



Citation for published version:

Radrizzani, S, Kudla, G, Izsvák, Z & Hurst, LD 2024, 'Selection on synonymous sites: the unwanted transcript hypothesis', *Nature Reviews Genetics*. <https://doi.org/10.1038/s41576-023-00686-7>

DOI:

[10.1038/s41576-023-00686-7](https://doi.org/10.1038/s41576-023-00686-7)

Publication date:

2024

Document Version

Peer reviewed version

[Link to publication](#)

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 Selection on synonymous sites: the unwanted transcript hypothesis

2
3 Sofia Radrizzani^{1,2}, Grzegorz Kudla³, Zsuzsanna Izsvák⁴ and Laurence D. Hurst^{1†}

4
5 ¹Milner Centre for Evolution, Department of Life Sciences, University of Bath, Bath, UK.

6 ²Milner Therapeutics Institute, Jeffrey Cheah Biomedical Centre, University of Cambridge,
7 Cambridge, UK.

8 ³MRC Human Genetics Unit, Institute for Genetics and Cancer, The University of Edinburgh,
9 Edinburgh, UK.

10 ⁴Max-Delbrück-Center for Molecular Medicine in the Helmholtz Society, Berlin, Germany.

11 †e-mail: bssldh@bath.ac.uk

12 13 Abstract

14 Although translational selection to favour codons that match the most abundant tRNAs is not
15 readily observed in humans, there is nonetheless selection in humans on synonymous
16 mutations. We hypothesise that much of this synonymous-site selection can be explained in
17 terms of protection against unwanted RNAs — spurious transcripts, mis-spliced forms or
18 RNAs derived from transposable elements or viruses. We propose that selection on
19 synonymous sites not only acts to reduce the rate of creation of unwanted transcripts (e.g. via
20 selection on exonic splice enhancers and cryptic splice sites), but, additionally, that high-GC
21 content (but low CpG content), along with intron presence and position, is both particular to
22 functional native mRNA and employed to recognize transcripts as native. Evidence supports
23 this hypothesis with transcription, nuclear export, liquid phase condensation and RNA
24 degradation all recently being shown to promote GC rich transcripts and suppress AU/CpG
25 rich ones. With traps set against AU/CpG rich transcripts, codon usage of native genes in turn
26 evolves to avoid such suppression. That parallel filters against AU/CpG rich transcripts also
27 affect the endosomal import of RNAs further supports the unwanted transcript hypothesis of
28 synonymous site selection and explains the similar design rules that have enabled the
29 successful use of transgenes and RNA vaccines.

30 31 [H1] Introduction

32 As they do not alter the amino acid that is specified, **synonymous mutations [G]** were originally
33 assumed to be **neutrally evolving [G]**¹. However, further analysis led to the well-accepted
34 notion that selection favours synonymous mutations forming codons that partner with the most
35 abundant iso-acceptor tRNAs in that organism²⁻⁴. This is thought to reflect selection operating
36 on translational velocity⁵, possibly owing to shorter ribosomal dwell times of ribosomes on the
37 mRNA **[Au:OK?]** (although other studies dispute this⁶), or on translational accuracy^{7,8}. Such

38 translational selection [G] is more pronounced in fast growing species⁹ and in those with large
39 effective population sizes [G] (N_e), these being subject to more efficient purifying selection to
40 remove deleterious alleles¹⁰ (**Box 1**).

41

42 Within species, translational selection is more pronounced in highly expressed genes, and is
43 thus typically inferred from stronger codon usage bias [G] in such genes and from positive
44 correlations between codon usage and tRNA copy numbers². In humans, as expected of a
45 species with small N_e (ref.¹¹) (**Box 1**), these hallmarks of translational selection are not readily
46 observed^{2,12}. It was thus supposed in early discussions² that selection does not operate on
47 human synonymous mutations nor cause genetic disease.

48

49 However, this supposition has gradually changed (as reviewed by refs^{13,14}). Although there is
50 still no evidence for commonplace translational selection in humans¹⁰, the supposition that
51 codon usage is determined solely by translational selection is now recognised to be wrong.
52 For example, increasing GC content at synonymous sites of transgenes in mammals leads to
53 increased protein levels but not through increased rates of translation¹⁵⁻¹⁷. Estimates of the
54 proportion of human synonymous sites that are under selection vary from a few percent^{18,19} to
55 between 10% and 20%²⁰⁻²². The largest database reports thousands of synonymous mutations
56 'associated' with ~1,500 diseases²³. There is now evidence for selection on synonymous
57 codon usage throughout the gene expression pathway, from transcription onwards, as we
58 discuss here.

59

60 Rather than simply review this evidence for selection on synonymous mutations, we ask why
61 such selection might be so apparently pervasive. More particularly, whereas the translational
62 selection model has dominated discussion of selection on synonymous codon usage for the
63 past half century, is there an alternative framework that can explain this selection in humans?
64 Here, we present the unwanted transcript hypothesis.

65

66 We argue that with low N_e and the resulting high proportion of 'junk' DNA in the human genome
67 (**Box 1**), the main problem for gene expression in humans is the presence of unwanted RNA
68 transcripts. Unwanted transcripts are defined as those for which their removal is
69 advantageous, even if the process of their creation is itself beneficial (for example, via
70 transcriptionally mediated modification of chromatin²⁴). We propose that much of the selection
71 on synonymous sites curtails the creation of unwanted transcripts (through effects on
72 transcription and splicing) or traps them to prevent translation (**Fig. 1**). These mechanisms of
73 selection use generic patterns, often nucleotide skews, present in non-functional or foreign
74 RNA sequences that discriminate them from functional native sequences. Understanding this

75 form of selection, and the resulting trends in synonymous codon usage, is important for
76 diagnostics and for the design of transgenes.

77

78 **[H1] The unwanted transcript hypothesis**

79 We start by laying out what we consider to be the problem of unwanted transcripts, arguing
80 that they will commonly be costly and thus provide the context for their suppression. We then
81 discuss possible solutions that might be broadly classified as either reduction in the rate of
82 their creation or the setting of traps. Both processes are affected by codon identity and thus
83 can be impacted by synonymous mutations.

84 **[H2] The problem**

85 There are multiple sources of unwanted transcripts. **Spurious transcripts [G]** are likely to be
86 common in humans as eukaryotic transcription factor (TF) binding sites are small, degenerate
87 and thus common²⁵. Of the 15.6 million TF binding sites in the human genome, fewer than
88 13% have sequence conservation²⁶. Thus, novel sequences could be a rich source of
89 transcripts. Two recent preprints report that naïve (i.e. not native) sequences introduced into
90 yeast (with mostly open chromatin) result in large amounts of spurious transcripts, many of
91 which are rare^{27,28}. Analysis of random sequences of 120 nucleotides integrated into yeast
92 reveals that half have promoter activity, and only 1–5% of intergenic transcripts in yeast cannot
93 be attributed to similar spurious transcriptional activity²⁹.

94

95 Some spurious transcription may be unavoidable. Most mRNA promoters are inherently
96 bidirectional, with transcription upstream of transcription start sites commonly generating non-
97 functional transcripts³⁰ known as promoter upstream transcripts (PROMPTs) that are
98 degraded, for example by the nuclear RNA **exosome complex [G]**. Other types of potentially
99 non-functional transcripts include long noncoding RNAs (lncRNAs), enhancer RNAs and
100 natural antisense transcripts. Although some have documented functions, many (especially
101 the rare transcripts^{31,32}) are poorly conserved in sequence^{31,33} and presence³⁴, and it seems
102 likely that many of them are functionless³⁵. Indeed, many are destroyed so soon after creation
103 that they are defined as cryptic transcripts³⁶, visible only after the degradation machinery is
104 removed. Similarly, lncRNA transcription commonly terminates prematurely, co-transcriptional
105 splicing is less efficient than for protein-coding transcripts and lncRNA transcripts tend to be
106 unstable with expression levels 10-fold lower than for mRNAs (reviewed in ref.³⁵). Functional
107 studies seem to reinforce the conclusion that few lncRNAs have evident function^{37,35}. lncRNA
108 transcripts that are rare and rapidly degraded are less likely to have any RNA sequence-
109 dependent functionality³⁸.

110

111 Much transcription is either parasitic or associated with remnants of parasitic sequences (such
112 as viruses and **transposable elements [G]**) as their successful genomic colonization requires
113 them to have TF binding sites³⁹. Even though only about 1% of transposable elements in the
114 human genome can still transpose, they are likely to be a rich source of unwanted transcripts.
115 Indeed, primate-specific TF binding sites tend to be unconserved, found within transposable
116 elements and not identified in genome-wide association studies and hence with no evident
117 phenotypic effects²⁶. More than 85% of primate-specific TF binding sites, and more than 20%
118 of all TF binding sites, are derived from transposable elements²⁶. Although some transposable
119 elements have become domesticated to serve novel functions in the host cell (for
120 example^{40,41}), many transposable element-derived transcripts are likely to be spurious. Thus,
121 humans may have particularly high levels of spurious transcription, as more than half of human
122 DNA is derived from transposable elements¹¹ (**Box 1**). These processes likely explain why
123 nearly all the DNA of the human genome is transcribed⁴² and why 83% of lncRNAs have
124 exonised transposable elements⁴³.

125

126 Immature functional transcripts can also pose a problem in humans as splicing is error
127 prone^{44,45}. This can be owing to mutation or to the intrinsic nature of the process, with human
128 transcripts often having weak splice sites⁴⁶. Experimental reports of the proportion of exonic
129 point mutations that disrupt splicing vary from 20% to nearly 100% (reviewed in ref.⁴⁷). Just as
130 small N_e is associated with weak translational selection, a recent preprint reports that it is
131 associated with increased transcriptome diversity⁴⁵. Alternative splicing is increased when N_e
132 is small and, importantly, this is largely owing to low-abundance splice forms, many of which
133 are frameshifted and not associated with selection to preserve the splice site⁴⁵. The increase
134 in splice forms when N_e is small is thus more likely owing to weaker selection on splicing⁴⁵
135 (**Box 1**) rather than reflecting functional splice forms associated with increasing organismal
136 complexity.

137

138 ***[H2] Is expression of unwanted transcripts costly?***

139 A key assumption of the unwanted transcript hypothesis is that such transcripts are costly and
140 hence that there is selection against them. Transcription and translation of unwanted RNAs
141 are expected to be costly because both, especially translation⁴⁸, are energetically
142 demanding⁴⁸. In addition, just as the translation of functionally irrelevant reporter genes
143 reduces cellular fitness^{49,50} by engaging ribosomes that could otherwise process required
144 transcripts, so too we expect the translation of unwanted RNAs to be deleterious. Depending
145 on the nature of the unwanted transcripts, there may be additional costs associated with their
146 expression, such as direct toxicity⁵¹, interference with RNA metabolism or eliciting of an

147 immune response⁵² (e.g. via inappropriate triggering of viral sensing RIG-1-like receptors
148 leading to secretion of type 1 interferons⁵³).

149

150 One line of evidence that unwanted transcription can be costly comes from the association
151 between mutations in pathways responsible for the removal of such transcripts (see also
152 below) and human genetic diseases⁵³⁻⁵⁵. For example, mutations in *RBM7*, part of the nuclear
153 exosome targeting (NEXT) complex that directs RNA to the nuclear exosome complex for
154 degradation, result in spinal motor neuropathy⁵⁶. Mutations in cyclin-dependent kinase 13
155 (CDK13), which modulate RNA stability, are associated with melanoma⁵⁷. This reflects the role
156 of CDK13 in phosphorylation of ZC3H14, which in turn is required for nuclear RNA
157 degradation. Failure to activate nuclear RNA surveillance owing to CDK13 mutation results in
158 aberrant prematurely terminated transcripts (ptRNAs) that are translated and, in turn, promote
159 melanoma progression⁵⁷. Mutations in genes encoding nuclear RNA surveillance components
160 are seen in many malignancies⁵⁷; a curiosity in this case is that the production and translation
161 of ptRNAs increases cellular growth rates.

162

163 For transpositionally competent transposable elements, additional costs will be associated
164 with damage to host DNA through insertion and double-strand breaks that are mutagenic, can
165 cause break-induced genome instability and are toxic⁵². Furthermore, transposable elements
166 can interfere with host mRNAs through multiple mechanisms and the accumulation of
167 transposable element transcripts can trigger immune responses (reviewed in ref.⁵²).
168 Unsurprisingly, there is antagonistic coevolution between transposable elements and
169 mechanisms to suppress them (for example, see refs^{58,59}), which is itself indicative of costs
170 associated with their unwanted activity. Transposable elements are silenced by multiple
171 mechanisms, notably DNA and histone methylation. Particular to mammals are nuclear
172 paraspeckles, ribonucleoprotein bodies that function as traps for double-stranded RNA (typical
173 of viruses) and transcripts with **inverted repeat Alu elements [G]** in the 3' untranslated region
174 (UTR), thereby preventing their nuclear export⁶⁰.

175

176 The importance of these mechanisms to suppress transposable elements is suggested by
177 pathologies that are associated with the impairment of a specific repression mechanism (for
178 example, see refs⁶¹⁻⁶⁵) and by the increased transcription of transposable elements in
179 autoimmune diseases, cancer and neurodegenerative pathologies^{61,66-69}. Typically, these
180 pathological conditions worsen with age or manifest with aging (such as Alzheimer disease,
181 age-related macular degeneration and tauopathies), and aging itself is associated with
182 increased transposable element transcription and activity (for example of the SIRT6–**LINE1**
183 **[G]** (L1) silencer⁷⁰).

184

185 However, it is often hard to be certain that some of the apparent costs of unwanted transcripts
186 are direct costs rather than an incidental correlate. For example, transposable element
187 upregulation may be an irrelevant consequence of weaker chromatin control or there may be
188 unrelated pleiotropic effects of mutations in genes associated with RNA metabolism and
189 surveillance, including a dearth of functional transcripts that are regulated by the same quality
190 control mechanisms. In some cases, the evidence for causality of the unwanted transcripts is
191 robust. For example, introduction of either mutant CDK13 or phosphorylation-modifying
192 mutations of ZC3H14 promote tumour progression and both are associated with the
193 generation of unusual transcripts⁵⁷.

194

195 In the case of splicing disruption, the cause–effect relationship can be clearer as there can be
196 mutations not only in the RNA-binding proteins (such as serine–arginine-rich (SR) proteins)
197 that coordinate splicing (which are long recognised as being associated with disease⁷¹) but
198 also in the target transcripts. For example, mutations in **exonic splice enhancers [G]** (ESEs)
199 in *BRCA1* cause mis-splicing and disease⁷². Recent evidence from analysis of the site
200 frequency spectrum suggests that there is strong purifying selection in humans against
201 synonymous mutations that disrupt ESEs²⁰.

202

203 Several other points of evidence would be hard to understand were there not a cost to
204 unwanted transcripts. For example, the increase in apparently functionless splice forms in
205 species with small N_e is most parsimoniously explained by inefficient selection against weakly
206 deleterious splice forms⁴⁵. It is also notable that the excess transcripts seen in species with
207 small N_e are predominantly rare transcripts⁴⁵. This can be explained by a balance between
208 selection against mis-splicing and the costs associated with mis-spliced transcripts, which we
209 assume to be lower for rarer transcripts. The fact that noise in gene expression is higher for
210 the rarer transcripts⁴⁵ is also consistent with a failure of selection to operate on weakly
211 deleterious, rare mis-spliced transcripts. Similarly, it is notable that in yeast, nearly all bases
212 of naive sequence are transcribed, generating multiple rare transcripts, whereas native
213 sequence seems to have much ‘cleaner’ transcription, with rare transcripts largely suppressed
214 and functional transcripts at high abundance^{27,28}. This would be the expected outcome of
215 selection against costly spurious transcription.

216

217 **[H2] The solutions**

218 Assuming that unwanted transcripts are on average costly, there are two modes by which
219 organisms can address this problem. One is to curtail their creation, which would involve either
220 suppression of transcription or improved accuracy of transcript processing, notably splicing.

221 Gene body methylation in humans is, for example, necessary to reduce rates of spurious
222 intragene transcriptional initiation from cryptic promoters⁷³. The other is to have a quality
223 control system to filter transcripts. Quality control could involve either active degradation (for
224 example, by RNAses⁷⁴) or physical isolation (for example, preventing nuclear export).

225

226 Humans have multiple such quality control filters throughout the gene expression pathway
227 (**Fig. 1**). For example, ptRNAs and upstream antisense RNAs are targeted for nuclear
228 degradation by the polyA exosome targeting (PAXT) complex^{75,76}, containing MTR4 and
229 ZFC3H1. MTR4 is also involved in the NEXT complex, which targets early and unprocessed
230 RNAs for destruction⁷⁷. The physical separation of translation (in the cytoplasm) and
231 transcription (in the nucleus) underpins much translational suppression; for example, retention
232 in the nucleus of ADAR-edited pairs of inverted repeat Alu elements⁷⁸. Even if a transcript can
233 escape the nucleus, there are further traps (such as cytoplasmic **processing bodies [G]**) and
234 mechanisms of translational⁷⁹ or post-translational suppression, including **nonsense-mediated**
235 **decay [G]** (NMD), **no-go mRNA decay [G]**, **codon optimality-mediated RNA decay [G]** and
236 **non-stop RNA decay [G]**. Cytoplasmic stress granules are considered a triage centre⁸⁰ for
237 differentiating wanted from unwanted transcripts. Multiple systems selectively target non-
238 native transcripts for destruction (such as DICER⁸¹, **zinc finger antiviral protein [G]** (ZAP)⁸²
239 and RNases⁷⁴). Even after translation has occurred, there are additional error checks such as
240 degradation of proteins with a hydrophobic carboxyl terminus⁸³ as these are more likely to be
241 products of spurious transcription.

242

243 *N*⁶-methyladenosine (m⁶A) deposition on RNA has a complex involvement in the quality control
244 process. It is associated with retrotransposon suppression⁸⁴, accumulation in stress granules
245 of potentially foreign, long-exon transcripts⁸⁵ and nuclear retention and decay when on introns
246 of transcripts with unprocessed 5' splice sites⁸⁶. Conversely, small exons typical of native
247 genes are protected from m⁶A deposition by the **exon-junction complex [G]**⁸⁷. Although this
248 implicates m⁶A deposition in RNA suppression, its effects are sensitive to the 'reader' proteins
249 with which it is associated (reviewed in ref.⁸⁸). It can, for example, also enable nuclear export
250 of endogenous transcripts with few or large exons¹⁷ and release from transcriptional
251 suppression of a domesticated endogenous retrovirus⁸⁹.

252

253 Our hypothesis is that much of the selection on synonymous sites reflects selection for
254 reduced production or improved quality control of unwanted transcripts. The core problem of
255 any such system is how to differentiate wanted from unwanted transcripts. This differentiation
256 must depend on the recognition of 'fingerprints' that would not be expected to occur in a
257 random sequence or in a newly arrived sequence (for example, a new retrovirus). Similar logic

258 is recognized in the context of innate immunity in which conserved pathogen-associated
259 molecular patterns are recognised by pattern-recognition receptors such as Toll-like receptors
260 (TLRs) to trigger an immune response⁹⁰. In this language, we propose that humans have a
261 similar set of pattern-recognition receptors for unwanted transcripts. At the level of
262 retroviruses, innate immunity and selection against unwanted transcripts combine, with the
263 HUSH complex involved in transcriptional silencing being considered a pattern-recognition
264 receptor (**Box 2**).

265

266 As introns are especially abundant in mammals, and vertebrates generally⁹¹, their presence
267 (and positioning) is indicative of a native transcript. Conversely, longer exons are rare in
268 humans (80% of human exons are less than 200bp⁹²) and potentially indicative of foreign
269 transcripts. Intron presence is indeed known to inform quality control. For example, the human
270 NMD system uses location relative to an exon-junction complex to infer prematurity⁹³, and
271 introns have long been recognised as important for human transgene expression⁹⁴. The most
272 'suspicious' transcripts are therefore transcripts with long exons and intronless transcripts.
273 Also, humans cannot rely on intron-mediated NMD for quality control in intronless transcripts⁹⁵.
274 However, as some intronless genes are derived from retrotransposition, and thereby inherit
275 the exon splice motifs of the parental multi-exon gene, SR protein binding to such motifs has
276 splicing-independent activity⁹⁶ that could be used for quality control (for example, nuclear
277 export related).

278

279 Importantly in the context of selection on codon usage, we suggest that nucleotide content
280 provides a further useful guide to wanted versus unwanted transcripts. As the rate of GC to
281 AT mutation is higher than the converse (known as a mutation bias)⁹⁷, a sequence that is
282 selectively neutral will evolve low-GC content⁹⁷. Human AT-rich gene deserts are close to
283 **mutational equilibrium [G]**⁹⁸ and consequently the dinucleotide content of the genome as a
284 whole is skewed towards A and T (**Fig. 2**). High-GC content can thus be an indication of
285 functional constraint, because proteins will commonly need amino acids whose underlying
286 codons are G/C-rich, there will be a dearth of A/T-rich stop codons, and A/T-rich codons tend
287 to specify amino acids that are more metabolically expensive to manufacture (i.e. in terms of
288 ATP consumption) and are thus counter-selected⁹⁹. In addition to mutation bias, A enrichment
289 of single-exon transcripts reflects the reliance of retroviruses on structurally poor A-rich RNA
290 necessary for cDNA generation by reverse transcription¹⁰⁰. We thus suggest that GC richness
291 at non-synonymous sites naturally marks native transcripts as opposed to spurious transcripts
292 from AT-rich intergenic sequences, and that synonymous site usage has evolved to reinforce
293 the high-GC content. In particular, third sites of a codon are the freest to evolve without
294 consequence, so are a preferable solution for indicating that a protein coding gene is native.

295 As would be assumed on this basis, human coding sequences are most enriched, compared
296 with genomic mean dinucleotide content, for G/C dinucleotides and most lacking in A/T
297 dinucleotides (**Fig. 2**).

298

299 The idea that high-GC content is evidence of functionality is complicated by the fact that **GC-**
300 **biased gene conversion [G]** (gBGC) acts in the opposite direction to the mutation bias¹⁰¹.
301 Thus, sequence that is GC-rich, especially at third codon sites, could be subject to gBGC
302 rather than selection¹⁰¹. However, this observation makes sense if humans use GC content as
303 a signal of nativity: one reason that gBGC is so biased in humans — but not yeast¹⁰² — could
304 be that it preserves beneficial high-GC content in the face of weak selection in mammals to
305 do the same. As gBGC in humans acts predominantly on gene-rich domains¹⁰¹, native genes
306 acquire high-GC content at third codon sites (GC3) whereas spurious transcripts from the
307 larger AT-rich part of the genome do not¹⁰¹. Additionally, with paralogs often in clustered arrays
308 inter-locus recombination and associated gBGC can boost their GC3 content, as for example
309 witnessed with histone genes¹⁰³. It is unclear whether the high-GC3 content of human exons
310 is entirely owing to gBGC as most exons have GC3 content higher than that of the flanking
311 introns¹⁰⁴. Although this can be explained in some part as being owing to AT-rich transposable
312 elements being more abundant in introns than in coding sequence¹⁰⁵, a difference between
313 GC3 content in exons and introns remains after transposable elements are masked¹⁰⁵.

314

315 In our proposal then, the evolution of suppression mechanisms against A/U rich transcripts,
316 favouring G/C rich ones, both forces selection on synonymous mutations and acts on the
317 evolution of the direction of gBGC's bias. This bias along with selection, given the GC
318 preference, then act to force synonymous site GC content above that expected at mutation-
319 selection equilibrium, this being highly AT biased. More generally, we expect that all quality
320 control mechanisms will co-evolve with native genes synonymous site nucleotide content.

321

322 The exception to high-GC content indicating native sequence [**Au:OK?**] is CpG, as methylated
323 CpG is hypermutable. Indeed, although CpG is more common in coding sequences than in
324 the genome as a whole (**Fig. 2**), CpG is the only C/G dinucleotide below the frequency that
325 would be expected given mononucleotide occurrence, and it has low frequency in coding
326 sequences, with only TpA being less frequent (**Fig. 2**). This same hypermutability forces
327 increased levels of TpG and CpA, which causes, for statistical reasons, under-representation
328 of both CpG and TpA¹⁰⁶. At the transcript level, high-CpG or high-UpA content is therefore
329 indicative of foreign transcripts. As viral sequences also need to encode amino acids, their
330 sequences could have high-G/C content but also high CpG (although some viruses, such as
331 SARS-CoV-2, have such extreme mutation bias favouring U that they are GC poor despite

332 selection operating against the mutation bias¹⁰⁷). For similar reasons, CpG should be used as
333 a suppressor of transcription (for example of transposable elements) mediated by methylation.

334

335 In sum, we expect that GC-rich and CpG-poor (including at synonymous sites) transcripts will
336 pass pre-translation quality control tests, whereas high CpG, high U/T and high-A content,
337 including high UpA, will lead to activity that prevents eventual translation. Similarly, intron
338 position and presence should be used as a guide to nativity. Intronless transcripts and those
339 with especially large exons should thus be regarded with suspicion. GC-rich multi- and small-
340 exon genes are top candidates for nativity, whereas AT-rich single- or long-exon transcripts
341 are most likely to be spurious. Functional viral sequences are likely to be caught by small-
342 exon, multi-exon, low-CpG, high-GC3 traps. In addition, we expect selection against
343 synonymous mutations to occur if these generate spurious splice forms of native transcripts.

344

345 **[H1] What can the hypothesis explain?**

346 Here, we note that the unwanted transcript hypothesis can explain many previously
347 unconnected observations, with nucleotide content being a guide to transcript nativity such
348 that high-GC content promotes gene expression, whereas high-CpG and high-AT content
349 suppress expression.

350

351 **[H2] Suppression of transcription**

352 In contrast to the prediction of the translational selection model, codon usage bias does affect
353 transcription rates. The first such evidence was derived from *Neurospora crassa*¹⁰⁸, in which
354 codon usage bias predicted both protein levels and intra-nuclear RNA levels. Similar results
355 have been replicated in numerous species, including humans¹⁰⁹, with various mechanisms
356 involved including histone 3 lysine 9 trimethylation (H3K9me3) and transcription factor
357 binding¹¹⁰. Most pertinent, however, seems to be the effects of 5' GC content on RNA
358 polymerase II (Pol II) processivity. After transcription initiation, processive elongation by Pol II
359 is necessary for completion of transcription rather than premature transcript termination that
360 generates micro 5' transcripts. In mammals, two features promote Pol II processivity: GC-rich
361 5' codons^{109,111} and a prospective splice site early in the transcript¹¹¹. We presume that
362 removing micro transcripts is preferable to the consequences of potentially extending
363 transcription of non-functional sequences.

364

365 Conversely, AT-rich genes are bound by the zinc-finger protein SALL4, which prevents
366 premature expression of such genes¹¹². Intronless transcripts are considered to be most
367 'suspicious' and it is notable that humans have a specialist system (the HUSH complex) to
368 mediate silencing of long intronless transcripts (or those with unusually long exons) that are

369 especially A-rich¹⁰⁰ (**Box 2**). Less than 0.4% of native exons, in about 3.5% of genes, have
370 high enough A content and are long enough (>1500bp) to be affected by HUSH (**Fig. 3**). About
371 a quarter of the affected genes encode zinc finger proteins, which suggests that HUSH and
372 zinc finger proteins might be alternative suppression systems, with the latter inactivated by the
373 former.

374

375 As expected, given that CpG is rare in native human coding sequences, the default state in
376 the non-coding genome is transcriptional inactivation via CpG methylation^{113,114}, which in part
377 is an anti-transposable element strategy.

378

379 **[H2] Prevention of mis-splicing**

380 The best studied mode of selection on synonymous sites in humans is to enable accurate
381 splicing. Humans have strong selection on synonymous mutations to preserve ESEs^{20,115,116}
382 near exon ends and such ESE disruption is associated with disease¹¹⁷. Similarly, selection on
383 synonymous mutations reflects selection against ESEs in unwanted (core) exonic locations¹¹⁸.
384 Unusual codon usage trends towards the ends of human exons can then be explained by
385 selection for synonymous sites that promote ESE incorporation^{19,119,120}. If codon usage
386 selection at exon ends needs to be balanced between accurate splicing and faster translation
387 (in species with translational selection), then accurate splicing dominates¹²¹. The avoidance
388 of **cryptic splice sites [G]** in proximity to exon junctions is similarly advantageous¹²² and
389 synonymous mutations that generate them are deleterious as they generate unwanted
390 transcripts¹²³.

391

392 **[H2] Quality control**

393 Broadly, in supposing that high-GC3 content enables quality control filters to be passed, this
394 model can explain why high-GC content at synonymous sites enables the expression of
395 intronless genes^{15,16}. It also resolves the enigma¹²⁴ of why non-recombining (hence not subject
396 to gBGC) intronless *Sry* has an unusually high synonymous GC3 content. Additionally, it can
397 explain why native intronless **retrogenes [G]** have higher GC3 content than their parental
398 (intron-containing) genes¹⁶. It is unknown to what extent a net flux of retrogenes¹²⁵ from the
399 AT-rich X chromosome to the GC-richer autosomes might explain the same trend.

400

401 More generally, for a gene to be expressed in humans it needs to have an intron near the 5'
402 end or be intronless but GC-rich¹⁶. Inspired by the finding that retrogenes have high-GC3
403 content, one study constructed many intronless variants encoding green fluorescent protein
404 that differed in GC3 content and found that the higher the GC3 content, the higher the level of
405 the protein expressed¹⁶. Both intronic and intronless human genes tend to be GC-rich at their

406 5' ends¹²⁶. This contrasts with bacteria, in which the 5' ends of sequences are AT rich, which
407 leads to low RNA stability and thus easier translation initiation^{49,127}. Although low 5' RNA
408 stability is a common feature across much of the evolutionary tree, it does not force high AT
409 content in mammals¹²⁸.

410

411 The increased protein expression of GC-rich transcripts described in GFP transgenes¹⁶ likely
412 involved both quality control measures and transcriptional effects. The most obvious location
413 for any form of 'passport' quality control of mRNA is nuclear export. Indeed, this study identified
414 that, at least in part, the GC-mediated expression effect is owing to increased relative
415 cytoplasmic presence of GC3-rich transcripts. This suggests a role for either nuclear export,
416 consistent with prior evidence¹²⁹, or some other intra-nuclear filter in restricting cytoplasmic
417 access of GC-poor transcripts.

418

419 Further studies subsequently identified the nuclear export pathway as being core to such
420 observations¹⁷. Of the two canonical nuclear RNA export pathways, NXF1 was required for
421 the export of single-exon transcripts, transcripts with long exons and A/U-rich, multi-exon
422 transcripts¹⁷, whereas the TREX complex facilitated export of spliced and G/C-rich transcripts.
423 It is, we suggest, no coincidence that the three transcript characteristics associated with NXF1-
424 mediated nuclear export are the three that we propose to be indicative of potentially unwanted
425 transcripts. Similar mechanisms seem to apply to m⁶A : on native transcripts it promotes NXF1-
426 mediated export of few-exon or large-exon transcripts¹⁷ but it is absent from small-exon
427 mRNAs (commonly exported via TREX) owing to the exon-junction complex⁸⁷. Intriguingly,
428 splice factor binding (in other words, SR protein binding) is conserved in intron-bearing and
429 intronless transcripts alike, as revealed by synonymous site conservation in ESEs⁹⁶.
430 Synonymous sites in ESE motifs are preserved in intronless genes as binding of SR proteins
431 to such motifs enables nucleocytoplasmic export via NFX1 (refs^{130,131}). That the HUSH
432 complex identifies A-rich transcripts is of particular note as nuclear export controls also restrict
433 transcripts with A-rich 5' ends¹²⁹ which suggests that multiple quality control filters have similar
434 effects.

435

436 Nuclear export control seems to be a key part of the multi-layered quality control mechanisms
437 and, more generally, retention and triage of unwanted transcripts may be relevant strategies.
438 Cytoplasmic processing bodies contain mostly AU-rich mRNAs, and as a result, AU-rich and
439 GC-rich mRNAs are subject to distinct translational control and decay pathways¹³². AU-rich
440 transcripts are also targeted in the cytoplasm by AU binding proteins that target transcripts to
441 the exosome complex¹³³. Similarly, stress granules inhibit viruses and transposable elements,
442 functioning as a triage⁸⁰ for retention of AU-rich transcripts^{134,135}. Stress granules also have a

443 preference for transcripts that are rare¹³⁵, have poor translatability¹³⁵ and have long m⁶A-
444 tagged exons⁸⁵, which further suggests that their function is part of a quality control
445 mechanism.

446

447 Dinucleotide content (or longer *k*-mers) is also central to RNA degradation in the manner
448 predicted by the unwanted transcript hypothesis. CpG-rich transcripts tend to be directly and
449 preferentially bound by ZAP to promote their decay¹³⁶. Similarly, UpA dinucleotides, which are
450 rare in native transcripts, are targets of ribonuclease L (RNase L)¹³⁷. Likewise, it makes sense
451 that the NEXT complex⁷⁷ targets unprocessed U-rich RNAs¹³⁸, particularly PROMPTs and
452 intron-terminating sequences¹³⁹.

453

454 Although translational selection is relatively unimportant in mammals, it may also be part of a
455 quality control mechanism. Codon optimality modulates ribosome velocity and hence protein
456 folding in yeast¹⁴⁰. Co-translational folding in mammals may also be modulated by ribosome
457 velocity, with phenotypic effects owing to synonymous mutations¹⁴¹. Thus, even if a transcript
458 passes prior quality control filters and makes it to the ribosome, the resulting protein may end
459 up misfolded and hence recycled. Furthermore, whereas it was commonly assumed that
460 codon optimality effects are mediated by ribosomal velocity or misincorporation rates, in yeast,
461 codon optimality is monitored by Dhh1p and Ccr4-Not, leading to degradation if the translation
462 process is too slow¹⁴²⁻¹⁴⁴. This may also be a quality control mechanism in that any transcript
463 being processed too slowly is considered to be 'suspicious' and hence translation is curtailed.
464 Codon composition is thought to affect RNA stability in mammals, possibly by similar
465 routes^{145,146}.

466

467 Can we be confident that these quality control filters are not simply means to regulate genes?
468 One line of evidence comes from genes that can co-evolve with these filters. If such filters
469 exist to counter unwanted transcripts, co-evolving genes would be under selection to counter-
470 adapt to these constraints. Consistent with this, CpG avoidance is commonly observed as an
471 adaptation to human hosts of SARS-CoV-2¹⁴⁷ and other viruses¹⁴⁸. As would be expected if
472 nuclear export is a key quality control filter, nuclear RNA viruses have higher GC content than
473 cytoplasmic RNA viruses¹²⁶.

474

475 A further line of evidence comes from examining the filters that restrict the import of nucleic
476 acids into the cell via endosomes. Such filters mirror those for the rejection of non-native
477 transcripts. Nuclear paraspeckles and cytoplasmic DICER and TLR3 capture long double-
478 stranded RNAs, nuclear export and cytoplasmic transcription both select against AT-rich
479 transcripts, as does cytoplasmic RIG-I (activated by polyU/UC), and ZAP and endosomal

480 TLR9 both filter out CpG-rich RNAs. These parallels we suggest are there because both sets
481 of processes are there to filter out non-native sequences, such as viruses

482

483 **[H1] What can the hypothesis not explain?**

484 There are features with which the unwanted transcript hypothesis does not seem to be
485 compatible. Although we predict that high-GC and low-CpG content function as signs of native
486 genes and, consistent with this, human promoters tend to have high-GC content¹⁴⁹, they can
487 also have more CpG than GpC in the form of CpG islands [G]. Human cells recognise these
488 CpG-rich sequences as native and do not methylate them. Consequently, CpG methylation
489 can be used as a suppressor of transcription (e.g. of transposable elements¹⁵⁰, or of intragene
490 cryptic promoters⁷³) and a high density of unmethylated CpG can be employed as an activator
491 of transcription (for example, as a binding site for the transcriptional activator CFP1 (ref.¹⁵¹)).
492 Unmethylated CpG is not hyper-mutagenic, so its preservation does not contradict our
493 hypothesis. However, it does seem contradictory that mammals have evolved systems to
494 suppress transcripts with high CpG content (for example, via ZAP) but have also evolved CpG
495 islands as indicators of transcript nativity. Why do mammals not methylate all occurrences of
496 CpG? Likewise, other effects of CpG on transcription are not as expected. For example, in
497 some genic contexts, CpG depletion in coding sequences results in reduced protein
498 expression owing to nucleosome positioning-mediated transcriptional effects¹⁵².

499

500 Whereas the unwanted transcript hypothesis can explain why active and successful human
501 transposable elements (Alu elements, short interspersed nuclear elements (SINEs) and
502 variable number of tandem repeat (VNTR) elements) tend to be GC rich, it is not clear how
503 intronless GC-poor L1, the only autonomously transposing transposable element, evades all
504 quality control filters. Indeed, even with the large number of host mechanisms that repress L1
505 at transcriptional, post-transcriptional and post-translational levels, *de novo* L1 insertions still
506 take place within somatic tissues¹⁵³. Part of the explanation is that open reading frame 1
507 (ORF1) of L1 is about 1kb less than the minimum length targeted by the HUSH complex¹⁰⁰.
508 The longer A-rich ORF2 of L1 is usually suppressed, but as transcription is needed for this
509 suppression¹⁰⁰ some leakage is perhaps inevitable.

510

511 Further explanations for the apparent exceptions to the predictions of the unwanted transcript
512 hypothesis might be provided, but whether these are truly 'special case' scenarios or instead
513 reflect a weakness of our model is unclear. One could argue that no quality control filter can
514 be absolute and that leaky expression is always expected to some extent. Indeed, we suggest
515 that the problem of why no filter should be perfect is a theoretically interesting one. Beyond
516 the notion that, as with mutation rates, perfect suppression of errors may well be too costly, a

517 perfect filter against ancestrally intronless transcripts would, for example, have prevented the
518 expression of histones and G-protein-coupled receptors.

519

520 Also, our model does not allow for the possibility that in early human embryogenesis, the
521 various filters against transposable elements and unwanted transcripts are not yet fully in
522 place^{52,154}. Although transposable elements are tightly suppressed in the zygote, they are
523 derepressed during early embryonic development and gametogenesis, when the genome is
524 globally demethylated to allow for epigenetic reprogramming⁵². Recent evidence suggests that
525 active retrotransposition at this stage is kept in check by the removal of highly affected cells
526 rather than by transcript-level processes¹⁵⁴.

527

528 **[H1] What do we not try to explain?**

529 The unwanted transcript hypothesis has little to say about selection on synonymous sites in
530 some contexts. Indeed, we do not suppose that all selection on synonymous sites in mammals
531 is associated with control of unwanted transcripts. For example, human coding sequences are
532 depleted of ultra-stable **G-quadruplexes [G]** with selection against synonymous mutations that
533 would generate them¹⁵⁵. This is likely a result of selection for unhindered transcript
534 processing¹⁵⁵, although an association with the inclusion of transposable elements in stress
535 granules (reviewed in ref.¹⁵⁶) might suggest that it also has a quality control function.

536

537 Selection on synonymous sites for or against micro RNA (miRNA) pairing sites¹⁵⁷⁻¹⁵⁹ may be
538 a further case in point. Notwithstanding the possibility that miRNA was originally an anti-
539 transposable element mechanism (and is used to control L1 (ref.¹⁶⁰)), it now functions also for
540 the control of native gene expression. Importantly, miRNA-mediated regulation is relatively
541 specific in the transcripts affected (a feature of miRNA–RNA pairing) and thus relatively ill-
542 suited as a general-purpose quality control filter. Indeed, it is striking that many of the human
543 quality control filters — HUSH being an exemplar (**Box 2**) — use general ‘fingerprints’ rather
544 than specific motifs for transcript recognition. The one example that uses specific motifs, the
545 preservation of ESEs, is an example of selection to prevent the creation of spurious forms not
546 the recognition of such transcripts.

547

548 **[H1] Testing the hypothesis**

549 How might we test the unwanted transcript hypothesis? Here we start by considering
550 necessary but not sufficient tests. These include possible tests of the assumptions of the
551 hypothesis, notably that such transcripts are costly, and transgene experiments that seek to
552 further define the nucleotide biases that lead to successful expression (and why). We then
553 consider tests of the effects of synonymous mutations. Here we consider both mechanism

554 derived predictions of the synonymous sites that will be under selection and predictions of
555 what happens when synonymous sites are mutated. Finally, we address comparative genomic
556 predictions.

557 **[H2] Testing the cost assumption**

558 Although there is robust evidence that unwanted transcripts pose fitness costs (see above),
559 experimental evolution provides a further means to test this assumption. Naive sequence
560 introduced into yeast produces a large number of presumably functionless transcripts^{27,28}. If
561 our model is right, we would expect that if we select over time for increasing cell fitness coupled
562 with retention of the naive sequence, there will be a reduction in the rate of production of
563 spurious transcripts, particularly those that were originally translated. In addition, given that
564 mutations in the quality control surveillance systems are associated with disease in humans
565 (see above), we expect that comparable mutations in cell lines that disable some quality
566 control filters should be both costly (reduced cell fitness) and lead to compensatory evolution.
567 Predicting what such compensation might involve in mechanistic terms is difficult but a
568 reduction in unwanted protein is a clearer prediction. These predictions come with the caveat
569 that we do not expect that accumulation of the translated products of ptRNAs should lead to
570 increased cell proliferation, although this is seen in humans leading to cancer⁵⁷.

571

572 **[H2] Testing the effects of synonymous mutations**

573 The above tests would not address whether synonymous codon usage affects fitness in a
574 predictable manner. In this respect, transgene codon randomization experiments to modulate
575 ESE usage and SR protein binding, as well as GC content, would be valuable. We predict that
576 codon usage that promotes binding in intronless transcripts of those SR proteins associated
577 with nuclear export will promote higher levels of protein expression. We also predict that
578 modifications of synonymous sites in mammals will mediate their effects in a manner that is
579 predicated on the mutation bias. Changing either A for T (or vice versa), or G for C (or vice
580 versa), thus having no effect on GC3 content, should have relatively little effect on expression
581 levels of a short protein from an intronless transcript (<1.5kb), allowing for alterations in UpA
582 and CpG content. By contrast, changing AT to GC, which increases GC3 content, should
583 increase protein expression, whereas the reverse should typically decrease protein levels. We
584 predict that such changes would be mediated by increased filtering (of AT-rich transcripts) or
585 decreased filtering (of GC-rich transcripts) by quality control mechanisms (rather than being
586 exclusively ribosome mediated), with control of nuclear export being the main such
587 mechanism. Furthermore, for intronless transcripts >1.5kb and so within the scope of HUSH-
588 mediated silencing, the A-to-T or T-to-A effects will be different than for transcripts <1.5kb as
589 the HUSH complex specifically targets A-rich transcripts (**Box 2**).

590

591 Arginine codon usage is potentially especially informative. Our hypothesis is that there should
592 be preference for G or C not just at third codon sites but more generally in native coding
593 sequence. Arginine is specified by AGA, AGG or CGN. We predict that in human cells in which
594 the anti-CpG systems such as ZAP are removed, transgenes with CGC or CGG codons should
595 express more protein than those bearing AGA (all else being equal). AGG-, CGA- and CGT-
596 bearing transgenes should have comparable and intermediate expression levels. Similar logic
597 applies to leucine (specified by UUA, UUC or CUN), with predicted protein-boosting effects of
598 the use of CUC and CUG codons compared with UUA. Serine usage may be less informative
599 as both the two-fold degenerate (AG[UC]) and four-fold degenerate (UCN) blocks are GC/AT-
600 matched at the first two sites, GC variation therefore being exclusive to codon third sites. The
601 evidence also suggests a primacy to 5' codon usage^{16,111}, such that these effects should be
602 more pronounced when modification occurs in the first 20 codons of a transcript. Examination
603 of 5'-to-3' trends of arginine and leucine codon usage in native genes may also be instructive
604 but comes with the caveat that CGN and UUA may attract anti-CpG or anti-UpA quality control
605 measures, respectively.

606

607 ***[H2] Analysis of selection on synonymous mutations***

608 Results from such transgene experiments will be necessary but not sufficient to test the
609 unwanted transcript hypothesis as they do not address the fitness consequences of
610 synonymous mutations in native genes. For predictions about native genes, there are two
611 alternative but related means to investigate the relationship between selection on synonymous
612 mutations and mechanisms of action. First, we can examine the effects of mutations at
613 synonymous sites that are presumed to be under selection¹⁸. Second, we can predict the
614 differential fitness effects of synonymous mutations given a particular mode of action. Fitness
615 effects can be inferred by consideration of between-species conservation or by reference to
616 within-species single nucleotide polymorphism (SNP) frequencies, as selection causes
617 deleterious SNPs to have low frequencies. From the frequency spectrum of SNPs one can, in
618 addition, derive a distribution of fitness effects^{161,20}, although care must be taken to exclude
619 artefacts owing to demography. Both conservation-based¹¹⁵ and SNP frequency-
620 based^{19,20,116,162} analyses report that synonymous sites in ESEs are under purifying selection.
621 Both statistical approaches can be augmented by considering association with genetic
622 diseases, either in well-worked single mutation analyses (**Supplementary Table 1**) or through
623 genome-wide association studies. Many synonymous mutations associated with disease are
624 known to cause deleterious novel splice forms (**Supplementary Table 1**).

625

626 Unfortunately, as there seem to be a large number of quality control filters that affect a range
627 of processes from transcription to translation via nuclear export, there is no comparable,

628 single-mechanism prediction for the unwanted transcript hypothesis as simple as that for
629 translational selection (that codons enriched in highly expressed genes match the most
630 abundant tRNAs). Rather, the unwanted transcript hypothesis proposes that multiple factors
631 affect the activity of synonymous mutations, broadly predicting that mutations that make an
632 RNA look more like foreign RNA should be deleterious, as should those that generate spurious
633 transcripts. Equally, the model predicts which foreign sequences should evade quality control
634 filters, such as GC-rich transposable elements.

635

636 Some effects of synonymous mutations may result from the simple gain or loss of key
637 nucleotides. For example, we predict that 5' mutations in single-exon genes converting GC to
638 AT would be deleterious, whereas the opposite would be advantageous. Generation of CpG
639 and TpA should similarly be deleterious. There should also be a fitness difference on average
640 between GC-to-AT mutations in single-exon genes and in multi-exon genes.

641

642 Effects of synonymous mutations may also manifest through changes to RNA stability and
643 structure. In rodents, synonymous sites that bear fixed substitutions have different effects on
644 *in silico*-predicted RNA stability than those that have no substitution¹⁶³. We can extend this
645 concept to more specifically predict that selection on synonymous mutations will sometimes
646 be mediated by selection to avoid long-stem RNA structures, as these trigger various sensing
647 mechanisms for viruses and transposable elements. It is of note that a synonymous mutation
648 in the *COMT* gene that is associated with reduced protein titre markedly alters predicted RNA
649 folding, giving rise to stable and long-stem structures¹⁶⁴.

650

651 In addition, although CpG and UpA rates are low in human genes, we expect that selection
652 will function to physically 'hide' these dinucleotides in native transcripts, probably mediated by
653 RNA secondary structures. Thus, we expect that there will mutations at synonymous
654 conserved sites that 'unhide' these dinucleotides by altering such structures. This process may
655 account for the observation that coding sequences subject to extensive selection on
656 synonymous mutations are enriched for TpA¹⁶⁵. If this reflects selection on synonymous sites
657 because they modulate RNA structures¹⁶⁵, a lack of UpA exposure is predicted to be disrupted
658 by mutation at the conserved synonymous sites. Differences in native RNA structures near
659 AGA and AGG versus CGN arginine codons would be predicted, the latter being more likely
660 to occur in stem structures than loop structures, assuming that anti-CpG mechanisms target
661 loop-exposed CpG.

662

663 Other aspects of the profile of conserved synonymous sites are predicted by the unwanted
664 transcript hypothesis, in line with the mechanisms that confer foreign versus native status to

665 transcripts. For example, with the same logic as outlined in the ESE-modifying transgene
666 experiment, we also expect (perhaps counter-intuitively) selection in functional retrogenes to
667 favour the retention of ESEs by conserving synonymous sites but only if these permit binding
668 of nuclear export-promoting SR proteins (derived from the parental gene). Similarly, if m⁶A can
669 offer protection to what might otherwise look like a foreign transcript (for example, a single-
670 exon transcript), synonymous mutations that disrupt m⁶A deposition in native single-exon or
671 few-exon transcripts near GC-rich motifs¹⁷ will lead to reduced protein titre owing to disrupted
672 NFX1-mediated nuclear export.

673

674 **[H2] Comparative genomic predictions**

675 In principle, the unwanted transcript hypothesis can make predictions about the need for
676 quality control across taxa. Species with 'lithe' genomes (small intergenic distance, high ratio
677 of coding sequence to genome size, few and small introns, few transposable elements) such
678 as *Saccharomyces cerevisiae* should have fewer problems than humans with unwanted
679 transcripts and so should need fewer quality control mechanisms or to invest less in them.
680 However, the obvious problem is how to assay the number of or investment in quality control
681 filters. It might be notable that the nuclear paraspeckles, which trap Alu elements and double-
682 stranded RNA, are specific to mammals⁶⁰. Also, higher ESE densities in species with larger
683 and more common introns are indicative of greater difficulties in suppressing noisy splicing
684 when N_e is low¹⁶⁶.

685

686 An alternative mode of testing is to examine the nature of the quality control mechanisms. For
687 example, whereas ZAP attacks CpG and RNase L attacks UpA, the two rarest dinucleotides
688 in human coding sequence, we predict that in mammals there should be no common
689 degradation system that targets the most common dinucleotides (GpG, CpC and GpC) prior
690 to translation. It is intriguing that over evolutionary time, the preference of multiple RNase A
691 family members has shifted from UpG (in fish) to UpA (in mammals)¹⁶⁷ consistent with
692 mammalian pressure to target RNA rich in the dinucleotides least common in human CDS.
693 The RNase preferences for A/U need not be the same in a species in which the mutation bias
694 is AT to GC⁹⁷, for which high-GC content could occur owing to neutrality. Unfortunately, we are
695 unaware of any species with both a small N_e and a mutation bias in the opposite direction to
696 that seen in humans, which would be an ideal testing ground for the unwanted transcript
697 hypothesis. However, we suggest that dependency of Pol II-mediated elongation on 5' introns
698 and high 5' GC content will not be universal, that the split functions of the nuclear export
699 channels might be taxonomically restricted to species with a high proportion of junk DNA in
700 their genomes, and that any quality control functionality of m⁶A would be particular to such
701 genomes.

702

703 [H1] Discussion

704 Although rapidly growing microorganisms show selection on codons affecting the translation
705 process (classical translational selection), in this Perspective article we put forward an
706 alternative form of selection based on selection against unwanted transcripts that we suggest
707 is more relevant for error-prone genomes with large amounts of junk DNA such as the human
708 genome. We suggest that it is no coincidence that all of the quality control filters that we have
709 observed function against transcripts rich in A, U, UpA and CpG, the nucleotides that are
710 under-used in native coding sequence in humans.

711

712 An objection to this hypothesis could be that we have mistaken components of gene regulatory
713 systems as components of filters for unwanted transcripts. We suggest that to suppose that
714 any such system must be classified as either a mode of regulation or a filter is a false
715 dichotomy. In evolutionary terms, it would be remarkable if a system that evolved to, for
716 example, suppress a given transposable element, would not also sometimes be repurposed
717 for spurious transcript control and for gene expression regulation (or vice versa). For example,
718 NMD is used to suppress transposable elements, remove erroneous splice forms and regulate
719 amino acid biosynthesis¹⁶⁸. Paraspeckles sequester Alu elements and are used to control
720 circadian gene expression¹⁶⁹. Similarly, gene regulation occurs via modification of the length
721 of the 3' UTR, by removing or adding the Alu element¹⁷⁰. Stress granules both sequester key
722 cell division transcripts for future release and inhibit viruses and transposable elements⁸⁰. In
723 the context of translational selection, that optimal codons are scanned and non-optimal
724 transcripts are degraded is often considered as evidence of a regulatory process controlling
725 protein levels of native transcripts, but it could equally be a means to suppress transcripts
726 without host codon usage patterns^{79,171}. Evidence for the regulatory control of some transcripts
727 by using what might otherwise be considered quality control filters doesn't then necessarily
728 support the above objection. We emphasise the importance of not assuming that everything
729 in a cell has an adaptive regulatory function.

730

731 Not only do we not presume that all selection on synonymous sites is mediated by filters or
732 prevention systems, we also do not presume that all prevention and filtering is mediated by
733 synonymous codon choice. Indeed, we expect there to be a wide range of systems to promote
734 the expression of wanted transcripts and suppress others. For example, mRNA cap structures
735 function as effective guides to wanted transcripts. For this reason, some viruses (such as
736 influenza A virus) have evolved cap-snatching from host RNAs¹⁷². The polyA tail of mRNAs is
737 similarly used for quality control¹⁷³. What we do suggest is that to understand codon usage
738 trends, it can be helpful to think in terms of keeping unwanted transcripts in check, rather than

739 focusing only on translational selection. Indeed, whereas codon optimization of synthetic
740 RNAs premised on the concept of tRNA-codon availability, is routinely carried out, the
741 assumptions of the method are not upheld in humans¹⁷⁴. It is instead likely that any effects of
742 codon optimization are mediated by the described quality control filters¹⁶.

743

744 Aside from the complex role of m⁶A in quality control, we highlight two major unknowns
745 concerning our model. First, is the extent to which the model has validity across taxa. We have
746 highlighted relevance of the hypothesis to humans as we are interested in human diseases
747 and possible cures, and because humans are an exemplar of a species that we expect to have
748 problems with unwanted transcripts and for which translational selection has mostly not been
749 observed. This is because humans have low N_e and hence inefficient selection, which can
750 explain our large genome with a low ratio of coding sequence to genome size (**Box 1**), a
751 genome in which we expect poorer control of processes such as splicing and transcription.
752 Despite the fact that many unwanted transcripts, such as PROMPTs, may be universal,
753 generally, we expect that as N_e increases so the weak selection associated with translational
754 selection should become more important in determining codon usage and the unwanted
755 transcript problem should diminish, as selection over the control of gene expression and
756 splicing should increase. Evidence that rare spurious transcripts are less common when N_e is
757 large⁴⁵ is consistent with this, as is, for example, evidence that non-optimal stop codon usage
758 is predicted by small N_e ¹⁷⁵ and that purifying selection acts on non-coding RNA in *Drosophila*
759 but not so obviously in humans³³.

760

761 One means to establish the phylogenetic extent of selection for synonymous sites to be
762 recognized as native may be through analysis of coding sequence 5' ends. In many taxa,
763 coding sequence 5' ends are under selection to generate unstable RNAs so as to mediate
764 more efficient translation initiation¹²⁸, thus explaining high 5' AT content¹⁷⁶, including at third
765 codon sites. In mammals, this selection for low RNA stability is not evident¹²⁸. We suggest that
766 this is because the counteracting selection for high GC content so as to appear as a native
767 transcript is dominant. If this logic is correct, cross-species analysis of 5' GC3 content may
768 help to determine the taxonomic relevance of selection based on wanted versus unwanted
769 transcripts.

770

771 The second major unknown is how this model relates to the mechanisms by which
772 synonymous mutations might induce genetic diseases. Selection on most synonymous
773 mutations is expected to be weak; indeed, whereas altering the GC3 content of a human
774 transgene has major expression effects^{15,16}, the impact of changing just one synonymous site
775 must, on average, be more subtle. By contrast, individual splice-disrupting mutations have

776 the potential to have major fitness consequences. Estimates suggest that ~10–30% of all
777 disease mutations mediate their effects through disruption of splice motifs such as ESEs (see
778 refs¹⁷⁷⁻¹⁷⁹ and references therein). Alternative mechanisms whereby a single synonymous
779 mutation has a large enough effect to cause disease are unclear; although examples can be
780 provided, their commonality is not known (**Supplementary Table 1; Fig. 4**). Although splicing
781 seems to be the dominant mode of disease-associated action, ascertainment bias is likely to
782 be profound.

783

784 The unwanted transcript hypothesis has implications for gene therapy, recombinant protein
785 production and the problem of how to make an intronless transcript. An obvious solution given
786 the above is to render at least the 5' end codons GC3 rich¹⁸⁰. CpG avoidance (or enrichment
787 to make attenuated live vaccines¹⁸¹) is a common strategy in multiple contexts^{114,182}. It is
788 notable that similar themes apply to RNA vaccines. Indeed, one breakthrough in their
789 development was recognition that RNA base modifications suppress TLR activity¹⁸³. The use
790 of pseudouridine in mRNA vaccines seems to enable RNA to be translated but simultaneously
791 evades the U-sensing filters¹⁸³. Parallels can be made to m⁶A deposition enabling the nuclear
792 export¹⁷ of what might otherwise appear to be foreign transcripts (those with few long exons
793 or single exons). It is similarly notable that reduced CpG content and high-GC content both
794 enable nuclear transgene activity¹⁶ and the efficacy of externally delivered mRNA¹⁸⁴. That
795 TLRs and the control of nuclear transcripts have evolved similar rules to detect that which is
796 foreign explains why RNA vaccines and transgenes benefit from the same sort of modifications
797 at synonymous sites and why, more generally, their design principles are similar.

798

799 **References**

- 800 1 King, J. L. & Jukes, T. H. Non-Darwinian evolution. *Science* **164**, 788-798, (1969).
- 801 2 Sharp, P. M., Averof, M., Lloyd, A. T., Matassi, G. & Peden, J. F. DNA sequence
802 evolution: the sounds of silence. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **349**, 241-
803 247, (1995).
- 804 3 Ikemura, T. Correlation between the abundance of *Escherichia coli* transfer
805 RNAs and the occurrence of the respective codons in its protein genes: a
806 proposal for a synonymous codon choice that is optimal for the E. coli
807 translational system. *J. Mol. Biol.* **151**, 389-409, (1981).
- 808 4 Ikemura, T. Codon usage and tRNA content in unicellular and multicellular
809 organisms. *Mol. Biol. Evol.* **2**, 13-34, (1985).
- 810 5 Sharp, P. M. & Li, W.-H. The codon adaptation index—a measure of directional
811 synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.*
812 **15**, 1281-1295, (1987).

813 6 Qian, W., Yang, J. R., Pearson, N. M., Maclean, C. & Zhang, J. Balanced codon
814 usage optimizes eukaryotic translational efficiency. *PLoS Genet.* **8**, e1002603,
815 (2012).

816 7 Akashi, H. Synonymous codon usage in *Drosophila melanogaster*: Natural
817 selection and translational accuracy. *Genetics* **136**, 927-935, (1994).

818 8 Stoletzki, N. & Eyre-Walker, A. Synonymous codon usage in *Escherichia coli*:
819 selection for translational accuracy. *Mol. Biol. Evol.* **24**, 374-381, (2007).

820 9 Sharp, P. M., Emery, L. R. & Zeng, K. Forces that influence the evolution of codon
821 bias. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **365**, 1203-1212, (2010).

822 10 dos Reis, M. & Wernisch, L. Estimating translational selection in eukaryotic
823 genomes. *Mol. Biol. Evol.* **26**, 451-461, (2009).

824 11 Lynch, M. & Conery, J. S. The origins of genome complexity. *Science* **302**, 1401-
825 1404, (2003).

826 12 Duret, L. Evolution of synonymous codon usage in metazoans. *Curr. Opin.*
827 *Genet. Dev.* **12**, 640-649, (2002).

828 13 Hunt, R. C., Simhadri, V. L., Iandoli, M., Sauna, Z. E. & Kimchi-Sarfaty, C. Exposing
829 synonymous mutations. *Trends Genet.* **30**, 308-321, (2014).

830 14 Bali, V. & Bebok, Z. Decoding mechanisms by which silent codon changes
831 influence protein biogenesis and function. *Int. J. Biochem. Cell Biol.* **64**, 58-74,
832 (2015).

833 15 Kudla, G., Lipinski, L., Caffin, F., Helwak, A. & Zylicz, M. High guanine and
834 cytosine content increases mRNA levels in mammalian cells. *PLoS Biol.* **4**, e180,
835 (2006).

836 16 Mordstein, C. *et al.* Codon usage and splicing jointly influence mRNA
837 localization. *Cell Syst.* **10**, 351-362 e358, (2020).

838 17 Zuckerman, B., Ron, M., Mikl, M., Segal, E. & Ulitsky, I. Gene architecture and
839 sequence composition underpin selective dependency of nuclear export of long
840 RNAs on NXF1 and the TREX Complex. *Mol. Cell* **79**, 251-267 e256, (2020).

841 18 Lin, M. F. *et al.* Locating protein-coding sequences under selection for
842 additional, overlapping functions in 29 mammalian genomes. *Genome Res.* **21**,
843 1916-1928, (2011).

844 19 Caceres, E. F. & Hurst, L. D. The evolution, impact and properties of exonic splice
845 enhancers. *Genome Biol* **14**, R143, (2013).

846 20 Savisaar, R. & Hurst, L. D. Exonic splice regulation imposes strong selection at
847 synonymous sites. *Genome Res.* **28**, 1442-1454, (2018).

848 21 Keightley, P. D. & Halligan, D. L. Inference of site frequency spectra from high-
849 throughput sequence data: quantification of selection on nonsynonymous and
850 synonymous sites in humans. *Genetics* **188**, 931-940, (2011).

851 22 Eory, L., Halligan, D. L. & Keightley, P. D. Distributions of selectively constrained
852 sites and deleterious mutation rates in the hominid and murid genomes. *Mol.*
853 *Biol. Evol.* **27**, 177-192, (2010).

854 23 Wen, P., Xiao, P. & Xia, J. dbDSM: a manually curated database for deleterious
855 synonymous mutations. *Bioinformatics* **32**, 1914-1916, (2016).

856 24 Statello, L., Guo, C.-J., Chen, L.-L. & Huarte, M. Gene regulation by long non-
857 coding RNAs and its biological functions. *Nat. Rev. Mol. Cell Biol.* **22**, 96-118,
858 (2021).

859 25 Hurst, L. D. Evolutionary genomics and the reach of selection. *J. Biol.* **8**, 12,
860 (2009).

861 26 Andrews, G. *et al.* Mammalian evolution of human cis-regulatory elements and
862 transcription factor binding sites. *Science* **380**, eabn7930, (2023).

863 27 Luthra, I. *et al.* Biochemical activity is the default DNA state in eukaryotes.
864 Preprint at bioRxiv <https://www.biorxiv.org/node/2906314.abstract> (2022).

865 28 Camellato, B., Brosh, R., Maurano, M. T. & Boeke, J. D. Genomic analysis of a
866 synthetic reversed sequence reveals default chromatin states in yeast and
867 mammalian cells. Preprint at bioRxiv
868 <https://www.biorxiv.org/content/10.1101/2023.06.20.545713v2> (2022).

869 29 Xu, H., Li, C., Xu, C. & Zhang, J. Chance promoter activities illuminate the origins
870 of eukaryotic intergenic transcriptions. *Nat Commun* **14**, 1826, (2023).

871 30 Preker, P. *et al.* PROMoter uPstream Transcripts share characteristics with
872 mRNAs and are produced upstream of all three major types of mammalian
873 promoters. *Nucleic Acids Res.* **39**, 7179-7193, (2011).

874 31 Schuler, A., Ghanbarian, A. T. & Hurst, L. D. Purifying selection on splice-related
875 motifs, not expression level nor RNA folding, explains nearly all constraint on
876 human lincRNAs. *Mol. Biol. Evol.* **31**, 3164-3183, (2014).

877 32 Managadze, D., Rogozin, I. B., Chernikova, D., Shabalina, S. A. & Koonin, E. V.
878 Negative correlation between expression level and evolutionary rate of long
879 intergenic noncoding RNAs. *Genome Biol Evol* **3**, 1390-1404, (2011).

880 33 Haerty, W. & Ponting, C. P. Mutations within lincRNAs are effectively selected
881 against in fruitfly but not in human. *Genome Biol* **14**, R49, (2013).

882 34 Johnsson, P., Lipovich, L., Grandér, D. & Morris, K. V. Evolutionary conservation
883 of long non-coding RNAs; sequence, structure, function. *Biochim Biophys Acta*
884 **1840**, 1063-1071, (2014).

885 35 Ponting, C. P. & Haerty, W. Genome-Wide Analysis of Human Long Noncoding
886 RNAs: A Provocative Review. *Annu Rev Genomics Hum Genet* **23**, 153-172,
887 (2022).

888 36 Wyers, F. *et al.* Cryptic Pol II transcripts are degraded by a nuclear quality control
889 pathway involving a new poly(A) polymerase. *Cell* **121**, 725-737, (2005).

890 37 Liu, S. J. *et al.* CRISPRi-based genome-scale identification of functional long
891 noncoding RNA loci in human cells. *Science* **355**, aah7111, (2017).

892 38 Schlackow, M. *et al.* Distinctive patterns of transcription and RNA processing for
893 Human lincRNAs. *Mol. Cell* **65**, 25-38, (2017).

894 39 Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable
895 elements: from conflicts to benefits. *Nat. Rev. Genet.* **18**, 71-86, (2017).

896 40 Wang, J. *et al.* Primate-specific endogenous retrovirus-driven transcription
897 defines naive-like stem cells. *Nature* **516**, 405-409, (2014).

898 41 Raskó, T. *et al.* A novel gene controls a new structure: PiggyBac Transposable
899 Element-Derived 1, unique to mammals, controls mammal-specific neuronal
900 paraspeckles. *Molecular Biology and Evolution* **39**, (2022).

901 42 Consortium, E. P. An integrated encyclopedia of DNA elements in the human
902 genome. *Nature* **489**, 57-74, (2012).

903 43 Carlevaro-Fita, J. *et al.* Ancient exapted transposable elements promote nuclear
904 enrichment of human long noncoding RNAs. *Genome Res.* **29**, 208-222, (2019).

905 44 Pickrell, J. K., Pai, A. A., Gilad, Y. & Pritchard, J. K. Noisy splicing drives mRNA
906 isoform diversity in human cells. *PLoS Genet.* **6**, e1001236, (2010).

907 45 Bénétière, F., Necsulea, A. & Duret, L. Random genetic drift sets an upper limit
908 on mRNA splicing accuracy in metazoans. Preprint at bioRxiv
909 <https://www.biorxiv.org/content/10.1101/2022.12.09.519597v5> (2023).

910 46 Irimia, M. *et al.* Complex selection on 5' splice sites in intron-rich organisms.
911 *Genome Res.* **19**, 2021-2027, (2009).

912 47 Savisaar, R. & Hurst, L. D. Estimating the prevalence of functional exonic splice
913 regulatory information. *Hum. Genet.* **136**, 1059-1078, (2017).

914 48 Wagner, A. Energy constraints on the evolution of gene expression. *Mol. Biol.*
915 *Evol.* **22**, 1365-1374, (2005).

916 49 Kudla, G., Murray, A. W., Tollervey, D. & Plotkin, J. B. Coding-sequence
917 determinants of gene expression in *Escherichia coli*. *Science* **324**, 255-258,
918 (2009).

919 50 Cambray, G., Guimaraes, J. C. & Arkin, A. P. Evaluation of 244,000 synthetic
920 sequences reveals design principles to optimize translation in *Escherichia coli*.
921 *Nat. Biotechnol.* **36**, 1005-1015, (2018).

922 51 Mittal, P., Brindle, J., Stephen, J., Plotkin, J. B. & Kudla, G. Codon usage influences
923 fitness through RNA toxicity. *Proc Natl Acad Sci U S A* **115**, 8639-8644, (2018).

924 52 Bourque, G. *et al.* Ten things you should know about transposable elements.
925 *Genome Biol* **19**, 199, (2018).

926 53 Lionetti, M. *et al.* A compendium of DIS3 mutations and associated
927 transcriptional signatures in plasma cell dyscrasias. *Oncotarget* **6**, (2015).

928 54 Fasken, M. B. *et al.* The RNA exosome and human disease. *Methods Mol Biol*
929 **2062**, 3-33, (2020).

930 55 Morton, D. J. *et al.* The RNA exosome and RNA exosome-linked disease. *RNA*
931 **24**, 127-142, (2018).

932 56 Giunta, M. *et al.* Altered RNA metabolism due to a homozygous RBM7 mutation
933 in a patient with spinal motor neuropathy. *Hum. Mol. Genet.* **25**, 2985-2996,
934 (2016).

935 57 Insko, M. L. *et al.* Oncogenic CDK13 mutations impede nuclear RNA surveillance.
936 *Science* **380**, eabn7625, (2023).

937 58 Luo, S. *et al.* The evolutionary arms race between transposable elements and
938 piRNAs in *Drosophila melanogaster*. *BMC Evol. Biol.* **20**, 14, (2020).

939 59 Bertozzi, T. M., Elmer, J. L., Macfarlan, T. S. & Ferguson-Smith, A. C. KRAB zinc
940 finger protein diversification drives mammalian interindividual methylation
941 variability. *Proc Natl Acad Sci U S A* **117**, 31290-31300, (2020).

942 60 Fox, A. H. & Lamond, A. I. Paraspeckles. *Cold Spring Harb. Perspect. Biol.* **2**,
943 a000687, (2010).

944 61 Kaneko, H. *et al.* DICER1 deficit induces Alu RNA toxicity in age-related macular
945 degeneration. *Nature* **471**, 325-330, (2011).

946 62 Muotri, A. R. *et al.* L1 retrotransposition in neurons is modulated by MeCP2.
947 *Nature* **468**, 443-446, (2010).

948 63 Tsvion-Visbord, H. *et al.* Increased RNA editing in maternal immune activation
949 model of neurodevelopmental disease. *Nat Commun* **11**, 5236, (2020).

950 64 Ansell, B. R. E. *et al.* A survey of RNA editing at single-cell resolution links
951 interneurons to schizophrenia and autism. *RNA* **27**, 1482-1496, (2021).

952 65 Li, P. *et al.* Aicardi-Goutieres syndrome protein TREX1 suppresses L1 and
953 maintains genome integrity through exonuclease-independent ORF1p
954 depletion. *Nucleic Acids Res.* **45**, 4619-4631, (2017).

955 66 Stearrett, N. *et al.* Expression of Human endogenous retroviruses in systemic
956 lupus erythematosus: multiomic integration with gene expression. *Front.*
957 *Immunol.* **12**, 661437, (2021).

958 67 Dembny, P. *et al.* Human endogenous retrovirus HERV-K(HML-2) RNA causes
959 neurodegeneration through Toll-like receptors. *JCI Insight* **5**, (2020).

960 68 Ramirez, P. *et al.* Pathogenic tau accelerates aging-associated activation of
961 transposable elements in the mouse central nervous system. *Prog. Neurobiol.*
962 **208**, 102181, (2022).

963 69 Grundy, E. E., Diab, N. & Chiappinelli, K. B. Transposable element regulation and
964 expression in cancer. *FEBS J.* **289**, 1160-1179, (2022).

965 70 Van Meter, M. *et al.* SIRT6 represses LINE1 retrotransposons by ribosylating
966 KAP1 but this repression fails with stress and age. *Nat Commun* **5**, 5011, (2014).

967 71 Hastings, M. L. & Krainer, A. R. Pre-mRNA splicing in the new millennium. *Curr.*
968 *Opin. Cell Biol.* **13**, 302-309, (2001).

969 72 Liu, H. X., Cartegni, L., Zhang, M. Q. & Krainer, A. R. A mechanism for exon
970 skipping caused by nonsense or missense mutations in BRCA1 and other genes.
971 *Nat. Genet.* **27**, 55-58, (2001).

972 73 Neri, F. *et al.* Intragenic DNA methylation prevents spurious transcription
973 initiation. *Nature* **543**, 72-77, (2017).

974 74 Ilinskaya, O. N. & Mahmud, R. S. Ribonucleases as antiviral agents. *Mol. Biol.* **48**,
975 615-623, (2014).

976 75 Meola, N. *et al.* Identification of a nuclear exosome decay pathway for processed
977 transcripts. *Mol. Cell* **64**, 520-533, (2016).

978 76 Ogami, K. *et al.* An Mtr4/ZFC3H1 complex facilitates turnover of unstable
979 nuclear RNAs to prevent their cytoplasmic transport and global translational
980 repression. *Genes Dev* **31**, 1257-1271, (2017).

981 77 Lubas, M. *et al.* Interaction profiling identifies the human nuclear exosome
982 targeting complex. *Mol. Cell* **43**, 624-637, (2011).

983 78 Chen, L. L., DeCerbo, J. N. & Carmichael, G. G. Alu element-mediated gene
984 silencing. *EMBO J.* **27**, 1694-1705, (2008).

985 79 Monaghan, L., Longman, D. & Cáceres, J. F. Translation-coupled mRNA quality
986 control mechanisms. *EMBO J.* **42**, e114378, (2023).

987 80 Anderson, P. & Kedersha, N. Stress granules: the Tao of RNA triage. *Trends*
988 *Biochem. Sci.* **33**, 141-150, (2008).

989 81 Ding, S. W. & Voinnet, O. Antiviral immunity directed by small RNAs. *Cell* **130**,
990 413-426, (2007).

991 82 Gao, G. X., Guo, X. M. & Goff, S. P. Inhibition of retroviral RNA production by
992 ZAP, a CCCH-type zinc finger protein. *Science* **297**, 1703-1706, (2002).

993 83 Kesner, J. S. *et al.* Noncoding translation mitigation. *Nature* **617**, 395-402,
994 (2023).

995 84 Liu, J. *et al.* The RNA m6A reader YTHDC1 silences retrotransposons and guards
996 ES cell identity. *Nature* **591**, 322-326, (2021).

997 85 Ries, R. J., Pickering, B. F., Poh, H. X., Namkoong, S. & Jaffrey, S. R. m6A governs
998 length-dependent enrichment of mRNAs in stress granules. *Nat. Struct. Mol.*
999 *Biol.*, (2023).

1000 86 Lee, E. S. *et al.* N-6-methyladenosine (m6A) promotes the nuclear retention of
1001 mRNAs with intact 5' splice site motifs. Preprint at bioRxiv
1002 <https://www.biorxiv.org/content/10.1101/2023.06.20.545713v2> (2023).

1003 87 He, P. C. *et al.* Exon architecture controls mRNA m(6)A suppression and gene
1004 expression. *Science* **379**, 677-682, (2023).

1005 88 Delaunay, S., Helm, M. & Frye, M. RNA modifications in physiology and disease:
1006 towards clinical applications. *Nat. Rev. Genet.*, doi.org/10.1038/s41576-41023-
1007 00645-41572, (2023).

1008 89 Sun, T. *et al.* Crosstalk between RNA m6A and DNA methylation regulates
1009 transposable element chromatin activation and cell fate in human pluripotent
1010 stem cells. *Nat. Genet.* **55**, 1324-1335, (2023).

1011 90 Janeway, C. A., Jr. Approaching the asymptote? Evolution and revolution in
1012 immunology. *Cold Spring Harb. Symp. Quant. Biol.* **54 Pt 1**, 1-13, (1989).

1013 91 Logsdon, J. M. The recent origins of spliceosomal introns revisited. *Curr. Opin.*
1014 *Genet. Dev.* **8**, 637-648, (1998).

1015 92 Sakharkar, M. K., Chow, V. T. & Kanguane, P. Distributions of exons and introns
1016 in the human genome. *In Silico Biol.* **4**, 387-393, (2004).

1017 93 Zhang, J., Sun, X. L., Qian, Y. M., LaDuca, J. P. & Maquat, L. E. At least one intron
1018 is required for the nonsense-mediated decay of triosephosphate isomerase
1019 mRNA: a possible link between nuclear splicing and cytoplasmic translation.
1020 *Mol Cell Biol* **18**, 5272-5283, (1998).

1021 94 Le Hir, H., Nott, A. & Moore, M. J. How introns influence and enhance eukaryotic
1022 gene expression. *Trends Biochem. Sci.* **28**, 215-220, (2003).

1023 95 Brocke, K. S., Neu-Yilik, G., Gehring, N. H., Hentze, M. W. & Kulozik, A. E. The
1024 human intronless melanocortin 4-receptor gene is NMD insensitive. *Hum. Mol.*
1025 *Genet.* **11**, 331-335, (2002).

1026 96 Savisaar, R. & Hurst, L. D. Purifying selection on Exonic Splice Enhancers in
1027 intronless genes. *Mol. Biol. Evol.* **33**, 1396-1418, (2016).

1028 97 Long, H. *et al.* Evolutionary determinants of genome-wide nucleotide
1029 composition. *Nat Ecol Evol* **2**, 237-240, (2018).

1030 98 Ho, A. T. & Hurst, L. D. Unusual mammalian usage of TGA stop codons reveals
1031 that sequence conservation need not imply purifying selection. *PLoS Biol.* **20**,
1032 e3001588, (2022).

1033 99 Charneski, C. A., Honti, F., Bryant, J. M., Hurst, L. D. & Feil, E. J. Atypical at skew
1034 in Firmicute genomes results from selection and not from mutation. *PLoS Genet.*
1035 **7**, e1002283, (2011).

1036 100 Seczynska, M., Bloor, S., Cuesta, S. M. & Lehner, P. J. Genome surveillance by
1037 HUSH-mediated silencing of intronless mobile elements. *Nature* **601**, 440-445,
1038 (2022).

1039 101 Duret, L. & Galtier, N. Biased gene conversion and the evolution of mammalian
1040 genomic landscapes. *Annu Rev Genomics Hum Genet* **10**, 285-311, (2009).

1041 102 Liu, H. *et al.* Tetrad analysis in plants and fungi finds large differences in gene
1042 conversion rates but no GC bias. *Nat Ecol Evol* **2**, 164-173, (2018).

1043 103 Galtier, N. Gene conversion drives GC content evolution in mammalian histones.
1044 *Trends Genet.* **19**, 65-68, (2003).

1045 104 D'Onofrio, G., Mouchiroud, D., Aissani, B., Gautier, C. & Bernardi, G. Correlations
1046 between the compositional properties of human genes, codon usage, and
1047 amino-acid-composition of proteins. *J. Mol. Evol.* **32**, 504-510, (1991).

1048 105 Duret, L. & Hurst, L. D. The elevated GC content at exonic third sites is not
1049 evidence against neutralist models of isochores evolution. *Mol. Biol. Evol.* **18**,
1050 757-762, (2001).

1051 106 Duret, L. & Galtier, N. The covariation between TpA deficiency, CpG deficiency,
1052 and G+C content of human isochores is due to a mathematical artifact. *Mol.*
1053 *Biol. Evol.* **17**, 1620-1625, (2000).

1054 107 Morales, A. C. *et al.* Causes and consequences of purifying selection on SARS-
1055 CoV-2. *Genome Biol Evol* **13**, 17, (2021).

1056 108 Zhou, Z. *et al.* Codon usage is an important determinant of gene expression
1057 levels largely through its effects on transcription. *Proc Natl Acad Sci U S A* **113**,
1058 E6117-E6125, (2016).

1059 109 Newman, Z. R., Young, J. M., Ingolia, N. T. & Barton, G. M. Differences in codon
1060 bias and GC content contribute to the balanced expression of TLR7 and TLR9.
1061 *Proc Natl Acad Sci U S A* **113**, E1362-1371, (2016).

1062 110 Zhao, F. *et al.* Genome-wide role of codon usage on transcription and
1063 identification of potential regulators. *Proc Natl Acad Sci U S A* **118**, (2021).

1064 111 Vlaming, H., Mimoso, C. A., Field, A. R., Martin, B. J. E. & Adelman, K. Screening
1065 thousands of transcribed coding and non-coding regions reveals sequence

1066 determinants of RNA polymerase II elongation potential. *Nat. Struct. Mol. Biol.*
1067 **29**, 613-620, (2022).

1068 112 Pantier, R. *et al.* SALL4 controls cell fate in response to DNA base composition.
1069 *Mol. Cell* **81**, 845-858, (2021).

1070 113 Hisano, M., Ohta, H., Nishimune, Y. & Nozaki, M. Methylation of CpG
1071 dinucleotides in the open reading frame of a testicular germ cell-specific
1072 intronless gene, Tact1/Actl7b, represses its expression in somatic cells. *Nucleic*
1073 *Acids Res.* **31**, 4797-4804, (2003).

1074 114 Hodges, B. L., Taylor, K. M., Joseph, M. F., Bourgeois, S. A. & Scheule, R. K. Long-
1075 term transgene expression from plasmid DNA gene therapy vectors is
1076 negatively affected by CpG dinucleotides. *Mol. Ther.* **10**, 269-278, (2004).

1077 115 Parmley, J. L., Chamary, J. V. & Hurst, L. D. Evidence for purifying selection
1078 against synonymous mutations in mammalian exonic splicing enhancers. *Mol.*
1079 *Biol. Evol.* **23**, 301-309, (2006).

1080 116 Carlini, D. B. & Genut, J. E. Synonymous SNPs provide evidence for selective
1081 constraint on human exonic splicing enhancers. *J. Mol. Evol.* **62**, 89-98, (2006).

1082 117 Sauna, Z. E. & Kimchi-Sarfaty, C. Understanding the contribution of synonymous
1083 mutations to human disease. *Nat. Rev. Genet.* **12**, 683-691, (2011).

1084 118 Savisaar, R. & Hurst, L. D. Both maintenance and avoidance of RNA-binding
1085 protein interactions constrain coding sequence evolution. *Mol. Biol. Evol.* **34**,
1086 1110-1126, (2017).

1087 119 Parmley, J. L. & Hurst, L. D. Exonic splicing regulatory elements skew
1088 synonymous codon usage near intron-exon boundaries in mammals. *Mol. Biol.*
1089 *Evol.* **24**, 1600-1603, (2007).

1090 120 Willie, E. & Majewski, J. Evidence for codon bias selection at the pre-mRNA level
1091 in eukaryotes. *Trends Genet.* **20**, 534-538, (2004).

1092 121 Warnecke, T. & Hurst, L. D. Evidence for a trade-off between translational
1093 efficiency and splicing regulation in determining synonymous codon usage in
1094 *Drosophila melanogaster*. *Mol. Biol. Evol.* **24**, 2755-2762, (2007).

1095 122 Eskesen, S. T., Eskesen, F. N. & Ruvinsky, A. Natural selection affects frequencies
1096 of AG and GT dinucleotides at the 5' and 3' ends of exons. *Genetics* **167**, 543-
1097 550, (2004).

1098 123 Livingstone, M. *et al.* Investigating DNA-, RNA-, and protein-based features as
1099 a means to discriminate pathogenic synonymous variants. *Hum. Mutat.* **38**,
1100 1336-1347, (2017).

1101 124 Eyre-Walker, A. & Hurst, L. D. The evolution of isochores. *Nat. Rev. Genet.* **2**,
1102 549-555, (2001).

1103 125 Potrzebowski, L. *et al.* Chromosomal gene movements reflect the recent origin
1104 and biology of therian sex chromosomes. *PLoS Biol.* **6**, e80, (2008).

1105 126 Mordstein, C. *et al.* Transcription, mRNA Export, and Immune Evasion Shape the
1106 Codon Usage of Viruses. *Genome Biol Evol* **13**, (2021).

1107 127 Goodman, D. B., Church, G. M. & Kosuri, S. Causes and effects of N-terminal
1108 codon bias in bacterial genes. *Science* **342**, 475-479, (2013).

1109 128 Gu, W., Zhou, T. & Wilke, C. O. A universal trend of reduced mRNA stability near
 1110 the translation-initiation site in prokaryotes and eukaryotes. *PLoS Comput Biol*
 1111 **6**, e1000664, (2010).
 1112 129 Palazzo, A. F. *et al.* The signal sequence coding region promotes nuclear export
 1113 of mRNA. *PLoS Biol.* **5**, e322, (2007).
 1114 130 Huang, Y., Gattoni, R., Stevenin, J. & Steitz, J. A. SR splicing factors serve as
 1115 adapter proteins for TAP-dependent mRNA export. *Mol. Cell* **11**, 837-843,
 1116 (2003).
 1117 131 Huang, Y. & Steitz, J. A. Splicing factors SRp20 and 9G8 promote the
 1118 nucleocytoplasmic export of mRNA. *Mol. Cell* **7**, 899-905, (2001).
 1119 132 Courel, M. *et al.* GC content shapes mRNA storage and decay in human cells.
 1120 *Elife* **8**, e49708, (2019).
 1121 133 Chen, C.-Y. *et al.* AU binding proteins recruit the exosome to degrade ARE-
 1122 containing mRNAs. *Cell* **107**, 451-464, (2001).
 1123 134 Namkoong, S., Ho, A., Woo, Y. M., Kwak, H. & Lee, J. H. Systematic
 1124 characterization of stress-induced RNA granulation. *Mol. Cell* **70**, 175-187 e178,
 1125 (2018).
 1126 135 Khong, A. *et al.* The stress granule transcriptome reveals principles of mRNA
 1127 accumulation in stress granules. *Mol. Cell* **68**, 808-820.e805, (2017).
 1128 136 Takata, M. A. *et al.* CG dinucleotide suppression enables antiviral defence
 1129 targeting non-self RNA. *Nature* **550**, 124-127, (2017).
 1130 137 Duan, J. & Antezana, M. A. Mammalian mutation pressure, synonymous codon
 1131 choice, and mRNA degradation. *J. Mol. Evol.* **57**, 694-701, (2003).
 1132 138 Hrossova, D. *et al.* RBM7 subunit of the NEXT complex binds U-rich sequences
 1133 and targets 3'-end extended forms of snRNAs. *Nucleic Acids Res.* **43**, 4236-4248,
 1134 (2015).
 1135 139 Lubas, M. *et al.* The human nuclear exosome targeting complex is loaded onto
 1136 newly synthesized RNA to direct early ribonucleolysis. *Cell Rep.* **10**, 178-192,
 1137 (2015).
 1138 140 Pechmann, S. & Frydman, J. Evolutionary conservation of codon optimality
 1139 reveals hidden signatures of cotranslational folding. *Nat. Struct. Mol. Biol.* **20**,
 1140 237-243, (2013).
 1141 141 Kimchi-Sarfaty, C. *et al.* A "silent" polymorphism in the MDR1 gene changes
 1142 substrate specificity. *Science* **315**, 525-528, (2007).
 1143 142 Radhakrishnan, A. *et al.* The DEAD-Box protein Dhh1p couples mRNA decay and
 1144 translation by monitoring codon optimality. *Cell* **167**, 122-132 e129, (2016).
 1145 143 Radhakrishnan, A. & Green, R. Connections underlying translation and mRNA
 1146 stability. *J. Mol. Biol.* **428**, 3558-3564, (2016).
 1147 144 Buschauer, R. *et al.* The Ccr4-Not complex monitors the translating ribosome
 1148 for codon optimality. *Science* **368**, (2020).
 1149 145 Medina-Munoz, S. G. *et al.* Crosstalk between codon optimality and cis-
 1150 regulatory elements dictates mRNA stability. *Genome Biol* **22**, 14, (2021).

- 1151 146 Shu, H. *et al.* FMRP links optimal codons to mRNA stability in neurons. *Proc Natl Acad Sci U S A* **117**, 30400-30411, (2020).
- 1152
- 1153 147 Kumar, A. *et al.* The slowing rate of CpG depletion in SARS-CoV-2 genomes is consistent with adaptations to the human host. *Mol. Biol. Evol.* **39**, (2022).
- 1154
- 1155 148 Ficarelli, M. *et al.* CpG dinucleotides inhibit HIV-1 replication through Zinc Finger Antiviral Protein (ZAP)-dependent and -independent mechanisms. *J. Virol.* **94**, e01337-01319, (2020).
- 1156
- 1157
- 1158 149 Hurst, L. D. *et al.* A simple metric of promoter architecture robustly predicts expression breadth of human genes suggesting that most transcription factors are positive regulators. *Genome Biol* **15**, 413, (2014).
- 1159
- 1160
- 1161 150 Bestor, T. H. DNA methylation: evolution of a bacterial immune function into a regulator of gene expression and genome structure in higher eukaryotes. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **326**, 179-187, (1990).
- 1162
- 1163
- 1164 151 Voo, K. S., Carlone, D. L., Jacobsen, B. M., Flodin, A. & Skalnik, D. G. Cloning of a mammalian transcriptional activator that binds unmethylated CpG motifs and shares a CXXC domain with DNA methyltransferase, human trithorax, and methyl-CpG binding domain protein 1. *Mol Cell Biol* **20**, 2108-2121, (2000).
- 1165
- 1166
- 1167
- 1168 152 Bauer, A. P. *et al.* The impact of intragenic CpG content on gene expression. *Nucleic Acids Res.* **38**, 3891-3908, (2010).
- 1169
- 1170 153 Singer, T., McConnell, M. J., Marchetto, M. C., Coufal, N. G. & Gage, F. H. LINE-1 retrotransposons: mediators of somatic variation in neuronal genomes? *Trends Neurosci.* **33**, 345-354, (2010).
- 1171
- 1172
- 1173 154 Singh, M. *et al.* A new human embryonic cell type associated with activity of young transposable elements allows definition of the inner cell mass. *PLoS Biol.* **21**, e3002162, (2023).
- 1174
- 1175
- 1176 155 Mirihana Arachchilage, G., Hetti Arachchilage, M., Venkataraman, A., Piontkivska, H. & Basu, S. Stable G-quadruplex enabling sequences are selected against by the context-dependent codon bias. *Gene* **696**, 149-161, (2019).
- 1177
- 1178
- 1179 156 Varshney, D., Spiegel, J., Zyner, K., Tannahill, D. & Balasubramanian, S. The regulation and functions of DNA and RNA G-quadruplexes. *Nat. Rev. Mol. Cell Biol.* **21**, 459-474, (2020).
- 1180
- 1181
- 1182 157 Wang, Y., Qiu, C. & Cui, Q. A large-scale analysis of the relationship of synonymous SNPs changing microRNA regulation with functionality and disease. *Int. J. Mol. Sci.* **16**, 23545-23555, (2015).
- 1183
- 1184
- 1185 158 Brest, P. *et al.* A synonymous variant in IRGM alters a binding site for miR-196 and causes deregulation of IRGM-dependent xenophagy in Crohn's disease. *Nat. Genet.* **43**, 242-245, (2011).
- 1186
- 1187
- 1188 159 Gartner, J. J. *et al.* Whole-genome sequencing identifies a recurrent functional synonymous mutation in melanoma. *Proc Natl Acad Sci U S A* **110**, 13481-13486, (2013).
- 1189
- 1190
- 1191 160 Hamdorf, M. *et al.* miR-128 represses L1 retrotransposition by binding directly to L1 RNA. *Nat. Struct. Mol. Biol.* **22**, 824-831, (2015).
- 1192

1193 161 Eyre-Walker, A. & Keightley, P. D. The distribution of fitness effects of new
1194 mutations. *Nat. Rev. Genet.* **8**, 610-618, (2007).

1195 162 Fairbrother, W. G., Holste, D., Burge, C. B. & Sharp, P. A. Single nucleotide
1196 polymorphism-based validation of exonic splicing enhancers. *PLoS Biol.* **2**, E268,
1197 (2004).

1198 163 Chamary, J. V. & Hurst, L. D. Evidence for selection on synonymous mutations
1199 affecting stability of mRNA secondary structure in mammals. *Genome Biol* **6**,
1200 R75, (2005).

1201 164 Nackley, A. G. *et al.* Human catechol-O-methyltransferase haplotypes modulate
1202 protein expression by altering mRNA secondary structure. *Science* **314**, 1930-
1203 1933, (2006).

1204 165 Schattner, P. & Diekhans, M. Regions of extreme synonymous codon selection
1205 in mammalian genes. *Nucleic Acids Res.* **34**, 1700-1710, (2006).

1206 166 Wu, X. M. & Hurst, L. D. Why selection might be stronger when populations are
1207 small: intron size and density predict within and between-species usage of
1208 exonic splice associated cis-motifs. *Mol. Biol. Evol.* **32**, 1847-1861, (2015).

1209 167 Prats-Ejarque, G., Lu, L., Salazar, V. A., Moussaoui, M. & Boix, E. Evolutionary
1210 trends in RNA base selectivity within the RNase A superfamily. *Front. Pharmacol.*
1211 **10**, 1170, (2019).

1212 168 Mendell, J. T., Sharifi, N. A., Meyers, J. L., Martinez-Murillo, F. & Dietz, H. C.
1213 Nonsense surveillance regulates expression of diverse classes of mammalian
1214 transcripts and mutes genomic noise. *Nat. Genet.* **36**, 1073-1078, (2004).

1215 169 Torres, M. *et al.* Paraspeckles as rhythmic nuclear mRNA anchorages responsible
1216 for circadian gene expression. *Nucleus (Calcutta)* **8**, 249-254, (2017).

1217 170 Prasanth, K. V. *et al.* Regulating gene expression through RNA nuclear retention.
1218 *Cell* **123**, 249-263, (2005).

1219 171 Lucks, J. B., Nelson, D. R., Kudla, G. R. & Plotkin, J. B. Genome landscapes and
1220 bacteriophage codon usage. *PLoS Comput Biol* **4**, e1000001, (2008).

1221 172 De Vlugt, C., Sikora, D. & Pelchat, M. Insight into influenza: A virus cap-
1222 snatching. *Viruses* **10**, (2018).

1223 173 Jalkanen, A. L., Coleman, S. J. & Wilusz, J. Determinants and implications of
1224 mRNA poly(A) tail size--does this protein make my tail look big? *Semin Cell Dev*
1225 *Biol* **34**, 24-32, (2014).

1226 174 Mauro, V. P. Codon optimization in the production of recombinant
1227 biotherapeutics: potential risks and considerations. *Biodrugs* **32**, 69-81, (2018).

1228 175 Ho, A. T. & Hurst, L. D. Effective population size predicts local rates but not local
1229 mitigation of read-through errors. *Mol. Biol. Evol.* **38**, 244-262, (2021).

1230 176 Allert, M., Cox, J. C. & Hellinga, H. W. Multifactorial determinants of protein
1231 expression in prokaryotic open reading frames. *J. Mol. Biol.* **402**, 905-918, (2010).

1232 177 Wu, X. & Hurst, L. D. Determinants of the usage of splice-associated cis-motifs
1233 predict the distribution of human pathogenic SNPs. *Mol. Biol. Evol.* **33**, 518-529,
1234 (2016).

- 1235 178 Abrahams, L. *et al.* Evidence in disease and non-disease contexts that nonsense
1236 mutations cause altered splicing via motif disruption. *Nucleic Acids Res.* **49**,
1237 9665-9685, (2021).
- 1238 179 Soemedi, R. *et al.* Pathogenic variants that alter protein code often disrupt
1239 splicing. *Nat. Genet.* **49**, 848-+, (2017).
- 1240 180 Mühlhausen, S. & Hurst, L. D. Transgene-design: a web application for the
1241 design of mammalian transgenes. *Bioinformatics*, btac139, (2022).
- 1242 181 Sharp, C. P. *et al.* CpG dinucleotide enrichment in the influenza A virus genome
1243 as a live attenuated vaccine development strategy. *PLoS Pathog* **19**, e1011357,
1244 (2023).
- 1245 182 Yew, N. S. *et al.* CpG-depleted plasmid DNA vectors with enhanced safety and
1246 long-term gene expression in vivo. *Mol. Ther.* **5**, 731-738, (2002).
- 1247 183 Kariko, K., Buckstein, M., Ni, H. & Weissman, D. Suppression of RNA recognition
1248 by Toll-like receptors: the impact of nucleoside modification and the
1249 evolutionary origin of RNA. *Immunity* **23**, 165-175, (2005).
- 1250 184 Vaidyanathan, S. *et al.* Uridine Depletion and Chemical Modification Increase
1251 Cas9 mRNA Activity and Reduce Immunogenicity without HPLC Purification.
1252 *Mol Ther Nucleic Acids* **12**, 530-542, (2018).
- 1253 185 Ohta, T. The nearly neutral theory of molecular evolution. *Ann Rev Ecol System*
1254 **23**, 263-286, (1992).
- 1255 186 Christmas, M. J. *et al.* Evolutionary constraint and innovation across hundreds
1256 of placental mammals. *Science* **380**, eabn3943, (2023).
- 1257 187 Lynch, M. *et al.* Genetic drift, selection and the evolution of the mutation rate.
1258 *Nat. Rev. Genet.* **17**, 704-714, (2016).
- 1259 188 Sung, W., Ackerman, M. S., Miller, S. F., Doak, T. G. & Lynch, M. Drift-barrier
1260 hypothesis and mutation-rate evolution. *Proc Natl Acad Sci U S A* **109**, 18488-
1261 18492, (2012).
- 1262 189 Seczynska, M. & Lehner, P. J. The sound of silence: mechanisms and implications
1263 of HUSH complex function. *Trends Genet.* **39**, 251-267, (2023).
- 1264

1265 **Acknowledgements**

1266 The authors thank the Humboldt Foundation for the award of the Humboldt Prize to LDH to enable him to spend time in Germany. SR is funded by the Evolution Education Trust. [\[Au: do](#)
1267 [you wish to include any Acknowledgements?\]](#)

1268 **Author contributions**

1269 LDH and SR researched data for the article. All authors contributed substantially to discussion of the content. LDH wrote the article. All authors reviewed and/or edited the manuscript
1270 before submission. ZI contributed to transposable element issues specifically. [\[Au: please provide relevant information\]](#)

1271 **Competing interests**

1272 The authors declare no competing interests. [\[Au:OK?\]](#)

1273 **Peer review information**

1274 *Nature Reviews Genetics* thanks the anonymous reviewers for their contribution to the peer review of this work.

1275 **Supplementary information**

1276 Supplementary information is available for this paper at <https://doi.org/10.1038/s415XX-XXX-XXXX-X>
1277

1278

1279

1280 Figure 1 | **Processes acting on transcription, processing and filtering of transcripts with**
1281 **emphasis on roles of codon content or intronic signals.** On the left is the regular gene
1282 expression pathway from transcription to translation, highlighting the features that allow transcripts
1283 to be recognised as native and correctly expressed as proteins. Arrows coming off towards the right
1284 are the features for which transcripts are recognised as “unwanted” and filtered from being translated.
1285 In bold are those steps for which codon content or intronic presence is known to be relevant. Relevant
1286 abbreviations: RNA Pol II (RNA Polymerase II), CFP1 (CxxC Zinc Finger Protein 1), HUSH (Human
1287 Silencing Hub), SALL4 (Spalt-Like) zinc finger protein, ESE (Exonic Splicing Enhancer), SR (Splicing
1288 Regulatory) protein, PROMPTs (Promoter Upstream Transcripts), lncRNA (long non-coding RNA), NATs
1289 (Natural Antisense Transcripts), dsRNA (double-stranded RNA), IRAlus (Inverte Repeat Alu) elements,
1290 ptRNA (prematurely terminated RNA), uaRNA (upstream antisense RNA), PAXT (PolyA Exosome
1291 Targeting) complex, RNase (ribonuclease), NEXT (Nuclear Exosome Targeting) complex, TREX
1292 (Transcription Exon) complex, NXF1/Tap (Nuclear RNA Export Factor 1), P bodies (Processing bodies),
1293 ZAP (Zinc-finger Antiviral Protein), NMD (Nonsense Mediated Decay), Dicer (endoribonuclease Dicer),
1294 Dhh1p (DEAD/H-box protein) Ccr4-Not (Carbon Catabolite Repression-Negative On TATA-less)
1295 transcription complex.

1296

1297 Figure 2 | **Dinucleotide frequencies of human coding sequences, single-exon coding**
1298 **sequences and genomic DNA, and as expected from mononucleotide frequencies. a,**
1299 The data are ordered according to the difference between dinucleotide frequency in the
1300 genome and that in coding sequence. Note that the four dinucleotides bearing only A and/or
1301 T are all rare in coding sequence, whereas the four dinucleotides bearing only G and/or C are
1302 the most enriched in coding sequence compared with genomic DNA. CpG is unique amongst
1303 these four in that its frequency in coding sequence is less than would be expected from its
1304 mononucleotide content, most probably owing to the mutagenicity of CpG. **b,** The absolute
1305 difference in dinucleotide content between coding sequence and genomic DNA. Bars are
1306 colour coded by G and/or C content, with darker colour indicating higher GC content. Data
1307 derived by the authors from human genome release GRCh38.

1308

1309 Figure 3 | **Native human exons predicted to be targeted by the HUSH complex.** Exons
1310 longer than ~1500bp and with an A content greater than ~27% are targeted for silencing by
1311 the HUSH complex (from Extended Data Fig. 3c of ref.¹⁰⁰). This population is shown as the
1312 top right quadrant and accounts for less than 0.4% of native exons and about 3.5% of native
1313 human genes. We highlight a few zinc finger (ZNF) genes and a proto-cadherin. Single-exon
1314 coding sequence genes are shown in yellow. Note the rarity of human exons likely to be

1315 subject to HUSH. This is consistent with both HUSH largely targeting foreign genes and
1316 forcing selection on the anatomy of native genes. Note too the commonality of ZNF genes in
1317 the group that are likely to human silenced by HUSH suggesting that HUSH and ZNF comprise
1318 alternative transposable element suppression systems, the one regulating the other. Data
1319 derived by the authors from human genome release GRCh38.

1320

1321 **Figure 4 | The mode of action of disease-associated synonymous mutations.**
1322 Synonymous mutations can interfere with post-transcriptional processing by disrupting
1323 constitutive splice sites or activating cryptic splice sites, therefore altering SR protein binding
1324 and the production of correctly spliced mature mRNAs. Synonymous mutations can also
1325 influence mRNA secondary structure and stability and miRNA-mediated gene regulation, in
1326 turn altering downstream protein expression. During translation, synonymous mutations can
1327 alter ribosome processivity, influencing translational kinetics and thus co-translational protein
1328 folding. Production of aberrant or truncated proteins can result in non-functional proteins or, in
1329 most cases, degradation. Examples of pathogenic synonymous mutations are indicated (in
1330 bold) next to their proposed mechanism of pathogenicity. See **Supplementary Table 1** for full
1331 list and references. Relevant abbreviations: RNA Pol II (RNA Polymerase II), ESE (Exonic
1332 Splicing Enhancer), NMD (Nonsense Mediated Decay).

1333

1334 **Box 1 | The nearly-neutral theory and the importance of effective population size**

1335 What is the fate of mutations that have little or no effect on fitness? One of the major advances
1336 in theoretical population genetics has been the realization that mutations with little effect on
1337 fitness have a different fate to those with no effect. The fate of mutations that have no effect
1338 is explained by the neutral (or strictly-neutral) theory, whereas the fate of mutations with little
1339 effect is explained by the nearly-neutral theory¹⁸⁵. Both theories assume that random changes
1340 in allele frequency (known as drift) are an influential process.

1341

1342 According to the strictly-neutral theory, the probability that one of two neutral alleles will drift
1343 to fixation is its current frequency. Thus, when a new mutation arrives in a diploid organism, it
1344 will be at a frequency $1/2N$, where N is the population size, and its probability of fixation is
1345 $1/2N$. The rate at which such mutations arrive is $2N \times mu$, where mu is the per generation
1346 mutation rate. Thus, the net rate of evolution is $2N \times mu/2N = mu$. Importantly, the rate of
1347 strictly-neutral evolution is independent of the population size.

1348

1349 In the nearly-neutral model, mutations reduce fitness by a very small degree, s . The key
1350 postulate of the nearly-neutral model is that as the population size increases (or, rather, the
1351 effective population size, N_e), selection is a more effective force, meaning that the chance that

1352 a deleterious mutation will be eliminated increases with increasing N_e . This is because in a
1353 small population, a mutation with small effect requires fewer chance increases in frequency to
1354 go from being rare to common, whereas in a large population, selection has longer (in terms
1355 of the number of generations) to remove the deleterious mutation from the population before
1356 it reaches fixation by chance.

1357

1358 A key question concerns the relationship between s , N_e and the rate of fixation. The nearly-
1359 neutral theory divides deleterious mutations into one of three classes. Some mutations are of
1360 such small effect that their rate of fixation is so close to μ that they can't be distinguished
1361 from strictly-neutral mutations. These 'effectively neutral' mutations occur when $s \ll 1/2 N_e$. A
1362 further class of mutations are so deleterious that they would never reach fixation; these
1363 'strongly deleterious' mutations occur when $s \gg 1/2 N_e$. The third class are those intermediary
1364 mutations that can reach fixation but do so at a slower rate than the mutation rate. These
1365 'weakly deleterious' mutations occur when $s \cong 1/2 N_e$.

1366

1367 A key result of the nearly-neutral model is that, if we assume a mutation has a small deleterious
1368 effect on fitness, its chance of fixation increases when N_e is small as it is more likely to be
1369 within the 'effectively neutral' or 'weakly deleterious' mutation classes. For example, insertion
1370 of a small segment of DNA with little phenotypic consequence could well be 'strongly
1371 deleterious' in a species with a large N_e and hence removed from the population, but
1372 'effectively neutral' or 'weakly deleterious' in a species with a small N_e and hence not always
1373 removed, despite s being the same in all cases. In this way, the nearly-neutral model explains¹¹
1374 why the human genome contains so much unconstrained 'junk' DNA (only about 10% of
1375 human DNA is constrained¹⁸⁶), why we have large introns that are hard to splice accurately
1376 and why much DNA is prone to spurious binding of transcription factors. In each case we
1377 evoke the idea that a harmful mutation happens, such as insertion of a transposable element,
1378 a mutation decreasing splicing fidelity or one providing a spurious TF binding site etc.. The
1379 fate of such weak effect mutations will be dependent on the efficiency of selection being able
1380 to remove the mutation from the population, efficiency that is lower when N_e is small. Junk
1381 DNA is DNA that is 'weakly deleterious' and thus cannot always be prevented from reaching
1382 fixation by drift. The model evokes a balance between drift and selection to keep genomes
1383 free of 'strongly deleterious' mutations. This drift–selection balance can, in addition, explain
1384 why the mutation rate is higher when N_e is low^{187,188}.

1385 Box 2 | **The HUSH complex**

1386 The human silencing hub (HUSH) is a transcriptional silencing system, constitutively
1387 expressed in most cell types, that, amongst other activities, safeguards against retroelement

1388 invasion and controls the flow of genetic information from RNA to DNA in mammalian
1389 genomes. HUSH represses retroviruses and retroelements (non-self) while largely avoiding
1390 the inappropriate repression of intron-containing host genes (self) (for review, see ref.¹⁸⁹).

1391

1392 HUSH is a complex of M-phase phosphoprotein 8 (MPP8), transcription activation suppressor
1393 (TASOR) and periphilin. In the canonical pathway, periphilin binds RNA and then guides HUSH
1394 to its target loci. At the targets, HUSH promotes deposition of histone 3 lysine 9 trimethylation
1395 (H3K9me3), as well as microorchidia CW-type zinc finger 2 (MORC2) and the nuclear exosome
1396 targeting (NEXT) complex. Although, canonically, HUSH is recruited via RNA (and hence
1397 requires some expression¹⁰⁰), it can also be recruited to DNA by DNA-binding protein nuclear
1398 protein 220 (NP220).

1399

1400 HUSH targets tend to be long exons¹⁰⁰ (>1.5kb) of either multi-exon genes or long intronless
1401 genes, neither of which are common for human transcripts. In addition, it targets mRNAs with
1402 high adenine (A) content. This is thought to enable the targeting of retroviruses and
1403 retrotransposons as retroposition is associated with high A content. Through its ability to
1404 distinguish self from non-self without the need for prior exposure (**Fig. 3**), HUSH may be
1405 considered a component of the innate immune system and intronless cDNA may be considered
1406 a pathogen-associated molecular pattern.

1407

1408 Given its mode of pattern recognition, HUSH suppresses retroelements from outside the cell
1409 (retroviruses such as HIV) as well as from within the cell (retrotransposons). HUSH deficiency
1410 phenotypes are likely owing to the release from repression of LINE1 transposons. HUSH thus
1411 guards against the activation of interferon and other immune genes that would otherwise be
1412 stimulated by LINE1 and associated double-stranded RNA (sensed by cytoplasmic RIG-I-like
1413 receptors).

1414

1415 Although HUSH is a system that has evolved to target mobile elements — and lentiviruses
1416 have co-evolved mechanisms to degrade HUSH by Vpr and Vpx — this suppressor system
1417 can be co-opted to control host genes. Native gene targets of HUSH include KRAB-zinc finger
1418 proteins (encoded by a single exon of 4.5kb) and some protocadherins (encoded by a single
1419 exon >2.5kb). The HUSH-mediated suppression of protocadherins is important for neuronal
1420 development.

1421

1422 HUSH can also repress recombinant adeno-associated virus (AAV), which is commonly used
1423 as a delivery vector in gene therapy. As HUSH recognizes unusually long exons and intronless
1424 transcripts, it presents an impediment to gene therapy for larger genes. That introns enable

1425 escape from HUSH-dependent silencing through splicing-independent routes (unspliced
1426 mutant introns also enable avoidance of HUSH) suggests that it may be possible to design
1427 transgenes that evade this silencing.

1428

1429 **Glossary**

1430 GC-biased gene conversion

1431 (gBGC). Biased gene conversion describes the recombination of short stretches of genetic
1432 material from a donor sequence to an acceptor sequence that is biased as to which is the
1433 donor strand and which is the acceptor strand. In gBGC, AT:GC mismatches in recombinant
1434 sections are resolved with a preference towards G or C.

1435

1436 Codon usage bias

1437 A pattern in which synonymous codons occur at frequencies different from some null
1438 expectation (often the null expectation is of equal frequency).

1439

1440 CpG islands

1441 Long (>200bp) stretches of DNA rich in CpG dinucleotides that commonly occur in the vicinity
1442 of mammalian promoters. They are commonly unmethylated if the associated gene is
1443 expressed.

1444

1445 Cryptic splice sites

1446 Splice sites within exons away from the canonical splice site that can result in alternative
1447 splicing.

1448

1449 Effective population size

1450 (N_e). The number of individuals that an idealized neutrally evolving population would require
1451 for it to have properties (such as genetic diversity) equivalent to the real population. N_e is often
1452 smaller than the census population size (N). In simple cases, it approximates to the number
1453 of breeding individuals.

1454

1455 Exonic splice enhancers

1456 (ESEs). Short (6–8bp) RNA motifs that enable more accurate splicing, particularly in the
1457 vicinity of weak splice sites. ESEs often function as RNA binding sites for serine–arginine-rich
1458 (SR) proteins.

1459

1460 Exon-junction complex

1461 Protein complex that forms on a pre-mRNA molecule at the junction between two exons joined
1462 by RNA splicing.

1463

1464 Exosome complex

1465 A multi-protein intracellular protein complex that degrades many types of RNA. In eukaryotes
1466 it is present in the cytoplasm, the nucleus and the nucleolus. In humans the cytoplasmic form
1467 is associated with the DIS3-like exonuclease 1 (DIS3L), while the nuclear complex employs DIS3.
1468 The nuclear form is termed the nuclear exosome complex, that which is targeted by NEXT
1469 complex. The (RNA) exosome complex is not to be confused with exosomes, extracellular vesicles
1470 generated by cells.

1471

1472 G-quadruplexes

1473 Stable secondary structures formed in guanine-rich DNA and RNA sequences.

1474

1475 Inverted repeat Alu elements

1476 Single-stranded Alu elements that are followed downstream by their reverse complement. This
1477 characteristic often allows inverted repeats to fold on themselves and form double-stranded
1478 structures.

1479

1480 LINE1

1481 (L1). LINE elements, of which L1 is an example, are the most abundant class of transposable
1482 elements in the human genome, formed by autonomous retrotransposition. [\[Au:OK?\]](#)

1483

1484 Mutational equilibrium

1485 The frequency (for example, of nucleotides) in a population at which, if only mutation bias and
1486 neutral evolution affect the frequency, the frequency does not change.

1487

1488 Neutral evolution

1489 The process by which allele frequency is determined by chance events alone operating on
1490 alleles of the same fitness.

1491

1492 Nonsense-mediated decay

1493 (NMD). Process by which mRNA molecules containing premature stop codons are degraded.

1494

1495 No-go mRNA decay

1496 (NGD). Process by which RNA with stacks of stalled ribosomes are degraded.

1497

1498 Codon optimality-mediated RNA decay
1499 (COMD). Process by which mRNA with abundant non-optimal codons (those matching rare
1500 tRNAs) are subject to decay, probably mediated by slow ribosomal progression.
1501
1502 Non-stop RNA decay
1503 (NSD). Process by which mRNA without a proper stop codon are identified and targeted for
1504 decay. In eukaryotes, NSD discharges ribosomes stalled at the 3' end, directing the mRNA to
1505 the exosome complex.
1506
1507 Processing bodies
1508 Ribonucleoprotein bodies found in the cytoplasm that retain mRNA molecules and contain
1509 proteins required for mRNA decay. A similar role is carried out by cytoplasmic stress granules.
1510
1511 Retrogenes
1512 Processed copies of genes formed from reverse transcription of mature mRNA molecules of
1513 a parental gene (hence without introns).
1514
1515 Spurious transcripts
1516 Transcripts generated by functionally irrelevant cellular events (such as TF binding to random
1517 sequence).
1518
1519 Synonymous mutations
1520 Single base-pair mutations in a protein-coding exon that change the codon to a different but
1521 synonymous one and hence do not a priori change the amino acid sequence of the encoded
1522 protein.
1523
1524 Translational selection
1525 Selection that favours synonymous mutations owing to their effects on translation, typically
1526 assumed to be mediated by faster or more accurate translation. Commonly evidenced by a
1527 codon usage bias that favours synonymous codons matching more abundant iso-acceptor
1528 tRNAs.
1529
1530 Transposable elements
1531 DNA sequences that can move within genomes and replicate without depending on gene
1532 replication of the host cell.
1533
1534 Zinc finger antiviral protein

1535 (ZAP). Protein involved in preventing retroviral infection by binding of CpG-rich sequences
1536 and recruitment of the RNA degradation machinery.

1537

1538

1539 **Table of Contents**

1540 Multiple mechanisms have evolved to prevent or trap deleterious unwanted transcripts. The
1541 unwanted transcript hypothesis proposes that selection at synonymous sites favours
1542 mutations that prevent the generation of transcriptional rubbish or that make native transcripts
1543 look “wanted” by being GC-rich.

1544

1545