts in N2102Ep. Functional intact LINEs colored in blue (BH corrected p-value
615 labeled for t-test, **** denotes adjusted p-value < 0.01). *ZCCHC3* enriched and depleted transcripts
616 selected by differential expression analysis between *ZCCHC3*-IP versus Mock/total with a $|\log_2FC|$
617 greater than 3 and adjusted p-value < 0.05.

618
619 **Figure 3 | Double-stranded forces in the human genome. A**, Histogram of x_{ds} calculated for mRNA
620 coding sequences, non-coding RNAs, inserts, consensus sequences of repeats, and sequences
621 obtained by randomly reshuffling mRNA coding sequences (yellow). **B**, Mean of x_{ds} calculated for
622 each family of repeats as a function of the mean Kimura distance of all inserts in a repeat family from
623 their consensus sequence. The solid line corresponds to mean value (and standard deviation from it)
624 for all families binned into the same distance from consensus. **C**, Complexity of sequences in

625 complementary regions found in the human genome as a function of segment length. Complementary
626 regions that overlap with known repeat element or ncRNA or mRNA are highlighted as gray dots, with
627 different contour colors depending on the specific family they overlap with. Dashed lines correspond
628 to the complexity of a completely random sequence (top line) and trivial region consisting of a single
629 nucleotide (bottom). Complexity of both complementary segments are similar, so we only include the
630 complexity of one of each complementary transcript. **D**, x_{ds} histograms in human genome (sliding
631 window with transcript of length of 3 kb) compared to MDA5 binding RNA transcripts. Enriched
632 transcripts have a positive log-enrichment with respect to the control experiment. Inverted repeat (IR)
633 transcripts are annotated repeats with another repeat of the same family in opposite genomic sense
634 within 3 kb. **E**, Similar to panel (**D**), for J2 binding transcripts. **F**, Type of repeat (as annotated in
635 RepeatMasker) with the longest overlapping sequence in complementary sequences for high- x_{ds} (x_{ds}
636 > 0.5) windows in hg38 (left), the MDA5 binding experiments (middle) and the J2 binding experiment
637 (right). Sequences are accounted as "IR" (Inverted Repeats) if the two complementary regions
638 overlap with repeats annotated in the database with the same name but inverted sense (+/- or -/+).
639 "Non-IR" indicates cases where the two repeats overlapping with the two complementary regions
640 have a different name. "Unannotated" indicates cases where one or both the two complementary
641 regions do not overlap with any repeat in the database.

642
643 **Figure 4 | Evolution and conservation of forces on PAMPs.** **A**, Complexity of sequences in
644 complementary regions found in the *Danio rerio* genome as a function of segment length. Dashed
645 lines correspond to the complexity of a completely random sequence (top line) and trivial region
646 consisting of a single nucleotide (bottom). **B**, Scatter plot of the overlap coefficient between the high-
647 x_{CpG} ($x_{CpG} > 0$) sequences in the human genome and those of other primates versus the most recent
648 common ancestor (MRCA) time⁵⁶. Two high- x_{CpG} sequences are considered overlapping if they result
649 as a hit from BLAST (Methods). The blue curve denotes an exponential fit. **C**, Barplot presenting
650 overlap with repeats of conserved high- x_{CpG} sequences. The x-axis indicates the MRCA time (0 Mya
651 are human sequences). Sequences are accounted as repeats if they overlap with annotations in the
652 DFAM database. The + or - sign after the repeat name indicates the sense in which the repeat is
653 annotated in the database. "Unannotated" sequences do not overlap with any repeat in the database.
654 **D**, Same analysis as (**C**), but with high- x_{ds} sequences ($x_{ds} > 0.5$). **E**, Barplot presenting overlap with
655 repeats of conserved high- x_{ds} sequences. The x-axis indicates the MRCA time (0 Mya are human
656 sequences). Sequences are accounted as repeats if they overlap with annotations in the
657 RepeatMasker database. Sequences are accounted as "IR" (Inverted Repeats) if the two
658 complementary regions overlap with two annotations in the RepeatMasker database with the same
659 name but inverted sense (+/- or -/+). Sequences are indicated as "Non-IR" if the two repeats
660 overlapping with the two complementary regions have a different name. "Unannotated" indicates
661 cases where one or both the two complementary regions do not overlap with any known repeat.

662 **Extended Data Figure 1 | A**, x_{CpG} versus Kimura distance from consensus sequence for each Alu
663 family. Solid lines indicate binned means and standard deviations. **B**, Genomic distribution of high-
664 x_{CpG} ($x_{\text{CpG}} > 0$) regions in the human genome (center), compared with the distribution of the full
665 genome (left bar). In the right bar we show the genomic distribution of high- x_{CpG} regions that do not
666 overlap with any repeat in the DFAM database. **C**, x_{CpG} on L1 ORF1p binding LINE transcripts in
667 N2102Ep. Functional intact LINES are colored in blue. BH corrected p-value is labeled for t-test. ****
668 denote adjusted p-value < 0.01 . L1 ORF1p enriched and depleted transcripts are selected by
669 differential expression analysis between L1 ORF1p-IP vs Mock/total with a $|\log_2\text{FC}|$ greater than 3
670 and adjusted p-value < 0.05 .

671
672 **Extended Data Figure 2 | A**, The mean of maximum lengths in a secondary structure in a single-
673 stranded RNA sequence (green line), and the mean maximum length of complementary segments
674 (blue line), along with respective fits (Methods). **B**, x_{ds} on repeat family Tigger4a. The force relaxation
675 evolutionary model fit shows the relaxation of the inserts compared to the relaxation simulated by
676 neutral Kimura model. **C**, x_{ds} histograms in human genome (sliding window with transcript of length
677 of 3 kb) compared to MDA5 binding RNA transcripts as experimentally found in⁵². Enriched transcripts
678 have a positive log-enrichment with respect to the control experiment. Inverted repeat (IR) transcripts
679 are annotated repeats with another repeat of the same family in opposite genomic sense within 3 kb.
680 **D**, Correlation between log-enrichment of reads aligning to each complementary sequence in MDA5-
681 binding experiment, and x_{ds} . The blue line shows the fit of a generalized additive model. **E**, relation
682 between log-enrichment of reads aligning to each complementary sequence in J2-binding
683 experiment, and x_{ds} . The blue line shows the fit of a generalized additive model. **F**, Type of repeat
684 with the longest overlapping sequence in complementary sequences with high MDA5 signal and high
685 x_{ds} ($x_{\text{ds}} > 0.5$) and complementary sequences with low MDA5 signal and low x_{ds} . **G**, Genomic
686 distribution of high- x_{ds} regions in the human genome (right bar), compared with the distribution of the
687 full genome (left bar). **H**, Type of repeat with the longest overlapping sequence in complementary
688 sequences with high J2 signal and high x_{ds} and complementary sequences with low J2 signal and
689 low x_{ds} .

690
691 **Extended Data Figure 3 |** Volcano plot of repeat element expression of elements with double
692 stranded force greater than 0.5 in H3B-8800 versus DMSO treated SF3B1-K700 mutant K562 cell
693 lines.

694
695 **Extended Data Figure 4 |** Distributions of x_{CpG} in several organism genomes (sliding window with
696 transcript of length of 3000). For human and mouse, we also show, in orange, the profile of the
697 histogram of x_{CpG} after excluding reads annotated as CpG islands or enhancers.

698
699 **Extended Data Figure 5 |** Distributions of x_{ds} in several organism genomes (sliding window with
700 transcript of length of 3000).

701
702 **Extended Data Figure 6 | A**, Standard deviation versus mean of x_{CpG} computed for each 3000-base
703 windows for each organism analyzed. Orange denotes points computed from primate genomes. **B**,
704 Skewness versus mean of x_{ds} computed for each 3000-base windows for each organism analyzed.
705 Orange denotes points computed from primate genomes.

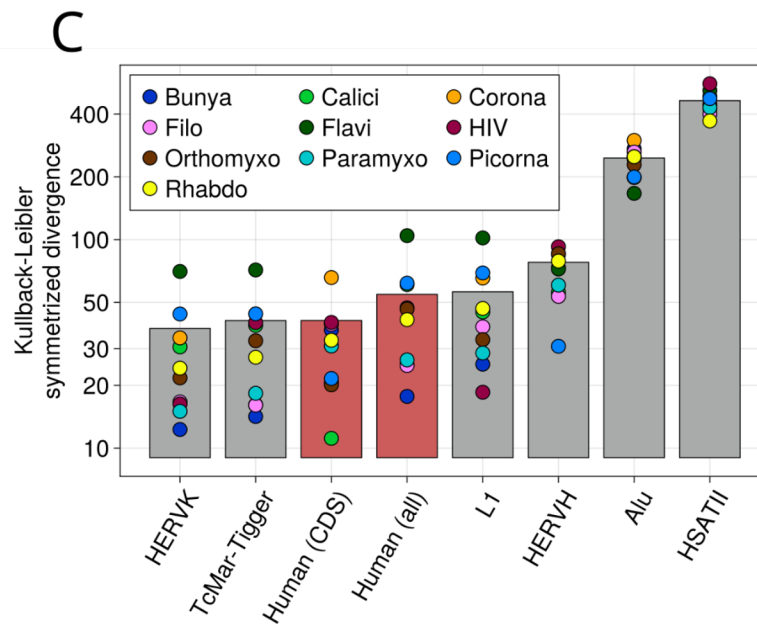
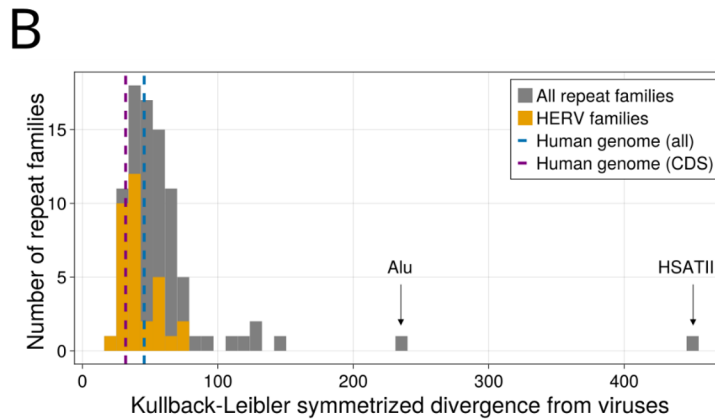
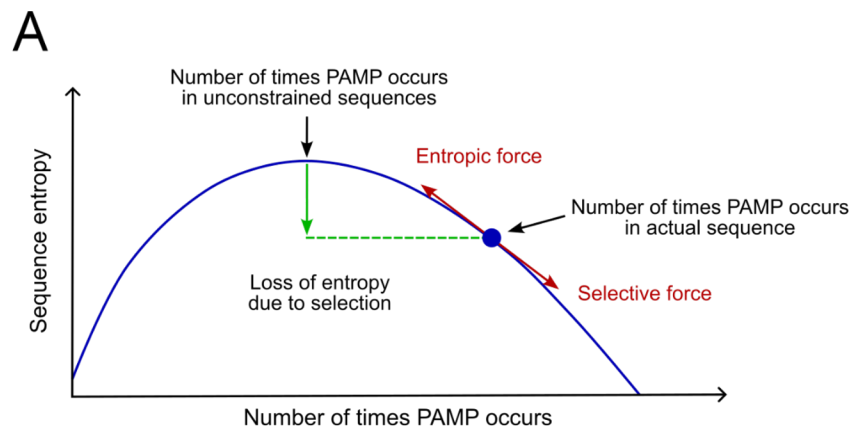


Figure 1 | Competition between selective and entropic forces define presence of pathogen associated patterns in the genome. A, Representation of selective versus entropic forces on a PAMP. Random sequences reproducing the nucleotide frequencies generically have numbers of occurrences of a PAMP significantly different from what is observed in an actual genomic sequence. Deviations imply constraints to enrich or avoid a PAMP, which we characterized by a, respectively, positive or negative selective force. In our Maximum Entropy framework, the selective force counterbalances the entropic force resulting from the loss of diversity (entropy) in sequences having statistically abnormal PAMP numbers. **B,** Comparison of all nucleotide biases between repeat and viral families. In this histogram each repeat family is compared with reference viral genomes of viral families by computing the symmetrized Kullback-Leibler divergence of the probability distributions associated to models inferred from viral and repeat sequences. Repeats with particularly high divergence values are indicated by arrows. The value "Human genome (all)" is obtained with a model trained on sequences randomly sampled from the human genome (hg38), while for "Human genome (CDS)" we only consider coding sequences. **C,** Detailed comparison of nucleotide biases between selected repeat and viral families. Each point is the symmetrized Kullback-Leibler divergence between a repeat family (x-axis) and a specific viral family (as indicated by the color). Bars represent, for each repeat family under consideration, the average value over the viral families.

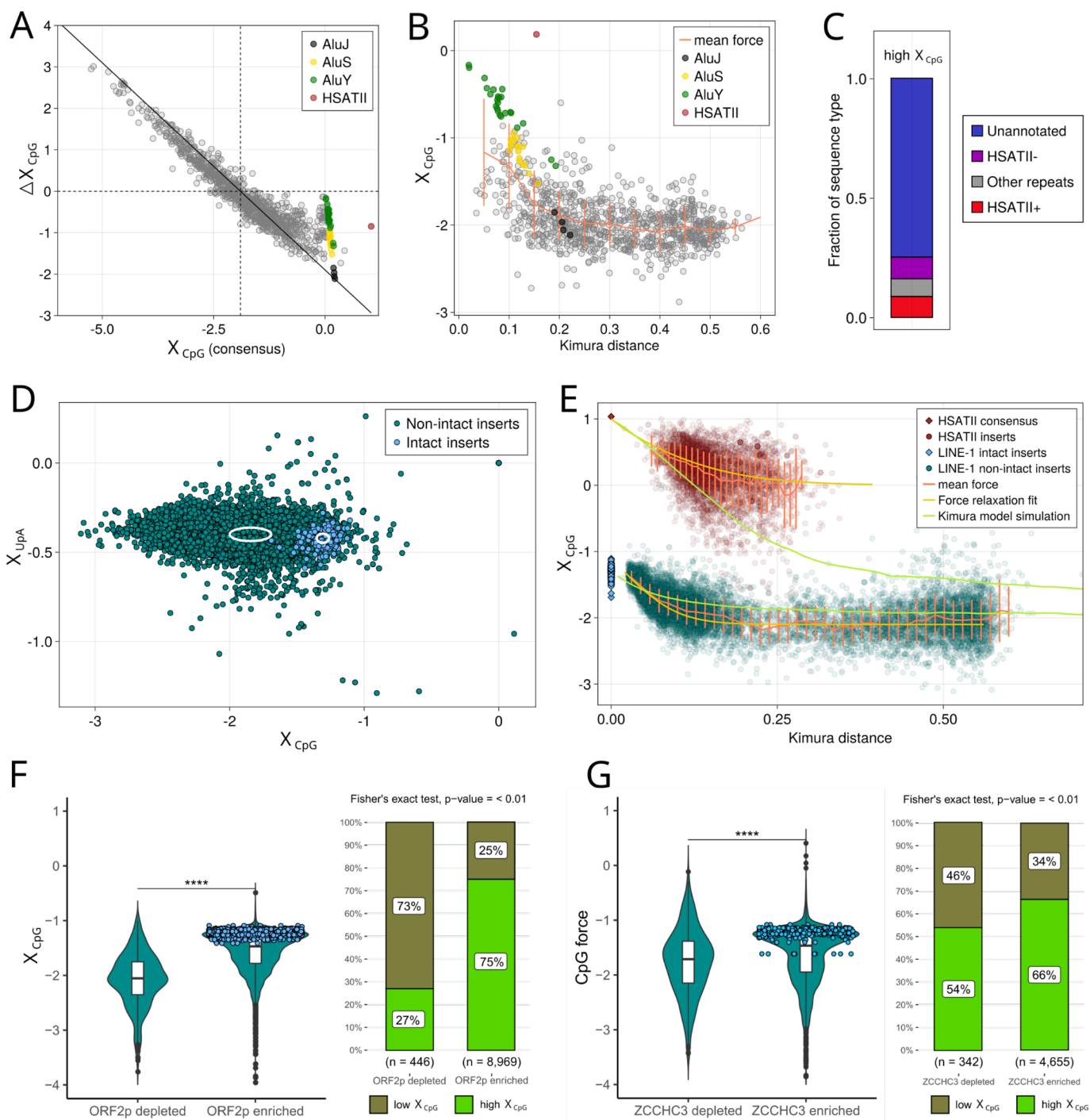


Figure 2 | Forces on CpG dinucleotides in the human genome. **A**, Change in x_{CpG} computed on all inserts annotated in hg38 for each repeat family, versus x_{CpG} of the consensus repeat reported in the DFAM database. Alus and HSATII are highlighted as exceptions to the general trend. **B**, The mean x_{CpG} of all inserts in a repeat family as a function of the Kimura distance from the consensus sequence for each family. **C**, Annotation of high- x_{CpG} ($x_{\text{CpG}} > 0$) sequences in the human genome according to their overlap with annotated repeats in the DFAM database. The + or - sign after the repeat name indicates the sense in which the repeat is annotated in the database. "Unannotated" sequences do not overlap with any repeat in the database. **D**, Scatter plot of x_{CpG} and x_{UpA} for LINE-1 functional (blue) and non-functional (green) elements in the human genome. The white ellipse corresponds to one standard deviation distance from the mean for x_{CpG} and x_{UpA} forces on FLI and FLnI LINE-1 inserts respectively. **E**, x_{CpG} for FLnI inserts of LINE-1 and HSAT-II in human genome as a function of average distance from the intact FLI sequences (for LINE-1) or the distance from the consensus sequence (for HSAT-II). The force relaxation evolutionary model fit is shown for both sequence families together with a Kimura (null) model fit. **F**, Distribution of x_{CpG} of L1 ORF2p binding L1 transcripts in embryonal carcinoma cell line (N2102Ep). Functional intact LINES are colored in blue (BH corrected p-value labeled for t-test, **** denotes adjusted p-value < 0.01). ORF2p enriched and depleted transcripts are selected by differential expression analysis between ORF2p-IP versus Mock/total with $|\log_2\text{FC}|$ greater than 3 and adjusted p-value < 0.05 for Fisher Exact test on proportion of x_{CpG} high versus x_{CpG} low of ORF2p enriched and depleted transcripts. **G**, x_{CpG} on ZCCHC3 binding LINE transcripts in N2102Ep. Functional intact LINES colored in blue (BH corrected p-value labeled for t-test, **** denotes adjusted p-value < 0.01). ZCCHC3 enriched and depleted transcripts selected by differential expression analysis between ZCCHC3-IP versus Mock/total with a $|\log_2\text{FC}|$ greater than 3 and adjusted p-value < 0.05 .

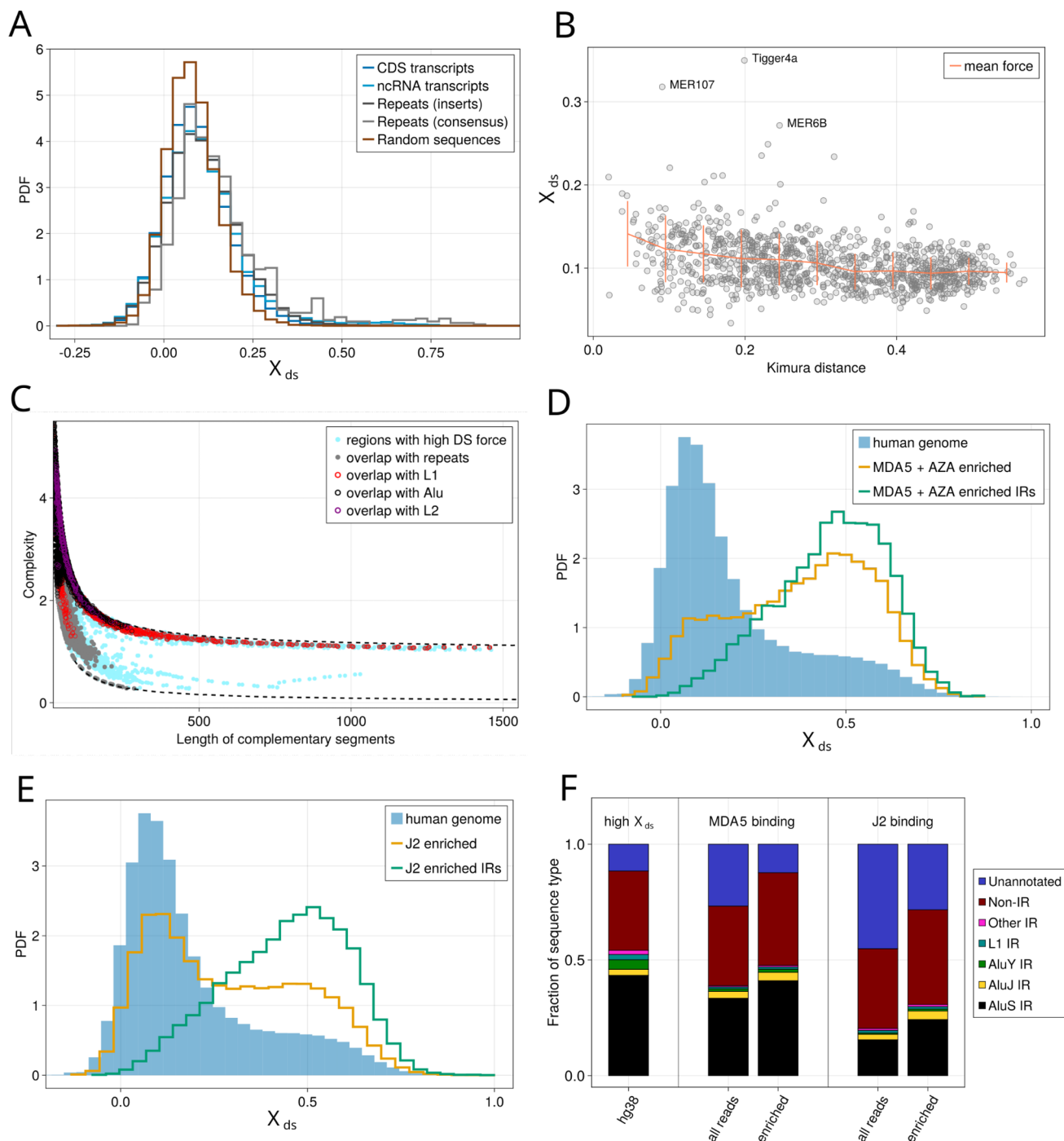


Figure 3 | Double-stranded forces in the human genome. **A**, Histogram of x_{ds} calculated for mRNA coding sequences, non-coding RNAs, inserts, consensus sequences of repeats, and sequences obtained by randomly reshuffling mRNA coding sequences (yellow). **B**, Mean of x_{ds} calculated for each family of repeats as a function of the mean Kimura distance of all inserts in a repeat family from their consensus sequence. The solid line corresponds to mean value (and standard deviation from it) for all families binned into the same distance from consensus. **C**, Complexity of sequences in complementary regions found in the human genome as a function of segment length. Complementary regions that overlap with known repeat element or ncRNA or mRNA are highlighted as gray dots, with different contour colors depending on the specific family they overlap with. Dashed lines correspond to the complexity of a completely random sequence (top line) and trivial region consisting of a single nucleotide (bottom). Complexity of both complementary segments are similar, so we only include the complexity of one of each complementary transcript. **D**, x_{ds} histograms in human genome (sliding window with transcript of length of 3 kb) compared to MDA5 binding RNA transcripts. Enriched transcripts have a positive log-enrichment with respect to the control experiment. Inverted repeat (IR) transcripts are annotated repeats with another repeat of the same family in opposite genomic sense within 3 kb. **E**, Similar to panel (D), for J2 binding transcripts. **F**, Type of repeat (as annotated in RepeatMasker) with the longest overlapping sequence in complementary sequences for high- x_{ds} ($x_{ds} > 0.5$) windows in hg38 (left), the MDA5 binding experiments (middle) and the J2 binding experiment (right). Sequences are accounted as "IR" (Inverted Repeats) if the two complementary regions overlap with repeats annotated in the database with the same name but inverted sense (+/- or -/+). "Non-IR" indicates cases where the two repeats overlapping with the two complementary regions have a different name. "Unannotated" indicates cases where one or both the two complementary regions do not overlap with any repeat in the database.

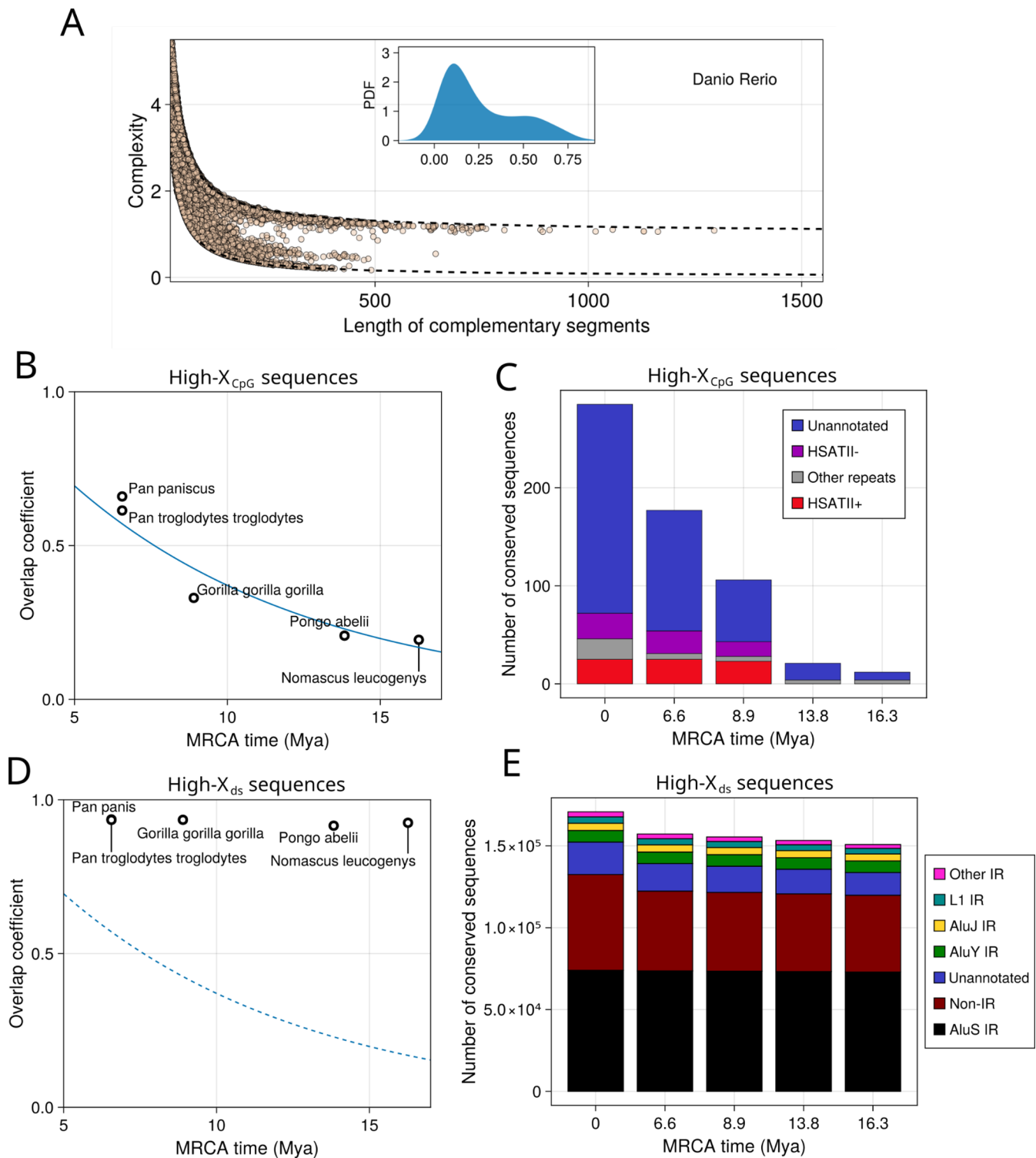
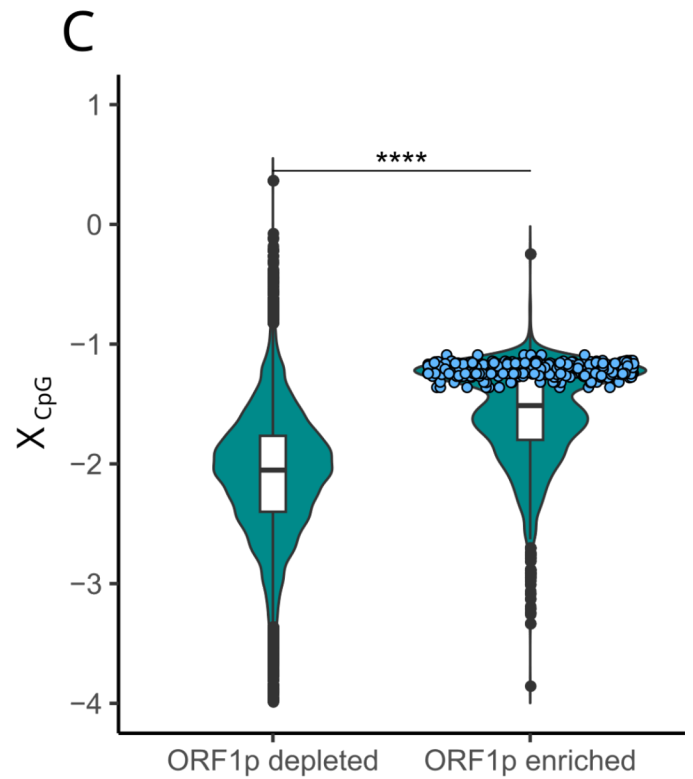
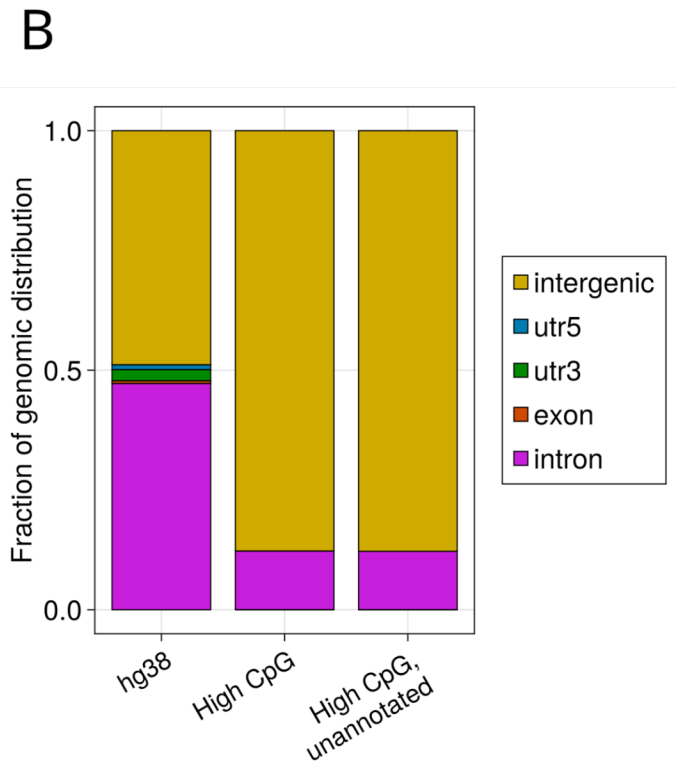
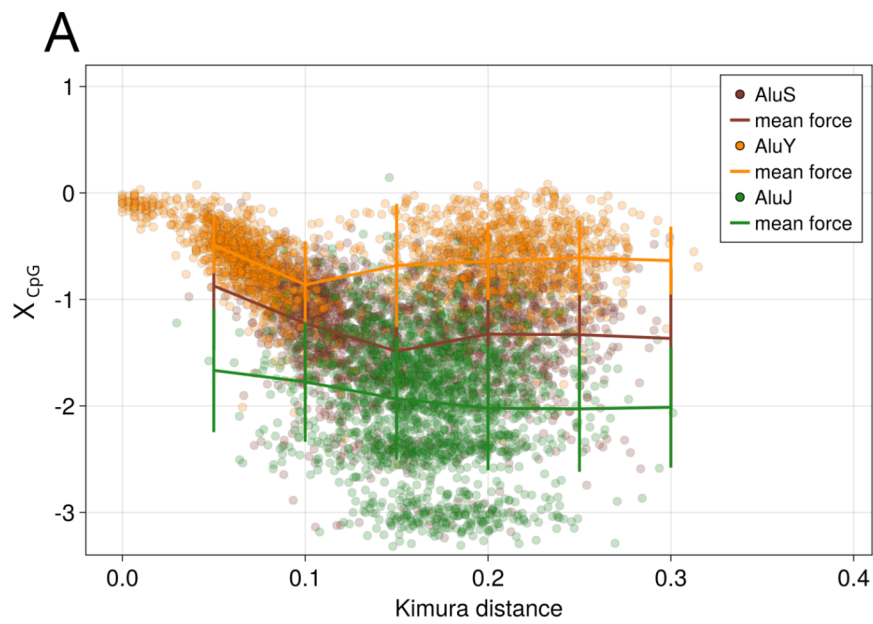
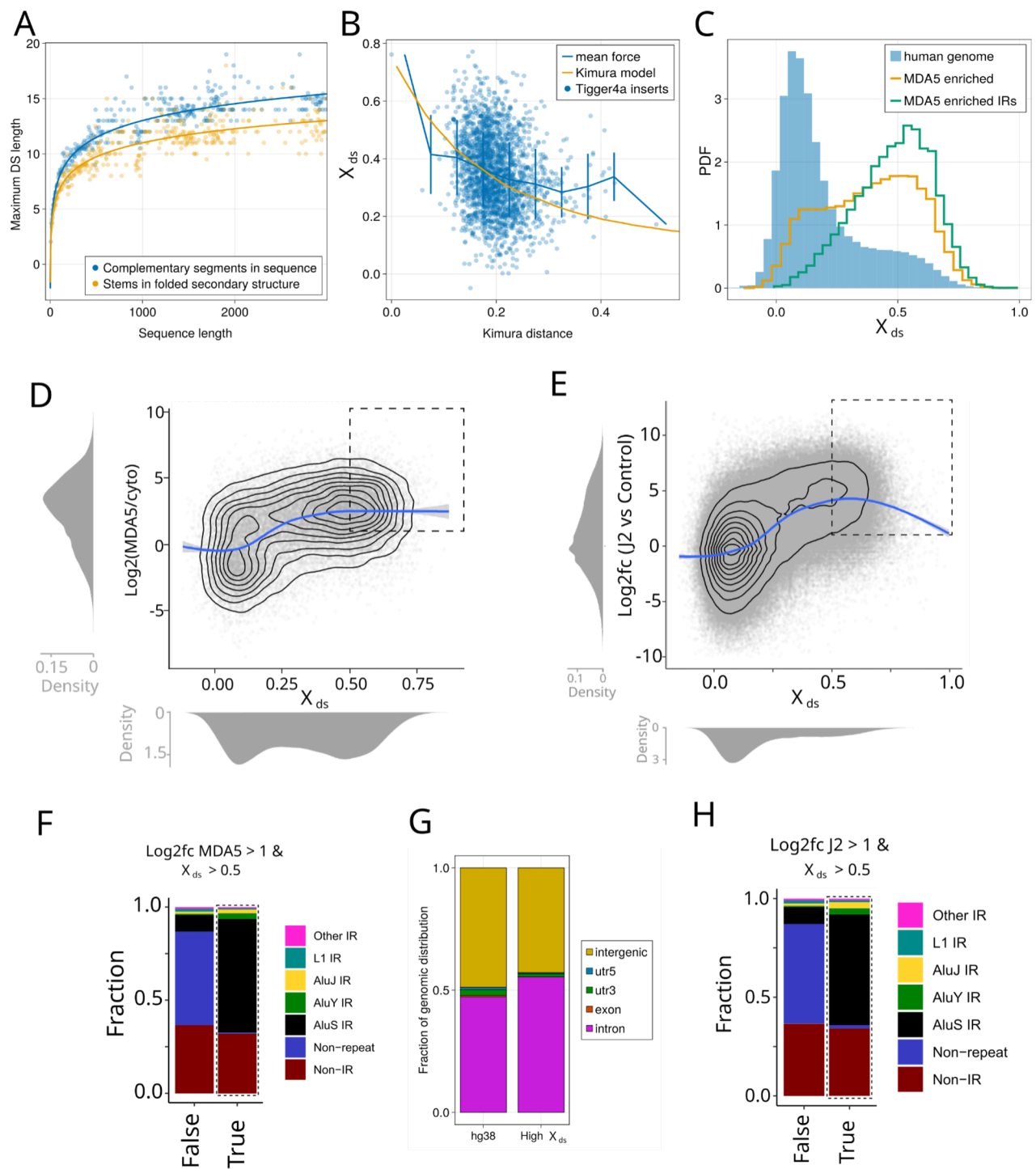


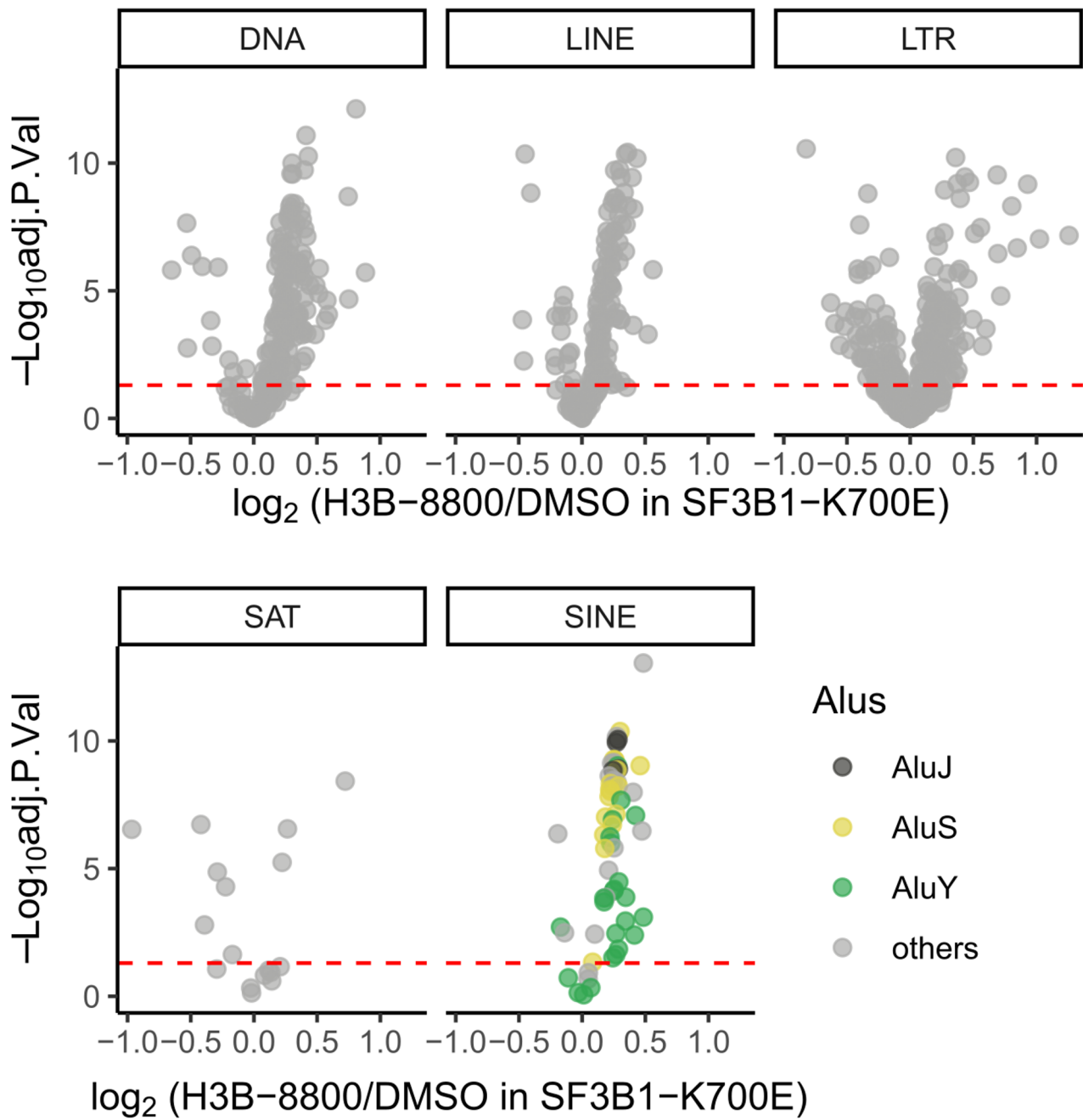
Figure 4 | Evolution and conservation of forces on PAMPs. **A**, Complexity of sequences in complementary regions found in the *Danio rerio* genome as a function of segment length. Dashed lines correspond to the complexity of a completely random sequence (top line) and trivial region consisting of a single nucleotide (bottom). **B**, Scatter plot of the overlap coefficient between the high- x_{CpG} ($x_{\text{CpG}} > 0$) sequences in the human genome and those of other primates versus the most recent common ancestor (MRCA) time⁵⁶. Two high- x_{CpG} sequences are considered overlapping if they result as a hit from BLAST (Methods). The blue curve denotes an exponential fit. **C**, Barplot presenting overlap with repeats of conserved high- sequences. The x-axis indicates the MRCA time (0 Mya are human sequences). Sequences are accounted as repeats if they overlap with annotations in the DFAM database. The + or - sign after the repeat name indicates the sense in which the repeat is annotated in the database. "Unannotated" sequences do not overlap with any repeat in the database. **D**, Same analysis as (C), but with high- x_{ds} sequences ($x_{\text{ds}} > 0.5$). **E**, Barplot presenting overlap with repeats of conserved high- x_{ds} sequences. The x-axis indicates the MRCA time (0 Mya are human sequences). Sequences are accounted as repeats if they overlap with annotations in the RepeatMasker database. Sequences are accounted as "IR" (Inverted Repeats) if the two complementary regions overlap with two annotations in the RepeatMasker database with the same name but inverted sense (+/- or -/+). Sequences are indicated as "Non-IR" if the two repeats overlapping with the two complementary regions have a different name. "Unannotated" indicates cases where one or both the two complementary regions do not overlap with any known repeat.



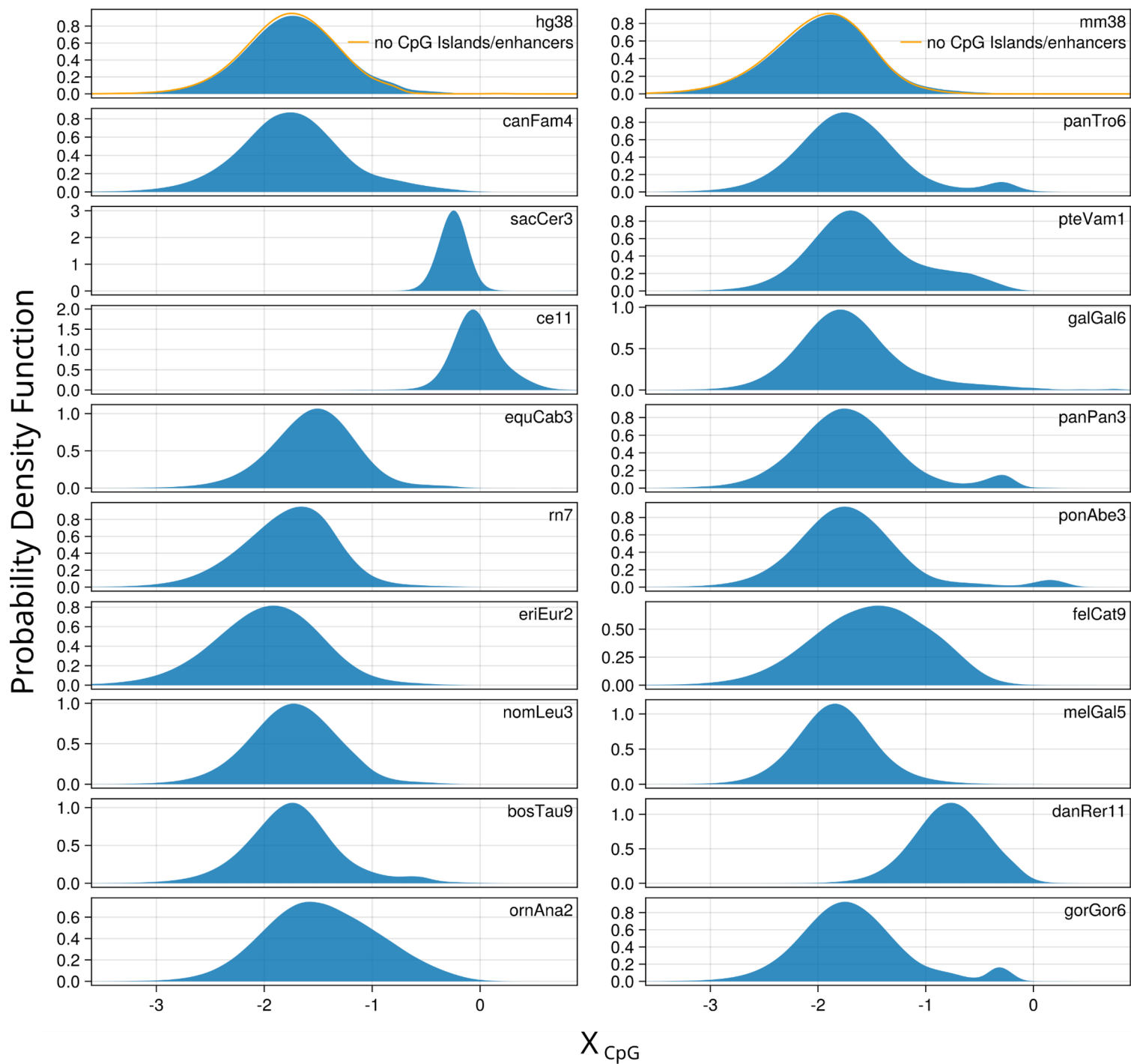
Extended Data Figure 1 | A, x_{CpG} versus Kimura distance from consensus sequence for each Alu family. Solid lines indicate binned means and standard deviations. **B**, Genomic distribution of high- x_{CpG} ($x_{CpG} > 0$) regions in the human genome (center), compared with the distribution of the full genome (left bar). In the right bar we show the genomic distribution of high- x_{CpG} regions that do not overlap with any repeat in the DFAM database. **C**, x_{CpG} on L1 ORF1p binding LINE transcripts in N2102Ep. Functional intact LINEs are colored in blue. BH corrected p-value is labeled for t-test. **** denote adjusted p-value < 0.01 . L1 ORF1p enriched and depleted transcripts are selected by differential expression analysis between L1 ORF1p-IP vs Mock/total with a $|\log_2FC|$ greater than 3 and adjusted p-value < 0.05 .



Extended Data Figure 2 | A, The mean of maximum lengths in a secondary structure in a single-stranded RNA sequence (green line), and the mean maximum length of complementary segments (blue line), along with respective fits (Methods). **B**, x_{ds} on repeat family Tigger4a. The force relaxation evolutionary model fit shows the relaxation of the inserts compared to the relaxation simulated by neutral Kimura model. **C**, x_{ds} histograms in human genome (sliding window with transcript of length of 3 kb) compared to MDA5 binding RNA transcripts as experimentally found in⁵². Enriched transcripts have a positive log-enrichment with respect to the control experiment. Inverted repeat (IR) transcripts are annotated repeats with another repeat of the same family in opposite genomic sense within 3 kb. **D**, Correlation between log-enrichment of reads aligning to each complementary sequence in MDA5-binding experiment, and x_{ds} . The blue line shows the fit of a generalized additive model. **E**, relation between log-enrichment of reads aligning to each complementary sequence in J2-binding experiment, and x_{ds} . The blue line shows the fit of a generalized additive model. **F**, Type of repeat with the longest overlapping sequence in complementary sequences with high MDA5 signal and high x_{ds} ($x_{ds} > 0.5$) and complementary sequences with low MDA5 signal and low x_{ds} . **G**, Genomic distribution of high- x_{ds} regions in the human genome (right bar), compared with the distribution of the full genome (left bar). **H**, Type of repeat with the longest overlapping sequence in complementary sequences with high J2 signal and high x_{ds} and complementary sequences with low J2 signal and low x_{ds} .

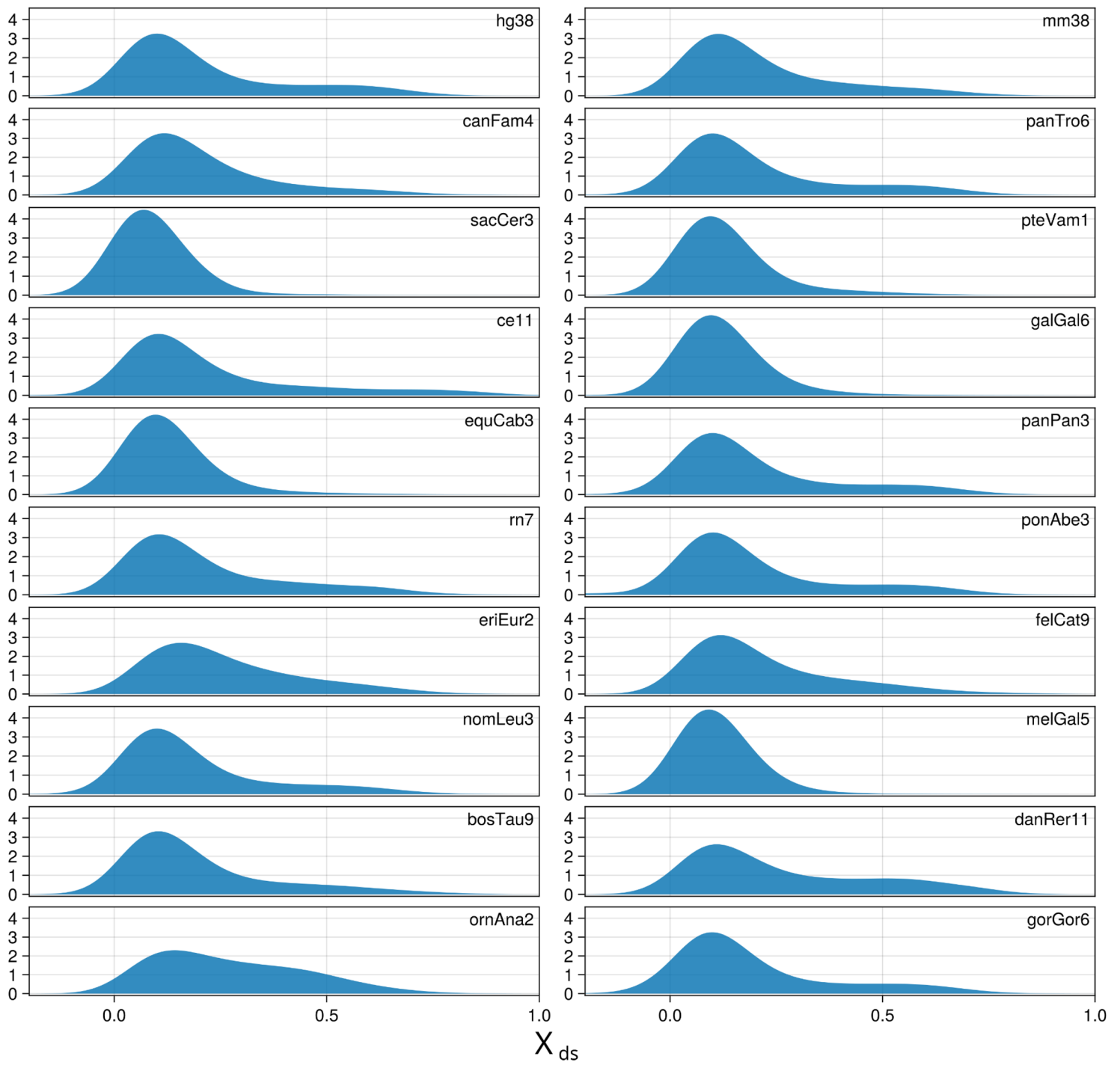


Extended Data Figure 3 | Volcano plot of repeat element expression of elements with double stranded force greater than 0.5 in H3B-8800 versus DMSO treated SF3B1-K700 mutant K562 cell lines.

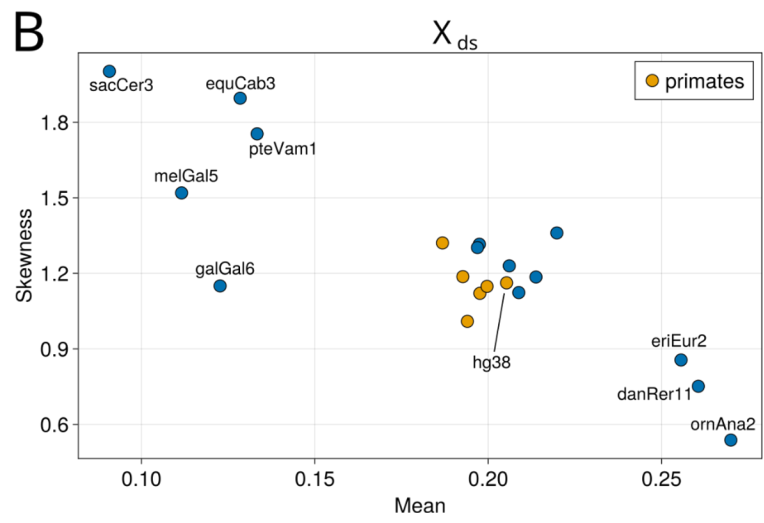
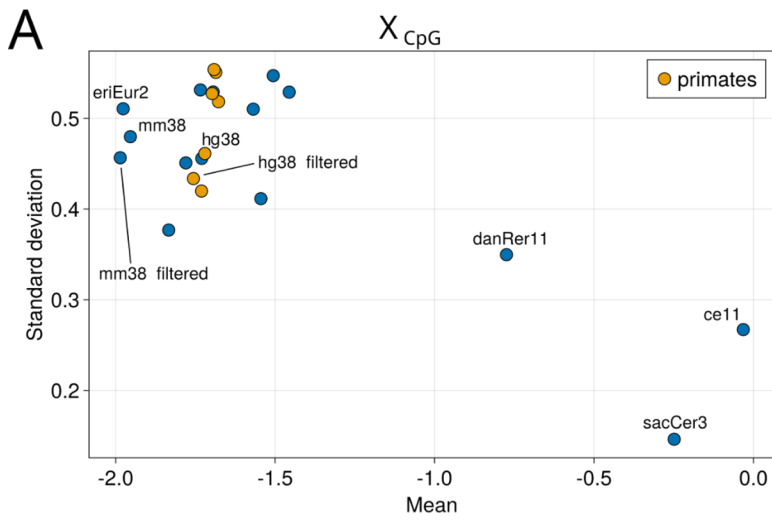


Extended Data Figure 4 | Distributions of x_{CpG} in several organism genomes (sliding window with transcript of length of 3000). For human and mouse, we also show, in orange, the profile of x_{CpG} the histogram of after excluding reads annotated as CpG islands or enhancers.

Probability Density Function



Extended Data Figure 5 | Distributions of x_{ds} in several organism genomes (sliding window with transcript of length of 3000).



Extended Data Figure 6 | A, Standard deviation versus mean of x_{CpG} computed for each 3000-base windows for each organism analyzed. Orange denotes points computed from primate genomes. **B**, Skewness versus mean of x_{ds} computed for each 3000-base windows for each organism analyzed. Orange denotes points computed from primate genomes.

Methods

Quantification of forces on sequence features

We define a Maximum Entropy (MaxEnt) framework [1] to determine the least constrained probability distribution over the set of sequences $s = \{s_1, s_2 \dots s_L\}$ of length L compatible with the observed occurrences (measurement) of a set of M features $\mathbf{N} = \{N_1, N_2 \dots N_m\}$. The distribution is written as

$$P(s|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\left(\sum_{m=1}^M x_m N_m(s)\right). \quad (1)$$

Here $Z(\mathbf{x})$ is a normalization factor:

$$Z(\mathbf{x}) = \sum_{s'} \exp\left(\sum_{m=1}^M x_m N_m(s')\right), \quad (2)$$

where the sum runs over all possible sequences s' having length L . The set of parameters $\mathbf{x} = \{x_1, x_2, \dots, x_M\}$, hereafter called *selective forces*, is chosen so that the average value of each feature over the distribution P matches the observed number of this feature N^{obs} in one or more reference sequences:

$$N_m^{\text{obs}} = \sum_s P(s|\mathbf{x}) N_m(s) = \frac{\partial \log Z(\mathbf{x})}{\partial x_m}. \quad (3)$$

The above equalities define a set of M coupled, nonlinear equations, with a unique solution due to the convexity of $\log Z(\mathbf{x})$. The forces are analogous to chemical potentials in statistical physics.

We use our formalism to calculate selective forces on sequences due constraints imposed by: (1) the force due to potentially immunogenic CpG motifs (the CpG force, x_{CpG}), (2) the forces on all nucleic acid motifs with length up to three nucleotides, and (3) the force on long complementary sequence stretches (the double-stranded force, x_{ds}). For the later case we use an equivalent direct approach to simplify the calculation.

Forces on CpG and other individual dinucleotides

To estimate the force on CpG motifs, x_{CpG} , arising from constraints on the usage of CpG dinucleotides we derive the MaxEnt distribution P of sequences with five forces chosen to reproduce the frequencies $f(\sigma)$ of the four nucleotides, $\sigma = \text{A,C,G,U}$, and the frequency of CpG dinucleotides only. Other dinucleotide forces are calculated in the same manner. We find that the single-nucleotide forces are, to a good accuracy, given by $x_\sigma = \log f(\sigma)$ (more specifically, the forces computed in this way are in a linear relationship with those computed with a full maximum entropy model, see **Methods Fig. 1A**), which allows us to approximate the normalization factor in Eq. (2) with

$$Z(x_{\text{CpG}}) = \sum_{\{s'_i=\text{A,T,C,G}\}} \left(\prod_{i=1}^L f(s'_i) \right) \exp(x_{\text{CpG}} N_{\text{CpG}}(s')). \quad (4)$$

A Newton's method based algorithm to efficiently (time $\sim O(L^2)$) calculate x_{CpG} , such that Eq. (3) (with $m = \text{CpG}$) is satisfied, was derived in [2]. A positive value of x_{CpG} for the observed feature count N_{CpG} in a given sequence indicates CpG motifs are enriched compared to what would be expected from a random sequence conditioned on single nucleotide usage only. A negative value corresponds to a depletion with respect to the null model.

Inference of general model based on mono-, di- and trinucleotide motif usage

To infer the full set of forces on each motif of length up to 3 nucleotides we extend the formalism used for individual dinucleotide forces alone. We model each viral or repeat family as a probability distribution over the sequences characterized by the frequency of each nucleotide, dinucleotide and trinucleotide motif, as described above. In the general case, this corresponds to an overall set of $M = 4 + 4^2 + 4^3 = 84$ forces to infer, such as x_{TCA} , x_{GGC} , x_{GA} , and x_T , respectively, on the motifs TCA, GGC, GA, and T. Note that, due to symmetries of the problems, only a subset of 39 parameters can vary freely, so we are left with their inference. For instance, the sum of N_A , N_C , N_G , and N_T is fixed (independent on the forces) and equal to L . However, many of these properties (such as the fact that $N_C = N_{AC} + N_{CC} + N_{GC} + N_{TC}$) strictly hold only for infinitely long sequences, but are very well approximated in longer sequences, such as those of length equal to 3000 nucleotides frequently studied here. Because of this and similar properties, we are free to fix a certain number of forces to an arbitrary value. For instance, we can set the force of each motif containing a T to zero, without losing generality.

To infer the remaining 39 forces we used a method analogous to the one developed for CpG forces, which allows for an efficient evaluation of the maximum entropy parameters through Eqs. (2) and (3).

To train the models on viral families, we used datasets for every RNA viral family and for HIV viruses collected from the Virus pathogen Database and Analysis Resource (<https://www.bv-brc.org/>) [3], after removing sequences with non-standard nucleotides (different from A, C, G, T) and duplicate sequences. Influenza A viral sequences were through the Influenza Research database (now housed at [3]) and filtered with the same criteria. The model for Influenza A viruses has been trained on the sequence obtained by joining the viral segments.

The models on repeat families have been trained in the same way, using consensus repeats from [4] and grouping them in family as annotated in [5]. Each model has been trained by computing the average frequencies of each motif for the viral or repeat family, then obtaining the number of motifs for the inference procedure through multiplication by the same length for each model (5000 nucleotides).

Once the force parameters of Eq. (3) have been fitted to match the motifs statistic of a set of viral genomes or repeats in the reference list, we quantify the similarity between viral and repeat families using the symmetrized Kullback-Leibler divergence between the corresponding probability distributions p_v and p_r given by $\frac{1}{2}(D_{KL}(p_v, p_r) + D_{KL}(p_r, p_v))$, where D_{KL} is the Kullback-Leibler (KL) divergence defined as

$$D_{KL}(p_v, p_r) = \sum_s p_v(s) \log \left(\frac{p_v(s)}{p_r(s)} \right). \quad (5)$$

D_{KL} is not an intensive quantity, as the probability distributions p_v and p_r depend on the length of the sequences modeled. In this work we fixed a reference length of 1000 nucleotides for all D_{KL} computations.

Computation of the Kullback-Leibler divergence

The model we define in Eq. (1) can be rewritten as

$$p(s) = \frac{1}{Z} e^{-E(s)}, \quad (6)$$

where $E(s)$ is an energy associated to each sequence, in the usual statistical physics sense. The KL divergence can be written as

$$D_{KL}(p_1, p_2) = \sum_s p_1(s) \log \left(\frac{p_1(s)}{p_2(s)} \right) = \log Z_2 - \log Z_1 + \sum_s p_1(s) (E_2(s) - E_1(s)). \quad (7)$$

As discussed in Methods, $\log Z_1$ and $\log Z_2$ can be computed exactly with the transfer matrix method. To compute the last term on the r.h.s. of Eq. (7) we define

$$Z_{12}(\lambda) = \sum_s e^{-E_1(s) + \lambda(E_2(s) - E_1(s))}, \quad (8)$$

and we have

$$\sum_s p_1(s) (E_2(s) - E_1(s)) = \left. \frac{\partial}{\partial \lambda} \log Z_{12}(\lambda) \right|_{\lambda=0}. \quad (9)$$

Forces on double-stranded RNA formation

We develop a framework to quantify the length of duplex strands. Given a reference sequence σ of length L we compute the frequencies $f(\sigma)$ of the single nucleotide $\sigma = A, C, G, U$. The feature $N_{\text{ds}}(s)$ is now the length of the longest subsequence of s whose complementary subsequence is also present in s , and which can therefore form a duplex of the same length. We present an intuitive derivation for the forces on double-stranded RNA, x_{ds} . That version is used preferentially in the text due to its interpretability. We also present the full MaxEnt approach, derived in an exactly parallel manner to that for forces on motifs. Due to the difficulty of computing exactly the corresponding $Z(x_{\text{ds}})$, which further justifies the intuitive approach, we utilize an approximate calculation. We show that both approaches give directly analogous results, and can therefore be used interchangeably. Both approaches are described below, along with their formal relationship.

Direct approach

Consider two subsequences $s = (s_1, \dots, s_K)$ and $s' = (s'_1, \dots, s'_K)$ of length K , with nucleotides drawn independently at random with the frequencies f . The probability that the two sequences are complementary is equal to $p^{\text{compl}}(K) = \alpha^K$, with

$$\alpha = \sum_{\sigma=A,C,G,U} f(\sigma) f(\sigma^{\text{compl}}). \quad (10)$$

where σ^{compl} denotes the complementary nucleotide to σ . In the presence of a biasing force x_{ds} acting on the length of the stretch the probability that the two sequences are complementary is modified into

$$p^{\text{compl}}(K, x_{\text{ds}}) = (\exp(x_{\text{ds}}) \alpha)^K = \tilde{\alpha}^K, \quad (11)$$

where $\tilde{\alpha} = \exp(x_{\text{ds}}) \alpha$. Positive and negative forces x_{ds} favor, respectively, longer and shorter complementary stretches than expected from the random nucleotide null model.

Consider now a sequence of length L , which we partition into $N = L/K$ subsequences of length K each. Under the simplifying assumption (that we check *a posteriori* in the following) that each pair of these subsequences is independent, the probability that none of them is fully complementary is

$$p_0(N) = (1 - \tilde{\alpha}^N)^M \quad (12)$$

where $M = N(N - 1)/2$ is the number of pairs of segments. Equivalently, $p_0(N)$ can be interpreted as the probability that the longest fully complementary segment has length $< N$. As a consequence the probability that the longest fully complementary segment is of length equal to N_{ds} reads

$$\begin{aligned} p(N_{\text{ds}}) &= (1 - \tilde{\alpha}^{N_{\text{ds}}+1})^M - (1 - \tilde{\alpha}^{N_{\text{ds}}-1})^M \simeq e^{-M\tilde{\alpha}^{N_{\text{ds}}}} \left(e^{-M\tilde{\alpha}^{N_{\text{ds}}}(\tilde{\alpha}-1)} - 1 \right) \\ &\propto \exp \left(N_{\text{ds}} \log(\tilde{\alpha}) + \log(M) - e^{N_{\text{ds}} \log(\tilde{\alpha}) + \log(M)} \right) = e^{-z - e^{-z}}, \end{aligned} \quad (13)$$

where $z = -N_{\text{ds}} \log(\tilde{\alpha}) - \log(M) = \frac{N_{\text{ds}} - \mu}{\beta}$ with $\beta = \frac{1}{\log(1/\tilde{\alpha})}$ and $\mu = \frac{\log(M)}{\log(1/\tilde{\alpha})}$. The approximations used are $\tilde{\alpha}^{N_{\text{ds}}} \ll 1$ and $M\tilde{\alpha}^{N_{\text{ds}}} \ll 1$, which are expected to hold in our case. Eq. (13) hence shows that the distribution of the longest fully complementary strand follows a Gumbel law, with mean $\mu + \beta\gamma$ (here γ is the Euler-Mascheroni constant) and variance $\beta^2\pi^2/6$.

As in Eq. (3), we can fix the force parameter x_{ds} by requiring that the average maximum length of fully complementary segment computed through the model, $\mu + \beta\gamma$, is equal to the value observed in a given sequence, $N_{\text{ds}}^{\text{obs}}$, and we obtain the equation

$$N_{\text{ds}}^{\text{obs}}(L, x_{\text{ds}}) \approx \frac{\log L}{\log \frac{1}{\sqrt{\exp(x_{\text{ds}})\alpha}}} + c, \quad (14)$$

for large values of L . With c being a correction term inserted to account for the set of simplifications done, and we estimate it directly from synthetic data as follows. As a check of the validity of the expression in Eq. (14), and to estimate the value of c , we fit the parameters x_{ds} and c from the set of maximum length of complementary segments in randomly generated RNA sequences of lengths ranging up to $L = 3000$ bases (**Extended Data Fig. 2A**). In this work, we consider both canonical Watson-Crick pairs and Wobble pairs as complementary basepairs. We obtain $c = -2.2$ and $x_{\text{ds}} = 0.06$, a value compatible with the zero force expected for this null model.

Eq. (14) with $c = -2.2$ can now be used to estimate x_{ds} for a reference sequence of length L (with nucleotidic frequencies f) and with longest complementary stretch of length $N_{\text{ds}}^{\text{obs}}$. We thus obtain a single metric to compare distribution of double-stranded segments across various RNA sequence ensembles and families diverse sequence statistics and lengths.

Maximum Entropy approach

We start from the probability that the longest fully complementary segment is of length equal to N_{ds} , $p(N_{\text{ds}})$. According to the maximum entropy principle, the probability distribution on sequences of length L which maximizes entropy while fixing the length of the maximum complementary segment is

$$p(\mathbf{s}) = \frac{1}{Z(\hat{x}_{\text{ds}})} \exp(\hat{x}_{\text{ds}} N_{\text{ds}}(\mathbf{s})) \quad (15)$$

with the normalization

$$Z(\hat{x}_{\text{ds}}) = \sum_{\mathbf{s}} \exp(\hat{x}_{\text{ds}} N_{\text{ds}}(\mathbf{s})) = S \sum_{N_{\text{ds}}} p(N_{\text{ds}}) \exp(\hat{x}_{\text{ds}} N_{\text{ds}}), \quad (16)$$

where \hat{x}_{ds} is the double-stranded force acting on the length of the longest complementary stretch using the MaxEnt approach, and S is the total number of sequences of length L . Under this model, the probability of observing a sequence with maximum complementary segment of length N_{ds} is

$$p(N_{\text{ds}}) \propto \exp(-z - e^{-z} + \hat{x}_{\text{ds}} N_{\text{ds}}), \quad (17)$$

where $z = -N_{\text{ds}} \log(\alpha) - \log(L(L-1)/2)$. For large L we can compute the average value of N_{ds} by integrating the continuous version of the probability distribution, and we obtain

$$\langle N_{\text{ds}} \rangle = \mu - \frac{\psi(\eta)}{\log(1/\alpha)}, \quad (18)$$

where $\mu = \frac{\log(M)}{\log(1/\alpha)}$, ψ is the digamma function and $\eta = 1 - \beta\hat{x}_{\text{ds}}$.

We can now substitute $\langle N_{\text{ds}} \rangle$ with the observed value of N_{ds} in the sequence under analysis, and add a constant c' to take care of the approximations done, and we obtain the equation

$$N_{\text{ds}}^{\text{obs}} = \mu - \frac{\psi(\eta)}{\log(1/\alpha)} + c' \quad (19)$$

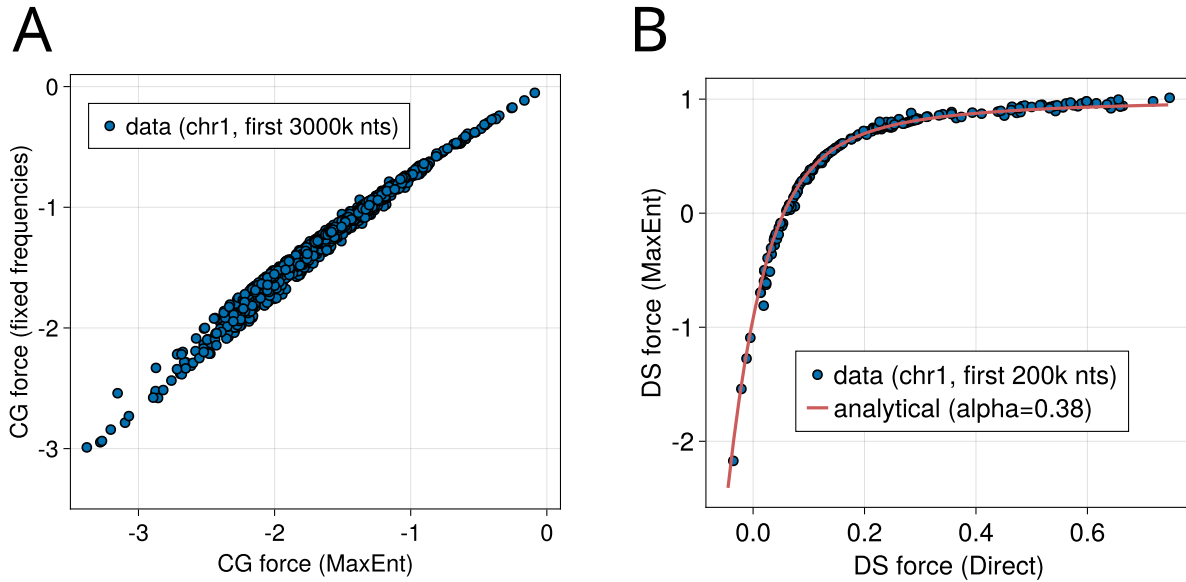
This equation can be used to estimate the MaxEnt double-strand force \hat{x}_{ds} . We estimated c' from randomly generated RNA sequences of lengths ranging up to $L = 3000$ bases and obtained the value of -1.79 .

Comparison of Direct and MaxEnt approaches

The direct approach to calculating x_{ds} is used in the manuscript. While the direct approach and the MaxEnt force are in general different, for a fixed value of α there is a monotonic relationship between these two quantities:

$$\hat{x}_{ds} = \gamma \left(1 - \psi^{-1} \left(-\frac{2 \log(L)}{\gamma} \frac{\gamma x_{ds}}{\gamma - x_{ds}} + \gamma (c' - c) \right) \right), \quad (20)$$

where $\gamma = \log(1/\alpha)$, c' and c are the constants inferred for the two approaches from synthetic data, \hat{x}_{ds} is the MaxEnt double-stranded force, x_{ds} is the double-stranded force computed via the direct approach, and ψ^{-1} is the inverse of the digamma function in the interval $(0, +\infty)$. In particular, we checked that even when α is computed for each 3000 nucleotide window, the relationship gives an extremely good approximation, thus any subsequence with double-stranded force larger than a given threshold would be equivalently characterized by high MaxEnt double-stranded force, as shown in **Methods Fig. 1B**.



Methods Figure 1 | A, Comparison of x_{CpG} computed with a full-maximum-entropy model (x axis) and with the method used in this work, Eq. (4). **B**, x_{ds} computed with the Direct approach (horizontal axis) versus the MaxEnt approach (vertical axis), for the sliding windows within the first 200k bases of the first chromosome of the human genome. The red line correspond to the analytical relationship between the two forces computed for $\alpha = 0.38$, which is the value obtained by considering the full human genome.

Compressing the x_{ds} table

We initially computed x_{ds} for each sliding window of 3kb in the human genome (hg38 assembly). The resulting table of windows and associated x_{ds} values, however, are not appropriate for some

of the successive analyses, mostly because the same pair of complementary sequences appears in many close-by windows. To deal with this issue we produced a new, compressed x_{ds} table with the following rules: (i) we discarded windows having the two complementary sequences less than 10 nucleotides away from the window ends; (ii) whenever more than one window have exactly the same pair of complementary sequences, we took the most upstream window and we discarded the others. Rule (i) prevents the edges of the windows from “cutting” one of the two complementary sequences, generating cluster of very similar pairs of complementary sequences in consecutive windows, while rule (ii) prevents the presence of multiple windows associated to the same pair of complementary sequences.

Evolutionary dynamics of a sequence motif with force relaxation formalism

One can harness the formalism developed in Eqs. (1) - (3) to study an evolutionary dynamics of number of motifs N_m , as it approaches the steady state (equilibrium) value N_m^{avg} [2]. As a sequence evolves it undergoes mutations, which cause changes in the number of motifs (and hence associated value of x_s). To model the evolutionary dynamics of sequences, we assume the number of motifs (N_m) evolves according to the relaxation dynamics given by

$$\tau \frac{dN_m(t)}{dt} = -x_s(N_m(t)) + x_m^{\text{eq}}, \quad (21)$$

where τ sets the timescale. The number of motifs reaches its stationary (equilibrium) value when $x_s = x_m^{\text{eq}}$, at which point the selective force is balanced by the entropic forces which randomize sequences. It is convenient to express (21) as

$$\tau \frac{dx_s}{dt} = -(-x_s(t) + x_m^{\text{eq}}) \text{var}(x_s|N_m), \quad (22)$$

where $\text{var}(x_s|N_m)$ is the variance of x_s for a given N_m .

If we can express $\text{var}(x_s|N_m)$ as a function of x_s , it is possible to obtain a solution of (22) that can then be fitted to the dataset with timescale τ , thus providing the approximation of relaxation dynamics, along with the estimate of the time it will take to $x_s(t)$ (and hence the number of the corresponding sequence motifs m) to reach its equilibrium value. For the case of HSATII and LINE-1 we fit $\text{var}(x_s|N_m)$ as a quadratic function of x_s .

Kimura-based model of population genetics for the evolution of sequence motifs

In addition to the force relaxation model introduced above, we present here a different approach to study the evolution of nucleotide sequence motifs based on the Kimura model of sequence evolution. We implement the model numerically, and evolve a set of sequences to provide a null model of neutral sequence evolution. For each simulation step, we pick a random base and mutate it to a randomly chosen different base with a given probability. We consider different possible mutation probabilities depending on the type of base it is mutating into, as well as on the context (identity of the bases in the neighborhood), as transversion (purine mutating to pyrimidine or vice versa) and transition (purine mutating to purine or pyrimidine mutating to pyrimidine) substitutions in sequences can have different likelihood [6].

Additionally, in vertebrates and plants, mutations in CpG context are known to be more common due to CpG hypermutability [7]. Hence, for the mutation rates in the model implementation, we use different ratios of mutation rates $\mu_{\text{TiCpG}}:\mu_{\text{TVCpG}}:\mu_{\text{Ti}}:\mu_{\text{TV}}$ (corresponding to nucleotide transitions and transversions in CpG context and to transitions and transversion in non-CpG context). In particular, we consider the ratios introduced in Ref. [6] and which are listed in Table 1. The increased mutation rate on CpG dinucleotides has the effect of reducing the expected CpG number after relaxation of a sequence, and it can thus be related to an equilibrium CpG force. To compute this force we considered the model without the transition-transversion bias (that cannot affect the number of CpGs at equilibrium) and we can write for the number of CpG in a sequence

$$N_{\text{CG}}^t = N_{\text{CG}}^{t-1} - 2\mu N_{\text{CG}}^{t-1} + \frac{1}{3}\gamma(N_{\text{C}}^{t-1} - N_{\text{CG}}^{t-1}) + \frac{1}{3}\gamma(N_{\text{G}}^{t-1} - N_{\text{CG}}^{t-1}) \quad (23)$$

while the number of C nucleotides N_{C} evolves as

$$N_{\text{C}}^t = N_{\text{C}}^{t-1} - \mu N_{\text{CG}}^{t-1} - \gamma(N_{\text{C}}^{t-1} - N_{\text{CG}}^{t-1}) + \frac{1}{3}\gamma(L - N_{\text{C}}^{t-1} - N_{\text{CG}}^{t-1}) + \frac{1}{3}\mu N_{\text{CG}}^{t-1} \quad (24)$$

$\mu_{TiCpG}:\mu_{TvCpG}:\mu_{Ti}:\mu_{Tv}$	x_{CpG}^{eq}
40:10:4:1	-1.9
40:4:4:1	-1.8
40:1:4:1	-1.7
4:4:4:1	-0.3
20:4:4:1	-1.2
27:2:4:1	-1.4

Table 1| Ratios of dinucleotide mutation rates (transition and transversion with and outside of CpG context) and a corresponding value of the equilibrium force on the CpG dinucleotide

and similarly for the number N_G of G nucleotides. In these equations, μ is the probability of a substitution happening in a CpG context, and γ is the probability of the mutation happening outside a CpG context, and L is the length of the sequence. At equilibrium, we find

$$\begin{aligned} \frac{N_{CG}}{L} &= \frac{1}{14r + 2} \\ \frac{N_C}{L} &= \frac{N_G}{L} = \frac{3r + 1}{14r + 2}, \end{aligned} \quad (25)$$

where $r = \mu/\gamma$. We can now compute the corresponding CpG force x_{CpG}^{eq} using the fact that the force is approximately equal to the logarithm of the relative frequency of the dinucleotide motif $x_{CpG} \approx \log(f(CpG)/f(C)f(G))$ [8]. The ratios 40:10:4:1, 40:4:4:1 and 40:1:4:1 provide the closest approximation to relaxation to the force observed in the genome. For the neutral model, we used the 40:10:4:1 ratio as it was closer to the saturated value of x_{CpG} of the LINE-1 elements.

Analysis of forces across species

Values of x_{CpG} or x_{ds} were computed for each 3000 kb sliding window. Non-numeric values were excluded. Windows with one or more ambiguous characters were excluded. To compute the distribution of a force across all windows, the values of that force were sorted numerically and every 100th entry was retained (entry number 50, 150, 250, ..., etc.). The distribution density was computed using a Gaussian kernel with bandwidth 0.05 as implemented in scikit-learn package [9]. The density was computed for all points within the target interval ([-5:2] for x_{CpG} , and [-2:3] for x_{ds}) with a step of 0.005. We compute the FDR as the ratio of the cumulative area of the null model distribution to the right of the cutoff to the cumulative area of the distribution to the right of the cutoff. The null model distribution is fitted as a Gaussian distribution with the peak at the point with the maximal density, standard deviation was computed using the 20 points to the right of the peak. We computed these values across the species: Pan troglodytes troglodytes, Pan paniscus, Gorilla gorilla gorilla, Pongo abelii, Nomascus leucogenys, Canis lupus, Danio rerio, Mus musculus, Rattus norvegicus, Equus caballus, Bos taurus, Gallus gallus, Felis catus, Pteropus vampyrus, Caenorhabditis elegans, Saccharomyces cerevisiae, Meleagris gallopavo, Erinaceus europaeus, and Ornithorhynchus anatinus, in addition to humans.

For humans and mice we also performed an additional analysis which excluded both enhancers and CpG islands. Coordinates of enhancers from the FANTOM database were lifted from hg19 to hg38, and mm9 to mm38 using LiftOver tool from UCSC ([10, 11], [12]). Coordinates of CpG islands for hg38 and mm38 were downloaded from UCSC. "Filtered" data for hg38 and mm38 in the CpG plot consist only of the windows which have zero overlap with CpG islands and enhancers.

Analysis of evolutionary conserved sequences with high x_{CpG} or x_{ds}

We considered the genomes of 5 species (Pan troglodytes troglodytes, Pan paniscus, Gorilla gorilla gorilla, Pongo abelii, Nomascus leucogenys) in addition to the human genome to look for conserved regions with high x_{CpG} or high x_{ds} . After computing both forces for each species exactly as we did for the human genome, we considered the set of $x_{\text{CpG}} > 0$ and $x_{\text{ds}} > 0.5$. We reduced the number of the CpG windows by clustering together all those which overlapped more than 1000 bases, and from each cluster we only considered the window with the highest value of x_{CpG} . For the x_{ds} windows, we first compressed them by excluding windows with cut-off complementary sequences or with identical complementary sequences (as discussed above), then for each window we extracted the subsequence spanning the pair of complementary sequences. Finally we clustered together all overlapping sequences and from each cluster we only considered the window with the highest value of x_{ds} .

We then ran BLAST to compare each of these sequences with high x_{CpG} or high x_{ds} extracted from the human genome and we retained any significant match (whenever one sequence of a given organism matched with more than one human sequences we only kept the match with the highest BLAST score) [13]. The result of this procedure consists in two sets of sequences for each organism that are alignable to human sequences, one for the high x_{CpG} and one for the high x_{ds} . We then computed the overlaps between the set A of high x_{CpG} (or high x_{ds}) organism sequences and that of the human, B , defined as

$$O(A, B) = \frac{|A \cap B|}{\min(|A|, |B|)}, \quad (26)$$

where a sequence of A belongs to the intersection $A \cap B$ if and only if it is alignable with significant score by BLAST to a sequence of B . $|A|$ is the size of the set A .

We next compared the overlaps we obtained with the time since the most recent common ancestor, taken from [14]: 6.6 million of years ago for Homo Sapiens, Pan troglodytes troglodytes and Pan paniscus; 8.9 for Homo Sapiens, Pan troglodytes troglodytes, Pan paniscus and Gorilla gorilla gorilla; 13.8 for Homo Sapiens, Pan troglodytes troglodytes, Pan paniscus, Gorilla gorilla gorilla and Pongo abelii; 16.3 for all the species considered here. Afterwards we focused on the annotations in RepeatMasker [5] and DFAM [4] databases to check for potential overlaps of the conserved sequences with annotated repeats in the human genome. For the set of high-CpG sequences we used the DFAM dataset as we found a better annotation for HSATII with respect to the one in RepeatMasker. We associated each conserved high-CpG sequence with a repeat if the corresponding window overlapped with its position as annotated in the database (for the windows overlapping with more than one annotated repeats, the one with the largest overlap was considered).

For the set of x_{ds} sequences we used the RepeatMasker dataset [5]. In this case, each x_{ds} window is characterized by two fully complementary sequences, and we searched for repeats overlapping with each of them (when overlaps with multiple repeats were found, the one with the largest overlap was considered). We observed 3 different cases: one or two of the two sequences do not overlap with any repeat annotated in RepeatMasker; each sequence overlapped with a repeat, both repeat being of the same family and annotated in the two strands of the genome, e.g. one sequence overlapping with AluS+ and one with AluS- (IR repeats); each sequence overlapped with a repeat, but of different families or of the same family but on the same strand of the genome (non IR repeats).

Sequence ensembles

The LINE-1 sequences were obtained from L1Base2 database [15]. We separately downloaded all the sequences annotated as full-length intact and hence are more likely to still be active (146 for human genome and 2811 for mouse genome), and sequences annotated as full-length non-intact (13148 for human genome and 14076 in mouse genome). We separately aligned each of the non-intact sequences with each of the respective intact sequences using pairwise alignment and calculated the Kimura distance between the sequences [16]. We then calculated the average distance for each of the non-intact sequences from the intact-sequences, and furthermore calculated the number of CpG motifs in each sequence.

Sequences of all inserts of HSATII and all other Human Genome repetitive elements considered in this work have been obtained from the DFAM database [17] (version introduced in 2016). Each family of sequences in the DFAM database contains sequences of all its inserts in the human genome and their consensus sequence, as well as with the hidden Markov Chain Model (HMM) that we use to align inserts with respect to the consensus sequence. For comparison of sequences of inserts with respect to their consensus sequence, we only consider inserts of length longer than 150 bases. To quantify the difference between the insert sequence and the consensus sequence, we use the Kimura distance [16] between the consensus and its inserts.

We note that we use the Kimura distance [18] from the consensus sequence (for inserts from DFAM) or from average of all full-length non-intact sequences (for LINE-1s from L1Base2) as a measure of time, assuming that it is proportional to the time since insertion of the particular transposable element into the species genome. All the sequences studied in this work have been obtained from hg38 genome assembly.

Search of long transcripts with complementary regions

We scanned the hg38 genome assembly for transcripts that can be possible source of long duplex formation. To this aim, for each window of length 3000 bases (taken in the positive sense of the read), we calculate the double-stranded force x_{ds} from Eq. (14), using the window-specific nucleotide frequencies to obtain α from Eq. (10). We considered windows resulting in $x_{ds} > 0.5$ as having a high double-stranded force when compared to the rest of the genome.

Sequence complexity quantification

We use an approximation of Kolmogorov complexity [19] to quantify how “non-trivial” complementary segments are. Adopting the approach from Ref. [20], we use the size (in bytes) of the sequence compressed with gzip software as a proxy of the Kolmogorov complexity. Simple sequences, e.g. poly(AT) or poly(C) and poly(G), will have low complexity, as they can be compressed to a smaller size than a completely random sequence of the same length (which would have maximum complexity).

Estimate of genome regions with high double stranded force

To estimate x_{ds} for a given repeat loci, we intersect each repeat loci with the calculated 3kb genomic windows that have high dsRNA forces ($x_{ds} > 0.5$). The Start and End coordinates of the corresponding dsRNA sequence pairs, which overlap with the repeat loci that match the criteria:

$|\log_2FC(\text{treated}/\text{untreated})| > 0.5$ and $FDR < 0.05$, were used to annotate different genomic features. We counted the genomic features of the predicted double-stranded RNA sequences that overlap with the upregulated repeats ($\log_2FC > 0.5$ and $FDR < 0.05$), and of those that overlap with the downregulated repeats ($\log_2FC < -0.5$ and $FDR < 0.05$). These counts have been compared with the genomic feature counts of all dsRNA sequences that overlap with the transcribed repeats to calculate the odds ratio and p-value using the Fisher Exact test.

Transcriptome analysis

Analysis of repeats from splice inhibitors

Raw RNAseq data (GSE95011) associated with the Seiler, et al., 2018 study [21] were downloaded from NCBI. Briefly, reads were trimmed and quality checked using skewer first and then mapped to the human genome (hg38) and repetitive elements from RepBase [22, 23]. In quality check, Illumina reads were trimmed to remove N's and bases with quality less than 20. After that, the quality scores of the remaining bases were sorted, and the quality at the 20th percentile was computed. Reads quality less than 15 at the 20th percentile or shorter than 40 bases were discarded. Only paired reads that passed the filtering step were retained. Only paired reads which both pass the quality check were mapped to the reference genome (hg38) using STAR (v2.7) with default parameters. Gene counts were assigned based on Gencode annotation using featureCounts (Subread package) with the external Ensembl annotation. Repeats counts per element subfamily were primarily quantified against RepeatMasker using featureCounts and then adding the counts of the unassigned reads that mapped to Repbase consensus sequence. Repeat counts of a given family is the sum of mapped reads to RepeatMasker and unmapped reads against Repbase.

Counts filtering, normalization and statistical analysis

Gene expression in terms of log₂-CPM (counts per million reads) was computed and normalized across samples using the TMM (trimmed-mean of M-values) method using the calcNormFactors() and cpm() functions from edgeR package [24]. These low-count values (CPM < 2) were removed before calculating the size factor for each sample. Then, filtered CPM was log₂ transformed and used in heat-map visualization and downstream statistical analysis. On the heatmap, genes (rows) were scaled by z-score scaling. Heat maps were generated by the R statistical programming package. Differential expression analysis was performed using limma package [25] between splicing modulator H3B-8800 treated versus DMSO treated SF3B1-K700 mutated cell line k562 for a given locus. The adjusted p-values were calculated using the Benjamini-Hochberg correction [26].

Analysis of MDA5 binding transcripts

We have used the double-stranded force calculation to score RNA-Seq transcripts identified in Ref. [27] to bind to MDA5 receptors. We have further looked if any of the identified transcripts that bind MDA5 from Ref. [27] are also overlapping with transcripts that have been identified as MDA5 ligands in Ref. [28].

RNA extractions and co-immunoprecipitations (RIP) from EC cells

Embryonal carcinoma cells were cultured at 37° C in humidified incubators maintained with 7% CO₂ atmosphere. N2102Ep Clone 2/A6 cells (Merk, #06011803) were cultured in DMEM (high glucose, no sodium pyruvate; Thermo Fisher, #11965092), supplemented with 10% (v/v) fetal bovine serum, 1x penicillin/streptavidin and 2 mM Glutamine (Thermo Fisher, #25030024). Large-scale growth, harvesting, cryo-milling, and co-IP was achieved as previously described [29–31], summarized as follows. α -ORF1p-, α -ORF2p-, and α -ZCCHC3-targeted co-IPs used in-house made [32] magnetic affinity media: for α -ORF1p [15 μ g antibody / mg magnetic beads], we used the 4H1 monoclonal antibody (Millipore Sigma, #MABC1152); for α -ORF2p [10 μ g / mg magnetic beads] we used the clone 9 monoclonal antibody [33]; for α -ZCCHC3 [10 μ g / mg magnetic beads] we used the rabbit polyclonal antibody (Proteintech, #29399-1-AP). Co-IPs were conducted using 100 mg cell powder, extracted at 25% (w/v) in 20 mM HEPES pH 7.4, 500 mM or 300 mM NaCl, 1% (v/v) Triton X-100, 1x protease inhibitors (Roche, #1187358001), and 0.4% (v/v) RNasin (Promega, #2515). Centrifugally clarified cell extracts were incubated with affinity medium (20 μ l of slurry for α -ORF1p and α -ORF2p, and 15 μ l of slurry for α -ZCCHC3) for 30 minutes at 4° C. The solutions were made with nuclease-free H₂O and experiments were conducted using nuclease-free tubes and pipette tips. Macromolecule extractions performed on this cell line as described typically yielded between 450 - 500 μ l of soluble extract at 6 - 8 mg/ml of protein as assessed by Bradford assay (Thermo Fisher, #23200). After target capture, washing the media was performed with the same solution without protease inhibitors and with RNasin at 0.1% (v/v). RNAs were eluted from the affinity media after RIP with 250 μ l of TRIzol Reagent (Thermo Fisher, #15596026). After adding chloroform to the TRIzol eluate, the separated aqueous phase (containing RNAs) was obtained using Phasemaker tubes (per manufacturer's instructions; Thermo Fisher, #A33248), and was then combined with an equal volume of ethanol and further purified using a spin-column according to the manufacturer's instructions (Zymo Research, #R2060). For α -ZCCHC3 co-IP, two 100 mg-scale preparations were pooled prior to spin column purification. Eluates from α -ORF1p and α -ORF2p co-IPs were not treated with DNase I on-column, this was done during the sequencing library preparation (described, below); eluates from α -ZCCHC3 were DNase I treated on-column. Purified nucleic acids from α -ORF1p and α -ORF2p co-IPs were eluted in 6ul of nuclease-free water; purified nucleic acids from α -ZCCHC3 IPs were eluted in 10ul of RNase-Free water; in all cases 1 μ l was used for quality analysis and the remainder conserved for RNA-seq. Mock RNA co-IP controls were prepared in an identical manner using either naïve polyclonal mouse IgG (control for α -ORF1p; Millipore Sigma, #I5381) or naïve polyclonal rabbit IgG (control for α -ORF2p: Innovative Research, #IRBIGGAP10MG; control for α -ZCCHC3: Millipore Sigma, #I5006). Total RNA controls were prepared by combining up to 35 μ l of the clarified cell extracts with up to 500 μ l of Trizol, vortex mixing for 1 min, then snap freezing in liquid N₂ - and then later proceeding as above.

cDNA library preparation and RNA-seq

All the sequenced samples/replicates that are reported in this study are listed in Supplementary Table 6.

α -ORF1p and α -ORF2p RIP-seq

RNA extractions were quantified and quality controlled using RNA Pico Chips (Cat. #5067-1513) on an Agilent 2100 BioAnalyzer. RNA-Seq cDNA libraries were prepared using the Trio RNA-Seq

Library Prep kit (Tecan, #0357-A01) with AnyDeplete Probe Mix-Human rRNA (Tecan, #S02305). DNase treatment preceded cDNA synthesis. cDNA synthesis: 3 - 5ng of input RNA from α -ORF1p RIPs and mock IPs, 1 ng from α -ORF2p RIPs and mock IPs, and 50 ng of total RNA were used with 8 (2+6) cycles of pre-depletion PCR library amplification and 8 (2+6) cycles of post-depletion amplification; the libraries were purified using Agencourt AMPure XP beads (Beckmann Coulter), quantified by qPCR, and the size distribution was checked using the Agilent TapeStation 2200 system. Final libraries were sequenced, paired-end, at 50 bp read-length on an Illumina NovaSeq 6000 v1.5 with 2% PhiX spike-in.

α -ZCCHC3 RIP-seq

RNA extractions were quantified and quality controlled using an Agilent TapeStation 4200 and High Sensitivity RNA ScreenTape (Agilent, #5067-5579). RNA-seq cDNA libraries were prepared using the SMARTer Stranded Total RNA-Seq Kit v3 - Pico Input Mammalian (Takara, #634485), including rRNA depletion during library construction. cDNA synthesis: 5 ng of input RNA from α -ZCCHC3 RIPs and total RNA, and 1 ng of input RNA from mock IPs were used with 5 cycles of pre-depletion PCR amplification and 12 cycles (α -ZCCHC3 RIPs and total RNA) or 14 cycles (mock IPs) of post-depletion amplification. Libraries were purified using NucleoMag beads supplied in the library preparation kit and subsequently quantified using the Qubit 4 Fluorometer and the Qubit dsDNA HS assay kit (Invitrogen, #Q32854). The size distribution was checked using the TapeStation 4200; noting that primer dimers (\sim 150bp) were persisted in the mock IP libraries (motivating an additional round of cleaning). To treat all libraries equally, they were pooled in a 4:4:1 ratio (α -ZCCHC3 RIPs:total RNA:mock IP) based on molarity of fragments of interest (range \sim 200-1000 bp). An additional round of cleanup with NucleoMag beads was done to remove the primer dimers. A size selection of the final library pool was performed on a 2% E-gel EX (Invitrogen, #G401002) to exclude small fragments (less than about 200bp) and the DNA was eluted from the gel slices using the Zymoclean Gel DNA Recovery Kit (Zymo, #D4001), followed by quantification (Qubit) and quality control (TapeStation). Final libraries were sequenced, paired-end, at 250bp read-length on an Illumina NovaSeq 6000 platform.

RIP-seq read mapping and quantification

Reads were trimmed and quality checked using skewer [34]. Briefly, ends of the reads were trimmed to remove Ns and bases with quality less than 20. After that, the quality scores of the remaining bases were sorted, and the quality at the 20th percentile was computed. Reads were discarded if their quality at the 20th percentile was less than 15. In addition, reads shorter than 40 bases after trimming were discarded. If at least 1 of the reads in the pair failed the quality check and had to be discarded, we discarded the mate as well. Quality filtered reads were mapped to annotated repeat loci in RepeatMasker using software: Quantifying Interspersed Repeat Expression (SQuIRE) (<https://github.com/wyang17/SQuIRE>) [35]. Briefly, the SQuIRE pipeline first obtains reference annotation files from RepeatMasker, then aligns reads using STAT and, lastly, quantify locus-specific repeat expression by redistributing multi-mapping read fractions in proportion to estimated TE expression with an expectation-maximization algorithm.

CpG quantification for repeats in hg38

Sequences of repeats were extracted from hg38 based on the coordinate annotation in RepeatMasker. x_{CpG} was then calculated for those hg38 derived repeat sequences. L1 inserts reported

in RepeatMasker are considered intact whenever they have minimum 80% overlap with inserts annotated as full-length intact in the L1Base database.

Selection of RIP-Seq enriched repeats/transcripts

Samples extracted at 300mM and 500mM NaCl were used as replicates to increase statistical power after checking for transcript composition similarity (via hierarchical clustering). Targeted protein enriched transcripts were selected by $\text{Log}_2(\text{Co-IP}/\text{mock}) > 3$, $\text{Log}_2(\text{Co-IP}/\text{total RNA}) > 3$, and Benjamini-Hochberg adjusted p-value < 0.05 . Similarly, target protein depleted transcripts were selected by $\text{Log}_2(\text{Co-IP}/\text{mock}) < -3$, $\text{Log}_2\text{FC}(\text{Co-IP}/\text{total RNA}) < -3$, and Benjamini-Hochberg adjusted pvalue < 0.05 .

Immunoprecipitation of double-stranded RNA by J2 antibody

Experimental protocol

Protein G Dynabeads were washed twice and resuspended in antibody conjugation buffer (1x PBS, 2mM EDTA, 0.1% BSA (w/v)). 5 μ g of anti-dsRNA mAb (J2) (SCICONS, cat# 10010500) were bound to 30 μ l of washed beads overnight at 4° C on a rotating wheel. 10^7 Patient-derived POP92 cells per IP were fixed with 0.1% paraformaldehyde at room temperature (RT) for 10 min. Immediately, cells were quenched by adding glycine and washed twice with cold PBS. Crosslinked cells were lysed in lysis buffer (20mM Tris [pH 7.5], 150mM NaCl, 10mM EDTA, 10% Glycerol, 0.1% NP-40, 0.5% Triton X-100, supplemented with protease inhibitor tablet) for 15 minutes on ice. Following a spin at 12,000g at 4° C for 15 minutes, supernatant was transferred to a new eppendorf. The lysate was then immunoprecipitated using 30 μ L antibody-conjugated Dynabeads per IP reaction overnight at 4° C in a rotator. Following magnetic separation, beads were washed three times with high salt wash buffer (20 mM Tris pH 7.5, 500 mM NaCl, 10 mM EDTA, 10% glycerol) and resuspended in 1X TBS. Per IP, 2 μ L of Promega RNasinPlus RNase inhibitor (Fisher Scientific, PRN2611) and 0.5 μ L of proteinase K (NEB, P8107S) was added. Decrosslinking was performed for all the IP samples at 65° C for 15 minutes. The Direct-zol RNA MiniPrep kit (Zymo Research, R2051) was used to extract RNA from IP supernatant. Samples were treated with Turbo DNase to remove any DNA contamination in the extracted RNA. Library prep was performed using Illumina Stranded total RNA ligation with Ribo Zero plus according to the manufacturers protocol. Samples were sequenced on a NovaSeq 6000 using paired end reads with 100 cycles.

Analysis of J2 immunoprecipitation RNAseq

RNAseq controls not enriched with J2 antibody from untreated POP92 cells were downloaded from GEO (submission GSE145639, samples: GSM4322694 and GSM4322693) [27]. 25 bp of J2 RNAseq reads were cut off with cutadapt to match the length of RNAseq control reads. All samples were aligned to the human genome hg38 using STAR with default settings [36]. BAM files were sorted using samtools. The compressed x_{ds} table was used to count fragments in the RNAseq data, see above for details on how this table was generated. The table describes windows for the genome as well as a complementary sequence which can form a double stranded sequence. Every complementary sequence has a seqA and a seqB part which was split into two different files. featureCounts was used to count the number of fragments aligning to seqA and seqB in J2 and RNAseq control BAM files including information about strand and reporting multimapping and multi-overlapping reads as fractional counts [37]. Counts for each seqA and seqB were then merged for each complementary sequence. For each complementary sequence, log2FC (mean J2 / mean control) was calculated and plotted against x_{ds} using geom_point, geom_density_2d and geom_smooth from the R package ggplot2 [38].

Analysis of MDA5 protection assays

Raw sequencing data was downloaded from GSE103539 and GSE145639 [28] and aligned to hg38 using STAR. The following settings were used to increase the mapping due to the repetitive nature of the data [27]: `-outFilterMultimapNmax 1000 -outSAMmultNmax 1 -outFilterMismatchNmax 10 -outMultimapperOrder Random -winAnchorMultimapNmax 1000`. After mapping, the data were processed as described above for J2 immunoprecipitation.

References for Methods Section

- [1] Jaynes, E. T. Information theory and statistical mechanics. *Physical review* **106**, 620 (1957).
- [2] Greenbaum, B. D., Cocco, S., Levine, A. J. & Monasson, R. Quantitative theory of entropic forces acting on constrained nucleotide sequences applied to viruses. *Proceedings of the National Academy of Sciences* **111**, 5054–5059 (2014).
- [3] Olson, R. D. *et al.* Introducing the bacterial and viral bioinformatics resource center (bv-brc): a resource combining patric, ird and vipr. *Nucleic acids research* **51**, D678–D689 (2023).
- [4] Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The dfam community resource of transposable element families, sequence models, and genome annotations. *Mobile DNA* **12** (2021).
- [5] Smit, A., Hubley, R. & Green, P. Repeatmasker open-4.0. 2013–2015 (2015).
- [6] Suzuki, Y., Gojobori, T. & Kumar, S. Methods for incorporating the hypermutability of cpg dinucleotides in detecting natural selection operating at the amino acid sequence level. *Molecular biology and evolution* **26**, 2275–2284 (2009).
- [7] Subramanian, S. & Kumar, S. Higher intensity of purifying selection on > 90% of the human genes revealed by the intrinsic replacement mutation rates. *Molecular biology and evolution* **23**, 2283–2287 (2006).
- [8] Di Gioacchino, A. *et al.* The heterogeneous landscape and early evolution of pathogen-associated cpg dinucleotides in sars-cov-2. *Molecular Biology and Evolution* **38**, 2428–2445 (2021).
- [9] Pedregosa, F. *et al.* Scikit-learn: Machine learning in python. *the Journal of machine Learning research* **12**, 2825–2830 (2011).
- [10] Lizio, M. *et al.* Update of the FANTOM web resource: expansion to provide additional transcriptome atlases. *Nucleic Acids Research* **47**, D752–D758 (2018). URL <https://doi.org/10.1093/nar/gky1099>. <https://academic.oup.com/nar/article-pdf/47/D1/D752/27437410/gky1099.pdf>.
- [11] Lizio, M. *et al.* Gateways to the fantom5 promoter level mammalian expression atlas. *Genome biology* **16**, 1–14 (2015).
- [12] Hinrichs, A. S. *et al.* The ucsc genome browser database: update 2006. *Nucleic acids research* **34**, D590–D598 (2006).
- [13] Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *Journal of molecular biology* **215**, 403–410 (1990).
- [14] Thinh, V. N. *et al.* Mitochondrial evidence for multiple radiations in the evolutionary history of small apes. *BMC Evolutionary Biology* **10**, 1–13 (2010).
- [15] Penzkofer, T. *et al.* L1base 2: more retrotransposition-active line-1s, more mammalian genomes. *Nucleic Acids Research* **45**, D68 (2017).
- [16] Durbin, R., Eddy, S. R., Krogh, A. & Mitchison, G. *Biological sequence analysis: probabilistic models of proteins and nucleic acids* (Cambridge university press, 1998).

- [17] Hubley, R. *et al.* The dfam database of repetitive dna families. *Nucleic acids research* **44**, D81–D89 (2016).
- [18] Kimura, M. & Weiss, G. H. The stepping stone model of population structure and the decrease of genetic correlation with distance. *Genetics* **49**, 561 (1964).
- [19] Li, M., Vitányi, P. *et al.* *An introduction to Kolmogorov complexity and its applications*, vol. 3 (Springer, 2008).
- [20] Dingle, K., Camargo, C. Q. & Louis, A. A. Input–output maps are strongly biased towards simple outputs. *Nature communications* **9**, 1–7 (2018).
- [21] Seiler, M. *et al.* Somatic mutational landscape of splicing factor genes and their functional consequences across 33 cancer types. *Cell reports* **23**, 282–296 (2018).
- [22] Bao, W., Kojima, K. K. & Kohany, O. Repbase update, a database of repetitive elements in eukaryotic genomes. *Mobile Dna* **6**, 1–6 (2015).
- [23] Jurka, J. *et al.* Repbase update, a database of eukaryotic repetitive elements. *Cytogenetic and genome research* **110**, 462–467 (2005).
- [24] Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of rna-seq data. *Genome biology* **11**, 1–9 (2010).
- [25] Ritchie, M. E. *et al.* limma powers differential expression analyses for rna-sequencing and microarray studies. *Nucleic acids research* **43**, e47–e47 (2015).
- [26] Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)* **57**, 289–300 (1995).
- [27] Mehdipour, P. *et al.* Epigenetic therapy induces transcription of inverted sines and adar1 dependency. *Nature* **588**, 169–173 (2020).
- [28] Ahmad, S. *et al.* Breaching self-tolerance to alu duplex rna underlies mda5-mediated inflammation. *Cell* **172**, 797–810 (2018).
- [29] LaCava, J., Jiang, H. & Rout, M. P. Affinity capture from cryomilled mammalian cells. *Journal of Visualized Experiments* **118**, 54518 (2016).
- [30] Taylor, M. S. *et al.* *Characterization of L1-Ribonucleoprotein Particles*, 311–338 (Springer New York, New York, NY, 2016). URL https://doi.org/10.1007/978-1-4939-3372-3_20.
- [31] Di Stefano, L. H. *et al.* *Affinity-Based Interactome Analysis of Endogenous LINE-1 Macromolecules*, 215–256 (Springer US, New York, NY, 2023). URL https://doi.org/10.1007/978-1-0716-2883-6_12.
- [32] Cristea, I. M. & Chait, B. T. Conjugation of magnetic beads for immunopurification of protein complexes. *Cold Spring Harbor Protocols* **2011**, pdb–prot5610 (2011).
- [33] Ardeljan, D. *et al.* Line-1 orf2p expression is nearly imperceptible in human cancers. *Mobile dna* **11**, 1–19 (2020).
- [34] Jiang, H., Lei, R., Ding, S.-W. & Zhu, S. Skewer: a fast and accurate adapter trimmer for next-generation sequencing paired-end reads. *BMC Bioinformatics* **15**, 182 (2014).

- [35] Yang, W. R., Ardeljan, D., Pacyna, C. N., Payer, L. M. & Burns, K. H. SQUIRE reveals locus-specific regulation of interspersed repeat expression. *Nucleic Acids Research* **47**, e27–e27 (2019). URL <https://doi.org/10.1093/nar/gky1301>. https://academic.oup.com/nar/article-pdf/47/5/e27/28041594/gky1301_supplemental_files.pdf.
- [36] Dobin, A. *et al.* Star: ultrafast universal rna-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- [37] Liao, Y., Smyth, G. K. & Shi, W. featurecounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- [38] Wickham, H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag New York, 2016). URL <https://ggplot2.tidyverse.org>.

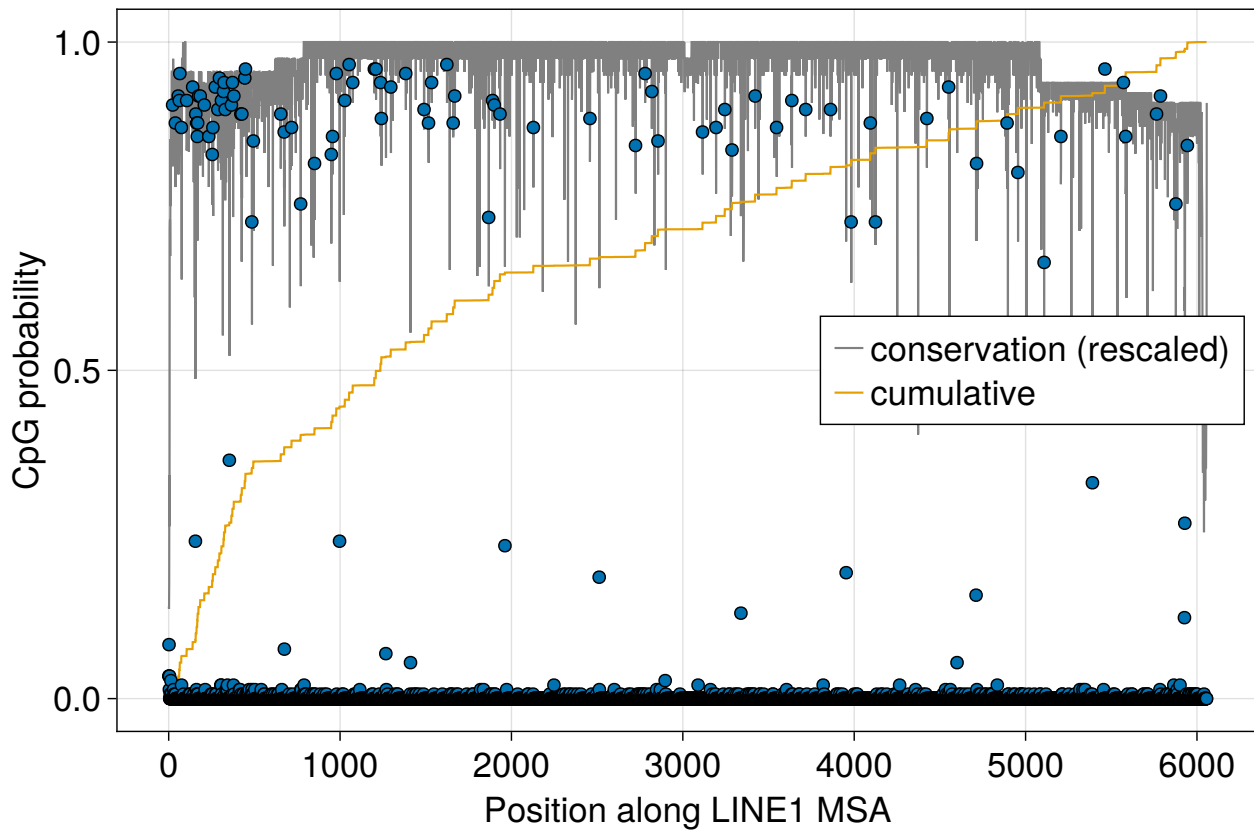
Supplementary Material

Supplementary Tables

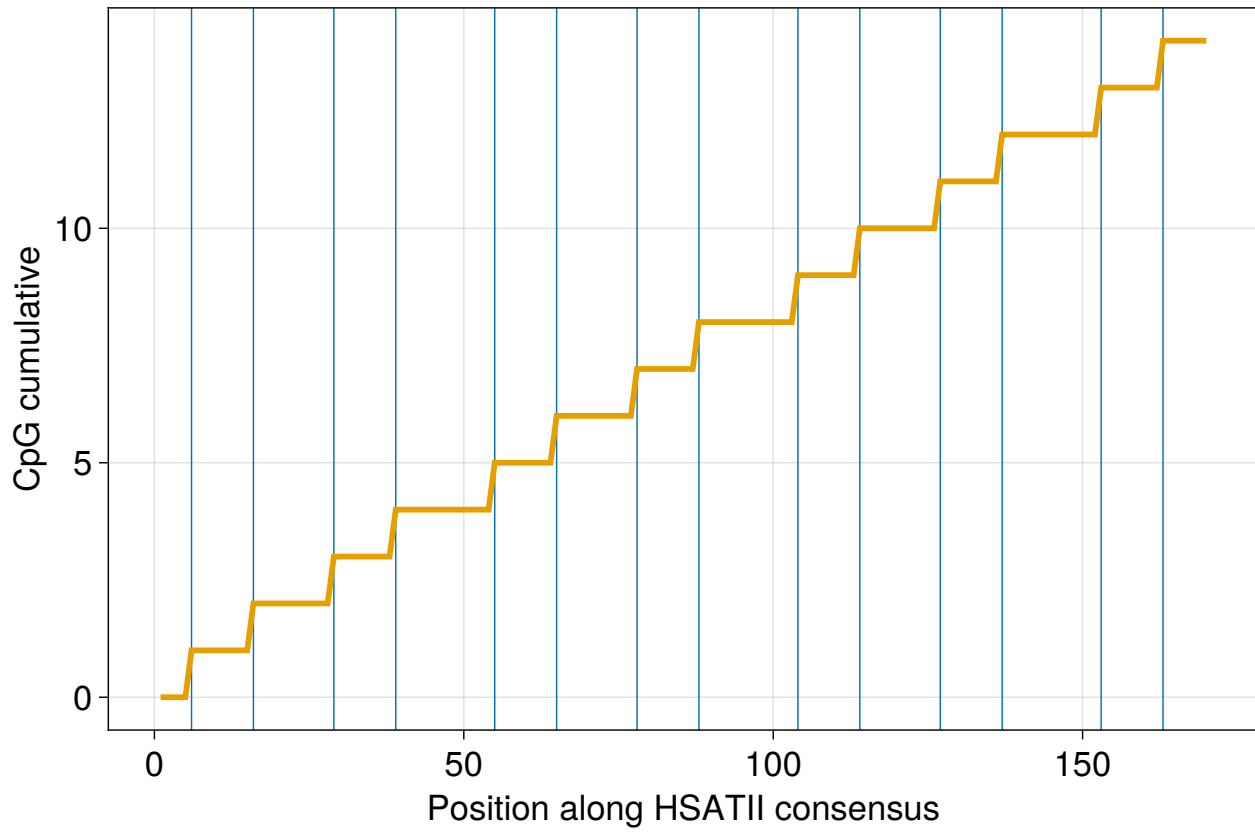
Supplementary Tables are available separately in Comma Separated Value (csv) format. The description of the respective tables is provided below:

- Supplementary Table 1 contains the x_{CpG} and x_{ds} values calculated for each repeat family annotated in the DFAM database. For each family, it includes the forces of the consensus sequence and the mean calculated for the inserts in the human genome.
- Supplementary Table 2 contains the list of high- x_{CpG} ($x_{\text{CpG}} > 0$) windows of 3 kb detected in the human genome, after filtering so that windows do not overlap more than 1 kb. The table also includes information about the repeat that maximally overlaps with each window as annotated in DFAM. Finally we report the most-recent common ancestor time of the primates for which we observed high- x_{CpG} sequences alignable by BLAST with each high- x_{CpG} window in the human genome.
- Supplementary Table 3 contains the list of the unique double-stranded segments resulting in windows with high- x_{ds} ($x_{\text{ds}} > 0.5$) detected in the human genome. The processing of data is described in Methods. The table also includes information about the repeats that maximally overlaps with each of the two double-strand forming segments as annotated in RepeatMasker. Finally we report the most-recent common ancestor time of the primates for which we observed pair of segments resulting in high- x_{ds} and alignable by BLAST with each pair of high- x_{ds} segments in the human genome.
- Supplementary Table 4 contains the list of the unique double-stranded segments resulting in windows with high- x_{ds} ($x_{\text{ds}} > 0.5$) detected in the zebrafish genome (danRer11).
- Supplementary Table 5 contains the results of differential expression analysis on all repetitive elements between treated versus untreated SF3B1 mutant samples.
- Supplementary Table 6 contains the list of annotated samples subject to L1 ORF1p/ORF2p RIP-seq analysis. Salt concentration, cell type, antibody type, replicate and sequencing details are listed.

Supplementary Figures



Supplementary Figure S1: CpG occurrence probability in a multiple sequence alignment of full-length intact LINE1 inserts in the human genome. The orange line denotes the cumulative probability, and the gray line denotes the conservation (in bits) rescaled between 1 and 0: $1 + \sum_{\sigma} f_i(\sigma) \log_2(f_i(\sigma)) / \log(5)$, where $f_i(\sigma)$ is the fraction of times the nucleotide σ (or a gap) appears in position i .



Supplementary Figure S2: CpG occurrence (vertical blue lines) in the consensus HSATII sequence as annotated in DFAM. The orange line denotes the cumulative number of CpGs.