

1 **Transfer RNA genes experience** 2 **exceptionally elevated mutation rates**

3 **Bryan Thornlow^a, Josh Hough^a, Jackie Roger^a, Henry Gong^a, Todd Lowe^{a,b,1}, and Russell**
4 **Corbett-Detig^{a,b,1}**

5 ^aDepartment of Biomolecular Engineering, University of California Santa Cruz, Santa Cruz CA 95064; ^bGenomics Institute,
6 University of California Santa Cruz, Santa Cruz CA 95064

7 ¹To whom correspondence should be addressed. E-mail: lowe@soe.ucsc.edu or russcd@gmail.com

8 This manuscript was compiled on December 6, 2017

9 H.G. performed initial analyses. B.T. performed research and analyzed data, with assistance from J.R.. J.H. performed the DFE analyses. B.T.,
10 T.L. and R.C.-D. conceived and designed research. B.T., T.L. and R.C.-D. wrote the paper. J.H., J.R., H.G., T.L. and R.C.-D. edited the paper.

11 The authors declare no conflict of interest.

12 **Transfer RNAs (tRNAs) are a central and necessary component for the biological synthesis of new**
13 **proteins, and they are among the most highly conserved and most frequently transcribed sequences**
14 **across all of life. Despite their clear significance for fundamental cellular processes, however, the forces**
15 **governing tRNA evolution are poorly understood. Here, we present evidence that transcription-**
16 **associated mutagenesis and strong purifying selection are key determinants of patterns of sequence**
17 **polymorphism and divergence within and surrounding tRNA genes across several diverse model**
18 **organisms. Remarkably, our results indicate that the mutation rate at broadly expressed tRNA loci is**
19 **between 8.7 and 13.8 times greater than the genome-wide average. Furthermore, evolutionary analyses**
20 **provide strong evidence that tRNA loci, but not their flanking sequences, experience strong purifying**
21 **selection, acting in direct response to this elevated mutation rate. Finally, we also find a highly**
22 **significant correlation between tRNA expression levels and the mutation rates in their immediate**
23 **flanking regions, suggesting the possibility of predicting gene expression levels based on relative**
24 **mutation rates and sequence variation data among tRNA gene loci. Our results provide novel insight**
25 **into individual tRNA gene evolution, and imply that tRNA loci contribute disproportionately to**
26 **mutational load in human populations.**

27 **Significance Statement**

28 tRNAs are essential for the production of all proteins in all tissues across life and are therefore
29 among the most highly transcribed loci in the genome. Our study shows that the frequent
30 transcription of tRNAs results in a highly elevated mutation rate at tRNA loci that is between 8-
31 and 14-fold higher for tRNAs than for the rest of the genome. We also show that the strength of
32 natural selection, which acts to remove sequence-altering mutations, is extremely strong in
33 tRNAs, but is relaxed in introns and regions flanking tRNAs. Finally, our results indicate that
34 mutation rates in non-functional tRNA flanking regions are similarly elevated, and levels of
35 genetic variation correlate strongly with expression. These observations suggest that a predictive
36 model could facilitate future studies of tRNA function.

37 Transfer RNAs (tRNAs) play an essential role in protein synthesis across all of life. Their

38 primary function is in the translation of the genetic code into the corresponding amino acid
39 sequences that make up proteins. Thus, tRNA molecules are critical for virtually all cellular
40 processes, and the genes encoding tRNA molecules have been highly conserved over
41 evolutionary time (1). The necessity of tRNAs in large quantities also makes them among the
42 most highly transcribed loci in the genome. Indeed, many tRNA genes may experience greater
43 levels of transcription than even the most highly transcribed protein-coding genes (2, 3). Such
44 high levels of transcription suggest that tRNA genes may experience high levels of transcription-
45 associated mutagenesis (TAM) compared to the rest of the genome, making the tRNA gene
46 family an excellent model system for studying the interplay between natural selection and
47 elevated mutation rates.

48 tRNA | transcription | mutagenesis | TAM

49 Transcription affects the mutation rates of transcribed genes (4) through the unwinding and
50 separation of complementary DNA strands (5). In particular, during transcription, a nascent RNA
51 strand forms a hybrid DNA-RNA complex with a template DNA strand. While the
52 complementary tract of non-template DNA is temporarily isolated, it is chemically reactive and
53 thus accessible by potential mutagens (5). In addition, if transcription and DNA replication occur
54 concomitantly at a particular locus, collisions between RNA Polymerase and the DNA
55 replication fork are possible, which may also result in damage to DNA (6). Several cellular
56 agents have also been shown to cause damage in highly expressed genes (7). Among the most
57 notable sources of mutation associated with high transcription is activation-induced cytidine
58 deaminase (AID) (8). AID accompanies RNA Polymerase II and converts cytosine to uracil,
59 causing a relative excess of cytosine to thymine substitutions in the non-template strand and
60 guanine to adenine substitutions on the template strand (9). More highly transcribed genes are
61 especially vulnerable to mutation by AID, making it a clear diagnostic signature of TAM (4).
62 Because tRNA loci are so often unwound for transcription, these regions are therefore expected
63 to experience elevated mutation rates due to TAM, with deamination via AID being one of the
64 primary mutational mechanisms.

65 In order to conserve mature tRNA sequence identity in the presence of an elevated mutation rate,
66 it is expected that tRNA genes should experience strong selective pressures. At the gene level,
67 tRNA transcription requires sequence-specific binding of transcription factors to the internal A
68 and B box promoter elements (10). Once transcribed, precursor tRNAs must fold properly to
69 undergo a complex process of maturation, which can be disrupted at any step by sequence-
70 altering mutations. The unique structure of tRNAs dictates processing by RNases, addition of an
71 assortment of modifications, accurate recognition by highly specific aminoacyl tRNA
72 synthetases, incorporation into the translating ribosome, and accurate positioning of the
73 anticodon relative to mRNA codons (11, 12). As a consequence of the need to maintain high
74 sequence-specificity, DNA encoding the mature portions of tRNAs are exceptionally well
75 conserved segments of the genome (11). Therefore, we expect that a large proportion of
76 mutations arising in tRNA genes will be deleterious, and therefore experience strong purifying
77 selection.

78 While most human tRNA genes do not have external promoters (10, 11), tRNA transcripts
79 generally include leader and trailer sequences, extending roughly 2-5 nucleotides upstream of the
80 annotated mature tRNA gene, and 5-15 nucleotides downstream of the mature tRNA sequence,

81 based on the position of the genomically encoded poly-T transcription termination sequence.
82 However, these sequences have limited functionality in most cases (13–16). For example, very
83 early in the tRNA maturation process, all tRNAs undergo removal of their 5' leader sequences
84 by RNase P (13, 14) and removal of their 3' trailers by RNase Z (17). Because these flanking
85 sequences are frequently unwound and therefore vulnerable to TAM, we expect that tRNA
86 flanking regions will experience similar mutation rates to tRNA genes. Whereas tRNAs should
87 experience purifying selection, we expect that flanking regions should be nearly neutral or under
88 very weak selection.

89 Despite these clear predictions for an elevated mutation rate of tRNA loci and strong purifying
90 selective pressure on tRNA genes compared to their flanks, there has been no attempt to quantify
91 the overall impact of mutation and selection on patterns of sequence variability in tRNAs or their
92 flanking regions. Here we investigate the patterns of conservation, divergence and within-species
93 variation of tRNAs in humans and several other model organisms (*Mus musculus*, *Arabidopsis*
94 *thaliana*, and *Drosophila melanogaster*).

95 **Results & Discussion**

96 **Flanking regions of tRNA genes are highly variable despite strong conservation of mature**
97 **tRNA sequences.** To estimate evolutionary conservation, we averaged phyloP data, a measure of
98 the conservation of each human genomic position across 100 vertebrate species (18), by position
99 within each tRNA locus (see Methods). To study the effects of evolution on a shorter time-scale,
100 we also aligned the human and *Macacca mulatta* (Rhesus macaque) genomes and counted the
101 non-gap nucleotide mismatches in the alignments as divergent positions for each tRNA locus.
102 Our analyses indicate that mature tRNA sequences are highly conserved at all positions, based
103 on both average phyloP score (18) (Figure 1A, Supplementary Table 1) and *M. mulatta*
104 alignment (Figure 1B). However, the “inner” flanking regions, defined based on inflection points
105 in the data (see Methods), show significant divergence by the same measures. The inner 5'
106 flanking region, defined as the 20 bases upstream of the tRNA, is the most divergent segment of
107 these regions on average, with roughly four times the rate of divergence between human and *M.*
108 *mulatta* as the untranscribed reference regions (Figure 1B). We found similarly increased rates of
109 divergence in the inner 3' flanking region, which was roughly three times as divergent between
110 human and *M. mulatta* as the untranscribed reference regions (Figure 1B). Both the outer 5'
111 flank (21-40 bases upstream of the tRNA) and outer 3' flank (11-40 bases downstream of the
112 tRNA) are also roughly 1.5 times as divergent as are the untranscribed reference regions.
113 Furthermore, we find that intergenic regions within clusters of active tRNAs (Figures S1A, S1B)
114 show similar patterns in their phyloP scores, with increased divergence extending hundreds to
115 thousands of bases up and downstream of each tRNA gene. Conversely, we find that the
116 intergenic regions in clusters of inactive tRNA genes do not show this pattern nearly as strongly
117 (Figure S1C).

118 To focus our evolutionary timescale further and eliminate any bias due to multi-species sequence
119 alignment errors, we also studied relative levels of tRNA variation within a human population by
120 observing the occurrence of low-frequency single nucleotide polymorphisms (SNPs) (minor
121 allele frequency < 0.05%) for each tRNA gene locus. Consistent with our phyloP and *M. mulatta*
122 divergence analyses, we find that low-frequency SNPs are more common across the entire tRNA

123 locus, including the mature sequence and flanking regions, relative to untranscribed reference
124 regions (Figure 1C). Although the inner 5' and inner 3' flanking regions are the most
125 polymorphic, the mature tRNA sequences have about twice as many low-frequency SNPs per
126 site as untranscribed reference regions. Overall, our results are remarkably consistent on multiple
127 timescales (across vertebrates, between primates, within human populations), indicating that
128 functional tRNA sequences are highly conserved across species but prone to mutations at the
129 individual level, and tRNA flanking regions are both more divergent and more polymorphic than
130 untranscribed, non-genic sequences.

131 **Transcription contributes to variation in tRNA and flanking regions.** We observed that
132 conservation and divergence patterns varied between tRNA loci, and hypothesized that highly
133 active tRNA genes would show the greatest mutation rates if transcription-associated
134 mutagenesis is a primary driver of variation among tRNA loci. Because tRNA transcript
135 abundance measures are often not attributable to individual loci for multiple reasons, including
136 redundant gene copies, variation in pre-tRNA processing and tRNA degradation rates, and
137 difficulty sequencing full-length tRNAs, we estimated relative transcriptional activity based on
138 chromatin state data from the Epigenomic Roadmap Project ((21); Cozen et al., in preparation).
139 We classified human tRNA genes as "active" if they were in regions of active chromatin and
140 near transcription start sites in at least 3% of the 127 tissues for which genome-wide epigenomic
141 data was available (Figure 2A, see Methods). We considered the remaining tRNAs "inactive".
142 No isoforms are over-represented in any of the epigenomic groups (Cozen, et al., in preparation).
143 We found that active tRNAs were significantly more conserved than inactive tRNAs (Mann-
144 Whitney U, $p < 8.40e-53$), and the flanking regions of active tRNAs were significantly more
145 divergent than the flanking regions of inactive tRNAs ($p < 7.98e-61$). Indeed, the peak measure
146 of divergence in the inner 5' flanking regions is roughly five times greater in active tRNAs than
147 in inactive tRNAs (Figure S2A). Active tRNAs also had significantly more low-frequency
148 polymorphisms per site than inactive tRNAs across the entire locus, including the tRNA and
149 flanking regions ($p < 3.72e-36$). Inactive tRNAs were still significantly more conserved ($p <$
150 $2.02e-12$) and polymorphic ($p < 0.007$) than the untranscribed reference regions, and their flanks
151 were significantly more divergent than the reference regions ($p < 1.36e-16$).

152 That the peak by all three measures is consistently about 12 to 13 nucleotides upstream of the
153 mature tRNA sequence is a curious result. At the most divergent position, roughly 55% of all
154 tRNA loci showed a difference between human and *M. mulatta* (Figure 1B) and roughly 15% of
155 tRNA loci have a low-frequency SNP at this site (Figure 1C). Furthermore, among tRNA loci
156 believed to be active, virtually all loci showed a change at this nucleotide between human and *M.*
157 *mulatta*, and roughly 25% have a low-frequency SNP at this site (Figure S2A, S2B). This implies
158 that this region either does not face uniform selective pressures or is not uniformly vulnerable to
159 TAM. While it has been suggested that distant flanking sequences may affect tRNA expression
160 levels in yeast (22), few studies to our knowledge have shown that the immediate flanking
161 regions have an effect on expression in humans or other higher eukaryotes (23). Importantly, the
162 duration of transcription initiation is long relative to the process of transcription itself (24, 25),
163 which would presumably lead to prolonged isolation of the non-template DNA strand at the
164 initiation site and increased vulnerability to TAM. A poised initiation complex might also
165 increase the likelihood of collisions between Pol3 and the replication fork (6). Thus, frequent
166 initiation at highly transcribed tRNA loci could contribute to the pattern of variation we observe.

167 This may also explain the increased variation in the outer 3' flank relative to the outer 5' flank,
168 as positioning of downstream transcription termination sites is highly variable among tRNA
169 genes (19, 26), whereas transcription start site positions are more consistent. Indeed, we find that
170 the TATA boxes for tRNA-SeC-TCA-1-1, RNase P and U6 RNA are all approximately 25
171 nucleotides upstream of the start of the gene (27). While most tRNAs do not have clear TATA
172 boxes, the TATA-Binding Protein (TBP) still binds non-specifically to the DNA duplex at this
173 position (28), which seems to coincide with the sudden decrease in variability. Furthermore, we
174 find that, while both flanking regions for many other Pol3-transcribed genes are divergent, the 5'
175 flanking regions are generally more divergent than the 3' flanking regions, suggesting that the
176 underlying mechanism is not tRNA-specific (Supplementary Table 1). However, additional
177 studies will be necessary to conclusively support the assertion that this strong mutation pattern is
178 due entirely or in large part to the process of transcription rather than due to a correlated process.

179 Two additional and orthogonal analyses strengthen the observed correlations between gene
180 expression and variation at tRNA loci. First, we found a significant correlation between the
181 TATA-Binding Protein (TBP) intensity peaks (29–31) (see Methods) and the average level of
182 divergence in the flanking regions (Spearman's $\rho = -0.64$, $p < 2.2e-16$) (Figure 2D), as well as
183 a correlation between the peaks and the average level of conservation of the mature tRNA
184 sequence (Spearman's $\rho = 0.64$, $p < 2.2e-16$) across all human tRNAs (Figure 2C). The TBP
185 peak data for these transcription factors provide an estimate of the level of transcription for each
186 tRNA, and are consistent with the idea that more highly transcribed tRNAs show higher levels of
187 variability in their transcribed regions.

188 Second, we found significant correlations between the mature tRNA sequence read counts and
189 tRNA conservation (Spearman's $\rho = 0.18$, $p < 0.001$) and flanking region divergence
190 (Spearman's $\rho = -0.61$, $p < 2.2e-16$) when we exclude mature tRNA sequences encoded for by
191 more than one gene (Figure 2E,F), as well as when we sum the average levels of tRNA
192 conservation (Spearman's $\rho = 0.12$, $p < 0.027$) and flanking region divergence (Spearman's
193 $\rho = -0.68$, $p < 2.2e-16$) for genes encoding identical tRNAs to account for correlations between
194 read count and gene copy-number (22, 32) (Figure S3). These read counts were collected from a
195 single human embryonic kidney cell line by Zheng et al (32) using DM-tRNA-seq, a specialized
196 sequencing method developed for tRNAs which overcomes modifications that impede standard
197 small RNA sequencing methods.

198 **Patterns of divergence and conservation can be leveraged to develop a predictive model for**
199 **tRNA gene expression.** Regardless of whether tRNA expression is estimated based on
200 epigenetic chromatin marks across many cell types, TBP transcription factor occupancy across
201 multiple cell lines, or by relative transcript abundance within one cell line, we find highly
202 significant correlations between gene expression and tRNA conservation, flanking region
203 divergence, and tRNA locus polymorphism. The consistency of these correlations indicates that
204 it may be possible to predict tRNA expression based solely on DNA sequence conservation
205 patterns. Genome-wide chromatin-IP and ChIP-seq data are resource-intensive to collect. As
206 sequencing technology is rapidly becoming more affordable and accessible, the prospect of
207 making predictions of tRNA gene expression levels through analysis of multiple alignments and
208 variant sites within populations is enticing. Creating and refining such a model would make
209 future tRNA gene annotation significantly easier and cost-effective.

210 Applicability of this proposed tool is likely best suited for tRNAs, other Pol3 genes, and unique
211 classes of highly expressed protein coding genes such as histones. For example, we find that
212 among the shortest histone protein coding genes (cutoff less than 1,000 nucleotides in length),
213 the average phyloP score per nucleotide is 3.4485, indicating a comparable level of conservation
214 to tRNA genes. Consistent with tRNA genes, their immediate 5' flanking regions are also more
215 divergent than are their immediate 3' flanking regions, on average. However, most genes
216 transcribed by RNA Pol2, including protein coding genes, lincRNAs, miRNAs, snoRNAs, and
217 others, generally do not appear to be good targets based on analysis of representatives of each.
218 For example, glyceraldehyde-3-phosphate dehydrogenase and ribosomal proteins are very highly
219 and very broadly transcribed (33). These genes have extremely well conserved exons, but their
220 introns and flanking regions are not nearly as divergent as tRNA flanking regions, based on
221 phyloP data (18, 27). It is possible that high intron and flanking region divergence in protein-
222 coding genes is still indicative of a high transcription rate, but the degree of variation in these
223 genes occupies a much smaller range, and would therefore be more difficult to incorporate into a
224 model. Additionally, microRNAs such as miR-21 and miR-25 are highly conserved and highly
225 abundant (27, 33), but they are processed out of longer pri-miRNA transcripts, and do not show
226 highly divergent flanking regions at fixed upstream positions, based on phyloP data
227 (Supplementary Table 1) (18, 27). That tRNAs are the best examples for studying signatures of
228 TAM can be attributed to the combination of several unique characteristics, including
229 consistently predictable transcript start and end sites, internal promoters, and extremely high
230 transcription rates. Other highly transcribed genes have conserved functional elements in their
231 flanking regions that may obscure the effects of TAM at these loci.

232 **tRNA flanking regions are among the least conserved sites in the human genome.** Upon
233 scanning the human genome for the least conserved base pairs, we found 247 sites in the genome
234 that had scores of -20, the lowest possible score on the phyloP scale (18, 34). Fifteen of these
235 sites were within 20 base pairs of an active tRNA, based on chromatin-IP data. Of these, 14 sites
236 were found in the inner 5' flanking region of the tRNA, between 10 and 15 base pairs upstream
237 of the first base of the mature tRNA sequence. We found that this set of minimum-phyloP-score
238 sites was enriched for sites within tRNA flanking regions (Hypergeometric test, $p < 1.65e-48$),
239 indicating that the least conserved sites in the genome are disproportionately found in tRNA
240 flanking regions. This indicates that the flanking regions of some active tRNA genes are among
241 the least conserved regions, and perhaps have among the highest mutation rates, of any in the
242 genome.

243 **Patterns of low-frequency SNPs indicate transcription-associated mutagenesis (TAM).** Prior
244 studies of TAM in protein coding genes indicate that transcription is a mutagenic process, in that
245 the untranscribed strand becomes more vulnerable to damage, either through collisions between
246 the DNA replication fork and RNA Polymerase, or by other molecules such as deaminases (4, 7,
247 9). It has been shown that repair pathways activated in response to deaminations lead to excess
248 conversions between guanine and adenine and between thymine and cytosine nucleotides on the
249 coding strand (4, 9). To test this prediction, we analyzed the relative frequencies of all low-
250 frequency SNPs for each region of tRNA loci. Across all tRNA loci, we found that the most
251 common low-frequency SNPs are C→T, G→A, T→C and A→G (transitions), and that these
252 mutations are significantly more common in both tRNA flanking regions and the tRNA gene,
253 relative to untranscribed reference regions (Fisher's exact test, $p < 0.05$ for all comparisons)

254 (Figure 3). Furthermore, the relative excesses of these SNPs are most pronounced in active tRNA
255 loci (Figure S4A). In contrast, and consistent with observed levels of divergence, these relative
256 changes are barely discernible when considering only inactive tRNA loci (Figure S4B).

257 It is important to note that, due to the necessity of preserving tRNA secondary structure, we
258 would expect transition mutations (e.g., A-U to G-U base pairs, C-G to U-G base pairs) to be
259 more common than transversions, regardless of the underlying mechanism, as they should impair
260 function less often. However, the strong mutational skew expected of regions affected by TAM
261 is even more pronounced in regions flanking tRNAs. While some pre-tRNAs may have extended
262 secondary structure that could influence the relative SNP frequencies, such pre-tRNA 5' leader
263 sequences tend to be a maximum of five nucleotides long in mammals (unpublished
264 observations).

265 Prior studies have implicated that CpG sites are significantly more prone to mutations than other
266 nucleotides (35). Therefore, to determine whether TAM was the primary cause of these relative
267 excesses, we repeated our analysis after excluding all CpG sites. We found that CpG sites had no
268 effect on the substitution patterns that we observed in the polymorphism data (Figure S5).

269 **tRNA flanking region variation in other model organisms is consistent with variation**
270 **observed in humans.** To test whether the patterns of polymorphism and divergence that we
271 observed in tRNAs and flanking regions also occurred in other species, we repeated our analyses
272 for tRNAs in *Mus musculus*, *Drosophila melanogaster* and *Arabidopsis thaliana*. Consistent
273 with our results from human data, we found similar patterns of sequence conservation of tRNA
274 loci across all species investigated (Figure S6). In particular, mature tRNA sequences were
275 highly conserved and the flanking regions were highly divergent (Figures S6A, S6D).
276 Particularly striking are the similarities in the outgroup comparisons in the inner 5' flank
277 (Figures S6B, S6E, S6G). The 5' flanks were more divergent than the 3' flanks and the most
278 divergent sites were roughly 10-15 bases upstream of the tRNA in all species. These results are
279 consistent with our human data (Figure 1) and suggest the possibility that an underlying
280 molecular mechanism drives these convergent patterns of polymorphism and divergence across
281 species.

282 We also tested whether the correlation between gene expression and variation was conserved
283 across species (Figure S7). To do this, we utilized chromatin-IP data across nine mouse tissues
284 and classified mouse tRNAs based on their breadth of expression. By this measure, active mouse
285 tRNAs were more strongly conserved than their inactive counterparts (Mann-Whitney U test, $p <$
286 $1.81e-19$), and their flanks were more divergent ($p < 7.04e-22$) (Figure S7A, S7D), consistent
287 with our results from the human data (Figure 2A,B). Active mouse tRNAs also had more low-
288 frequency SNPs in their flanking regions than did inactive mouse tRNAs ($p < 2.23e-4$) (Figures
289 S7C, S7F). Consistent with the human data, inactive mouse tRNAs were also more conserved (p
290 $< 1.76e-8$) and their flanking regions more divergent ($p < 2.37e-4$) than the untranscribed
291 reference regions (Figure S7D). Such consistency indicates that the mechanism underlying these
292 patterns works similarly in human and mouse.

293 The patterns of low-frequency SNPs are also consistent across all species. The greatest levels of
294 polymorphism are found in the inner 5' flanking regions for all species studied. The frequency

295 spectra of the low-frequency SNPs also show excess A→G, G→A, C→T and T→C SNPs on the
296 coding strand in all species analyzed (Figure S8) Additionally, as was observed in humans,
297 active mouse tRNAs show a greater excess of these SNPs (Figure S9A) than do inactive mouse
298 tRNAs (Figure S9B). Consistent with our analysis of human tRNAs, these patterns suggest that
299 deamination of the non-coding strand due to TAM and the DNA repair mechanisms acting in
300 response to deamination are especially common at these loci (4, 9, 36).

301 In humans, we do not have chromatin-IP data for germline tissues and cannot correct for the fact
302 that only mutations in these tissues are heritable, but we have no evidence that active tRNAs are
303 suppressed in the germline. However, the nine mouse tissues for which we had chromatin-IP data
304 included testes and mouse embryonic stem cells. Virtually all of the tRNAs that are inactive in
305 both stem cells and testes are also inactive in the other tissues. Because only germline mutations
306 are heritable, we expect that only germline expression causes elevated mutation rates at these
307 loci. That virtually all active tRNAs are expressed in the germline and that those not expressed in
308 the germline are inactive is consistent with our findings and suggests that estimates of tRNA
309 expression derived from somatic tissues are sufficient for studying the genetic consequences of
310 exceptional transcription rates.

311 **Functional tRNA sequences experience strong purifying selection in all species studied.** Our
312 analysis of the distribution of fitness effects (DFE) of deleterious mutations demonstrates that
313 tRNAs evolve under strong purifying selection in all of the species we analyzed. In contrast,
314 regions flanking tRNAs were inferred to be either neutral or subject to weak selection ($N_e S < 10$)
315 (Figure 4). These results are consistent with our estimates of evolutionary conservation in tRNA
316 regions, as well as elevated levels of polymorphism observed in the flanks (Figure 1). Our
317 estimates of the proportions of new mutations falling into each $N_e S$ range of the DFE for tRNAs
318 indicated that there were far fewer nearly neutral mutations ($N_e S < 1$) and substantially more
319 strongly deleterious mutations ($N_e S > 100$) in *D. melanogaster* and *A. thaliana* than in human or
320 mouse populations (Figure 4). Given that estimates of N_e in human (7, 000; 37) and mouse (25,
321 000 – 120, 000; 38) are substantially lower than in *A. thaliana* (300,000; 39) and especially *D.*
322 *melanogaster* (> 1,000,000; 40), this difference in the inferred strength of selection may be due
323 to differences in long term N_e . That the *A. thaliana* life cycle involves selfing may also
324 contribute to these differences. In addition, the human and mouse genomes contain far more
325 tRNAs (610 and 471, respectively) than *D. melanogaster* (295) (19), and this increased
326 redundancy could also affect the inferred fitness effects across these species, as mutations in high
327 copy-number tRNAs are potentially less deleterious than those affecting unique tRNAs.
328 However, given that there are 700 tRNA genes in the *A. thaliana* genome (19, 41), redundancy
329 alone is unlikely to fully account for the between-species differences in the DFE.

330 Several tRNAs are known to contain introns (19). We analyzed the introns separately and found
331 that intronic variation correlates with flanking variation in tRNAs; that is, tRNAs with the most
332 variable flanks also had the most variable introns (Figure S10). We considered using introns as
333 selectively neutral regions for estimating DFE, but found that these regions comprised only 619
334 nucleotides in total, fewer than the total number of human tRNAs (19). As such, this sample size
335 was too small to reliably use in our DFE analysis.

336 **tRNA loci contribute disproportionately to mutational load.** Our discovery of a highly

337 elevated mutation rate at tRNA loci suggests that tRNA genes may contribute disproportionately
338 to segregating mutational load in humans. To obtain an estimate of the contribution of tRNA loci
339 to this load, we used the ratio of the rate of low-frequency SNPs in tRNA flanking regions to that
340 in untranscribed reference regions (between 8.7 and 13.8) to estimate the tRNA mutation rate
341 relative to the neutral mutation rate in humans ($1.45e-8$, (42)). Given that there are 25,852 base
342 pairs of tRNA sequence for active tRNAs in the human genome, we estimate that the per
343 generation rate of deleterious mutation arising from tRNAs per diploid genome (U_{tRNA}) is 0.01.
344 Using previous estimates of the rate of deleterious amino acid mutation per diploid genome (0.35,
345 (43)), this implies that tRNAs may contribute 2.3% of deleterious mutations as protein coding
346 sequences. Given that tRNAs make up only 0.0009% of the human genome (19), this further
347 implies that mutations in tRNAs contribute to mutational load, the reduction in individual fitness
348 due to segregating deleterious mutations (44, 45), with an effect disproportionate to their total
349 lengths. Although such calculations are clearly approximate, they nevertheless highlight that
350 mutations at tRNA loci are likely an important source of fitness and disease variation in human
351 populations.

352 **Conclusions**

353 Our findings demonstrate the fundamental importance of tRNA sequences, which are highly
354 conserved despite the continual influx of mutations by TAM at a higher rate than anywhere else
355 in the genome. Our results are consistent across a broad range of taxonomically diverse species,
356 indicating that elevated mutation rates due to TAM and strong purifying selection are widespread
357 across life, and may be a good predictor of relative tRNA gene transcription levels. The conflict
358 between extreme TAM and strong purifying selection at tRNA loci is potentially an
359 unappreciated source of genetic disease, and may have a profound impact on the fitness of
360 human populations.

361 **Materials and Methods**

362 **Defining tRNA loci and flanking regions.** We used tRNA coordinates from GtRNAdb (19) for the human, *M.*
363 *musculus*, *D. melanogaster*, and *A. thaliana* genomes. For each species, we defined untranscribed reference
364 regions to use as negative controls. To find these regions in the human genome, we searched 10 kilobases
365 upstream of each tRNA and selected a 200-nucleotide tract. If this tract was within a highly transcribed region
366 of the genome (as determined by genome-wide chromatin-IP data (21)), overlapped a conserved element
367 (defined as a region with a phastCons log odds score greater than 0 (18)), was within 1,000 nucleotides of a
368 known gene (27), or overlapped an untranscribed reference region assigned to another tRNA, we selected a
369 new tract 1,000 bases further upstream, and repeated until we found an acceptable region. For the mouse
370 genome, we checked only known genes, previously assigned untranscribed reference regions, and conserved
371 elements, as analogous genome-wide chromatin-IP data of the caliber used for humans was not readily
372 available for other species. For the *D. melanogaster* and *A. thaliana* genomes, we began our searches only
373 1,000 bases upstream of each tRNA, and searched for 200-nucleotide tracts that were at least 100 nucleotides
374 away from any annotated genetic element (46, 47). This adjustment was made due to the relatively high
375 functional densities of the genomes of these species.

376 For each tRNA in all species, we defined the inner 5' flank as the 20 bases immediately upstream of the 5' end
377 of the tRNA on the coding strand, and the outer 5' flank as the 20 bases directly upstream of the inner 5' flank.
378 Likewise, the inner 3' flank refers to the 10 bases directly downstream of the tRNA on the coding strand, and
379 the outer 3' flank refers to the 30 bases downstream of these 10 bases. We made these decisions based on

380 inflection points in our data, as the flanking regions up to 20 bases upstream and 10 bases downstream of
381 tRNA genes seemed to have less variation. Further, while no studies to our knowledge report the length of
382 tRNA leader sequences in eukaryotes in general, we found that transcription usually ends about 10 bases
383 downstream of mature tRNA sequences (26, 48).

384 **Classifying tRNAs based on breadth of expression.** The Roadmap Epigenomics Consortium compiled
385 genome-wide epigenomic data across 127 human tissues and cell lines in order to characterize the state of
386 chromatin across the genome (21). Cozen, et al. (in preparation) analyzed the regions surrounding each tRNA
387 in each epigenome sample, and performed a clustering analysis to classify each genomic region according to
388 its most common epigenomic state. They then classified all human tRNAs based on the epigenomic state
389 annotation in the genome. In the corresponding model, regions in state 1 are near transcription start sites, and
390 regions in states 4 and 5 are not near transcription start sites but are nonetheless likely to be transcribed. tRNAs
391 in state 1 in at least 3% of tissues are referred to here as “active tRNAs”, and we consider the remaining
392 tRNAs to be “inactive”.

393 We followed a similar approach to classify mouse tRNAs. We used data from a 15-state Hidden Markov
394 Model based on chromatin-IP data in which states 5 and 7 corresponded to regions proximal to active
395 promoters (49). tRNAs in genomic regions annotated as state 5 or 7 in at least 3% of tissues were considered to
396 be "active", and all other tRNAs were considered "inactive". These classifications were not conducted in other
397 species due to lack of available data.

398 **Aligning tRNAs.** We aligned all tRNAs across all species using covariance models (41) and assigned
399 coordinates to each position in each tRNA and flank based on the Sprinzl numbering system (20). Using these
400 alignments, we created files assigning a Sprinzl coordinate to each genomic coordinate within tRNA sequences
401 or flanking regions for each species studied. For example, the first nucleotide at the 5' end of each tRNA was
402 assigned Sprinzl coordinate 1. To create Figures 1 and 2A and B, we averaged the phyloP, divergence and low-
403 frequency SNP data for all sites assigned to the same Sprinzl coordinate for their respective tRNA loci.
404 Because some tRNAs have insertions, deletions and variations in structure (e.g. Leucine tRNAs often have an
405 extended V-loop (19)), this alignment was necessary for position-wise comparisons between tRNAs.
406 Additionally, some low-scoring tRNAs did not align well using these methods, and Sprinzl coordinates could
407 not be properly assigned. We set a filter such that tRNAs with fewer than 50 aligned bases were excluded.

408 Some tRNAs are known to have extended leading or trailing sequences that are well conserved across species
409 and potentially contribute to the secondary structure of the tRNAs (16). We determined whether any conserved
410 elements (regions with a phastCons log odds score greater than 0 (18)) were present by using the Vertebrate
411 Multiz Alignment & Conservation track in the UCSC Genome Browser (27) for the regions 4-10 bases up or
412 downstream of each human tRNA. If a conserved element was present within this region, the tRNA was
413 excluded from our analyses, as these flanking regions might contribute to the secondary structure of mature
414 tRNAs, and would therefore be subject to higher levels of selection than the vast majority of tRNA flanking
415 regions.

416 We also excluded nuclear-encoded mitochondrial tRNAs from our analyses. These tRNAs are transferred from
417 mitochondrial genomes and therefore are not subject to the same evolutionary pressures as the vast majority of
418 tRNAs. Additionally, alignments of these tRNAs across species is dubious, as these transfers likely occurred
419 following speciation, and many of these genes are without true orthologs in other species. Therefore, excluding
420 these genes would better explain the patterns of mutation and selection affecting most tRNA genes.

421 **Parsing variation data.** We analyzed human variation data from the African superpopulation of humans,
422 consisting of 661 individuals from Kenya, Nigeria, Sierra Leone, The Gambia and Barbados, from Phase 3 of
423 the 1000 Genomes Project (50). For *D. melanogaster*, we acquired variation data for the Siavonga, Zambia

424 populations from the *Drosophila* Genome Nexus Database (46, 47). *Mus musculus castaneus* raw data were
425 obtained from Waterston, et al (51) and the *A. thaliana* data were obtained from the Arabidopsis Genome
426 Initiative (52). All non-human data were aligned and genotypes curated as described in Corbett-Detig et al (53).

427 Within each tRNA, flank, or untranscribed reference region, we considered positions with minor allele
428 frequencies greater than 0 but less than 0.05 to be low-frequency single nucleotide polymorphisms (SNPs).
429 This is based on the idea that SNPs with low minor allele frequencies are generally due to new mutations, on
430 which selection is less of a factor (54). Therefore, these are expected to more closely reflect the neutral
431 mutation rate and spectrum. We also determined the frequency of the 12 possible classes of mutations (e.g.
432 A→G, T→A) within each region of each tRNA where the identity of each base is defined according to the
433 coding strand sequence. Using the alignments, we found the frequency of divergences and low-frequency
434 SNPs by position across all tRNAs and flanking regions, and we obtained 95% confidence intervals for each
435 point estimate by non-parametric bootstrapping across tRNA loci.

436 For conservation studies across multiple species, we used the phyloP track (18) (across 100 vertebrate species
437 for the human data, across 60 vertebrate species for mouse), and across 27 insect species for the *D.*
438 *melanogaster* data) from the UCSC Genome Browser (27, 34) and calculated the average score for each
439 position within the tRNAs and flanking regions. The phyloP track assigns scores to each nucleotide in the
440 genome based on alignments to other species, where the score represents the $-\log$ p-values under a null model
441 of neutral evolution. Positive scores indicate strong conservation, negative scores indicate accelerated
442 evolution, and sites with scores of zero are undergoing change at a rate consistent with neutral genetic drift
443 (18). When plotting this data, we multiplied the average phyloP scores by negative one, such that sites
444 undergoing accelerated change would have high positive scores, and sites that were strongly conserved would
445 have negative scores (Figures 1A and D, 2A and B). We also performed non-parametric bootstrapping across
446 tRNA loci to determine 95% confidence intervals for all positions. No analogous genome-wide phyloP data
447 was available for *A. thaliana* (18).

448 For direct comparisons between the species of interest and an outgroup, we used the Multiz Alignment &
449 Conservation track from the UCSC Table Browser (34) and the Stitch MAFs tool from Galaxy (55) to create
450 sequence alignments of the regions of interest in the human and mouse genomes. For the human genome, we
451 downloaded the hg19 human reference genome from the UCSC Genome Browser and aligned to the *Macacca*
452 *mulatta* reference genome (rheMac2) (56), also from the UCSC Genome Browser (27). We also compared the
453 mouse (*Mus musculus*, mm10) and rat (*Rattus norvegicus*, rn6) genomes (34), and the *A. thaliana* (TAIR10)
454 and *A. lyrata* (v.1.0) genomes (57, 58) using the same methods. For *D. melanogaster*, we used an alignment of
455 the dm6 genome to the droYak2 (*D. yakuba*) genome (59). Non-gap nucleotide mismatches in the alignments
456 were classified as divergent sites. To account for the possibility that multiple substitutions occurred at a single
457 site, we applied a Jukes-Cantor correction to the average rate of divergence at each position (60).

458 **Transcription factor binding.** The ENCODE Project Consortium used ChIP-Seq data to identify binding
459 regions for regulatory factors (29–31), including the TATA-binding protein (TBP) and several Pol3
460 transcription factors in the human genome (10). These data were taken from the UCSC Genome Browser (27)
461 in the form of peak calls, in which the intensity of a given peak correlates with a greater frequency of
462 transcription factor binding to that region. For each human tRNA, we found the strongest TBP peak in the 50
463 base pairs immediately upstream of the tRNA, across the GM12878, H1-hESC, HeLa-S3, HepG2 and K562
464 cell lines. We chose to use TBP instead of Pol3 peaks because, although TBP is not specific to Pol3 genes, we
465 found that this data was a stronger and more reliable indicator of transcriptional activity. We also calculated
466 the average phyloP score across the flanking regions for each tRNA (18), and performed a Spearman's rank
467 correlation test to quantify the relationship between these data. We repeated this test for the maximum peaks
468 for BDP1 and RPC155 as well, but searched for peaks within the mature tRNA sequence instead, as this is
469 where these transcription factors bind (10). Further, we used only HeLa-S3 and K562 cell data for the BDP1

470 and RPC155 tests, as this was the only data available for these peaks (31). This ChIP-Seq data was available
471 only for the human genome, so other species were excluded from this part of our analysis.

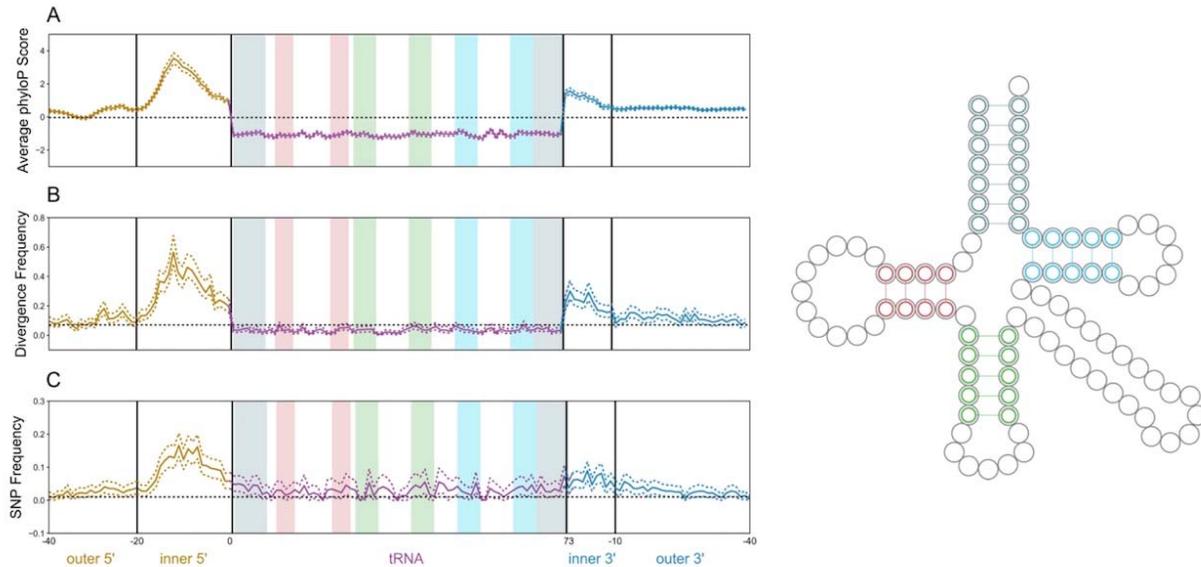
472 **Correlating variation to cell-line read counts.** Zheng, et al (32) developed a high-throughput demethylation
473 sequencing pipeline in order to efficiently detect tRNAs within human embryonic kidney (HEK293T) cells (32,
474 61). We performed Spearman's rank correlation tests to determine the relationship between their mature tRNA
475 read counts and tRNA gene and flanking region conservation. Because Zheng, et al (32) sequenced mature
476 tRNA sequences, which are sometimes encoded by multiple genes, and because we had variation data by gene,
477 we needed to account for this discrepancy. For example, tRNA-Val-CAC-1-1 and tRNA-Val-CAC-1-2 are two
478 distinct genes with different degrees of variation, but they encode the same mature tRNA. For tRNAs encoded
479 at only one locus, we evaluated the correlation between the read counts and the levels of variation at that locus.
480 However, for tRNAs encoded at multiple loci, we took two approaches. First, we excluded these genes entirely,
481 to eliminate the need to control for the correlation between gene copy-number and overall expression (Figure
482 2E, F) (22, 32). In a separate analysis, we summed the average phyloP scores at these loci, and evaluated the
483 correlation between these totals to the tRNA read counts (Figure S3).

484 **Finding genome-wide minimum phyloP scores.** We used the UCSC Table Browser (34) to determine which
485 sites in the human genome had phyloP scores of -20, the lowest possible phyloP score. These are the least
486 conserved sites in the human genome across the 100 vertebrate species compared in this track (18). Using
487 genomic coordinates of tRNAs from GtRNAdb (19), we determined what proportion of these sites overlapped
488 tRNA genes or flanking regions and performed hypergeometric tests to quantify associations between these
489 data.

490 **Estimating the distribution of fitness effects.** We estimated the distribution of fitness effects (DFE) for each
491 species by maximum likelihood using the method of Keightley et al (62), implemented in the DFE- α software.
492 The method is based on site frequency spectra (SFS) obtained from within-species SNP data, and assumes a
493 simple model of recent demographic change to correct the SFS at functional sites for possible skews caused by
494 demography. We used a two-epoch model of demographic change and estimated the DFEs for tRNAs, inner 3'
495 and inner 5' flanking regions for each species. Each of these classes of sites was assumed to be subject to
496 mutation, selection and drift, with gamma-distributed DFEs and an initial shape parameter (β) of 0.5. We also
497 estimated the DFE for sites that are likely to be evolving neutrally (outer 5' flanking regions), which were used
498 as the presumably untranscribed reference regions for generating the expected allele frequency distributions.
499 For each class of putatively selected sites, we analyzed folded site frequency spectra, and the fitness effects of
500 new deleterious mutations were estimated on a scale of $N_e S$, where N_e is a measure of the recent effective
501 population size and S is the strength of selection on a new mutation.

502 **ACKNOWLEDGMENTS.** We would like to thank Craig Mello for helpful discussions and input on this
503 project, Brian Lin for supplying the legend for Figure 1, as well as Andrew Holmes for sending the mouse
504 chromatin-IP data and for suggestions on this project. We also thank the rest of the Corbett and Lowe Labs for
505 their suggestions, which made our manuscript clearer and made our research more efficient. B.T. was funded
506 by a T-32 training grant (T32 HG008345) during this study.

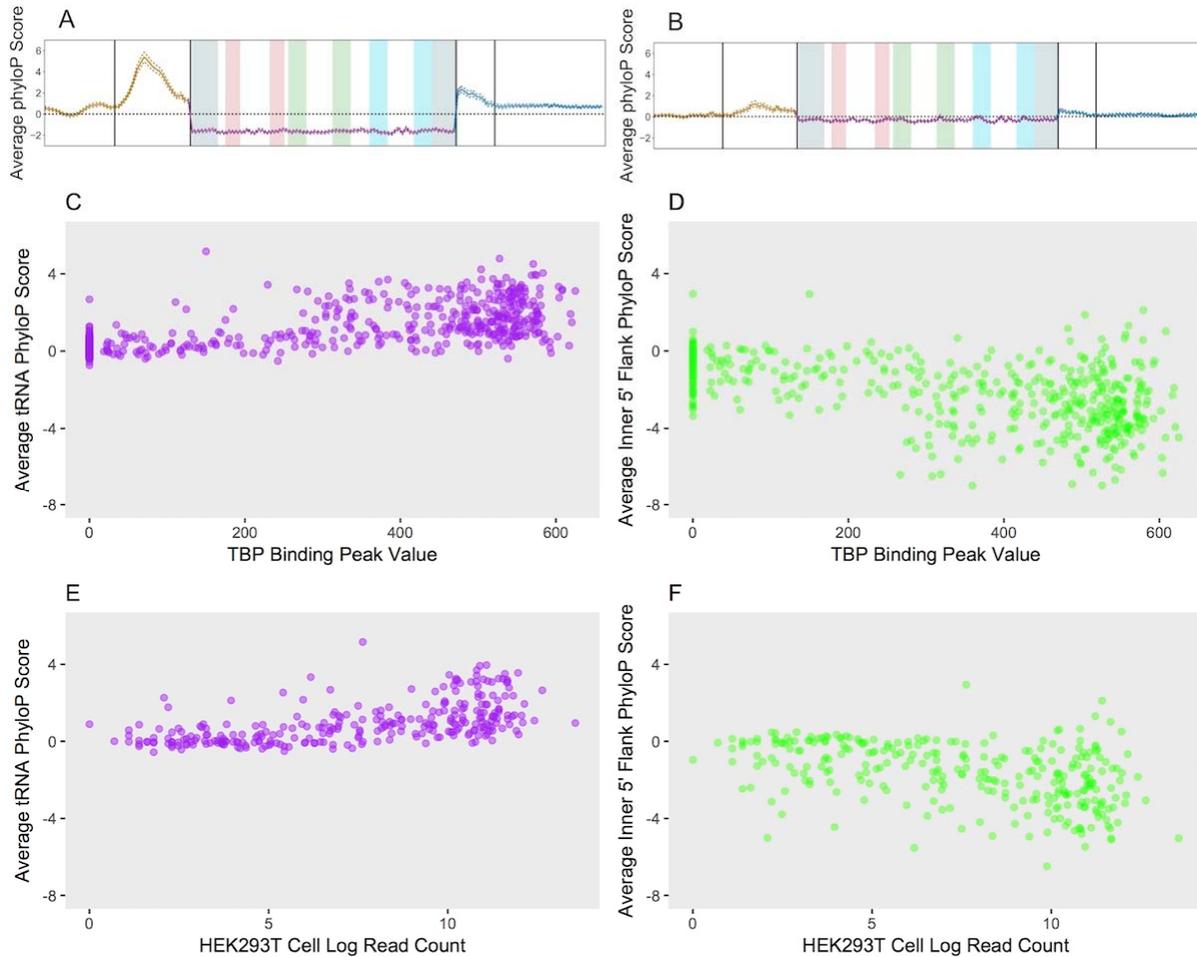
507



508

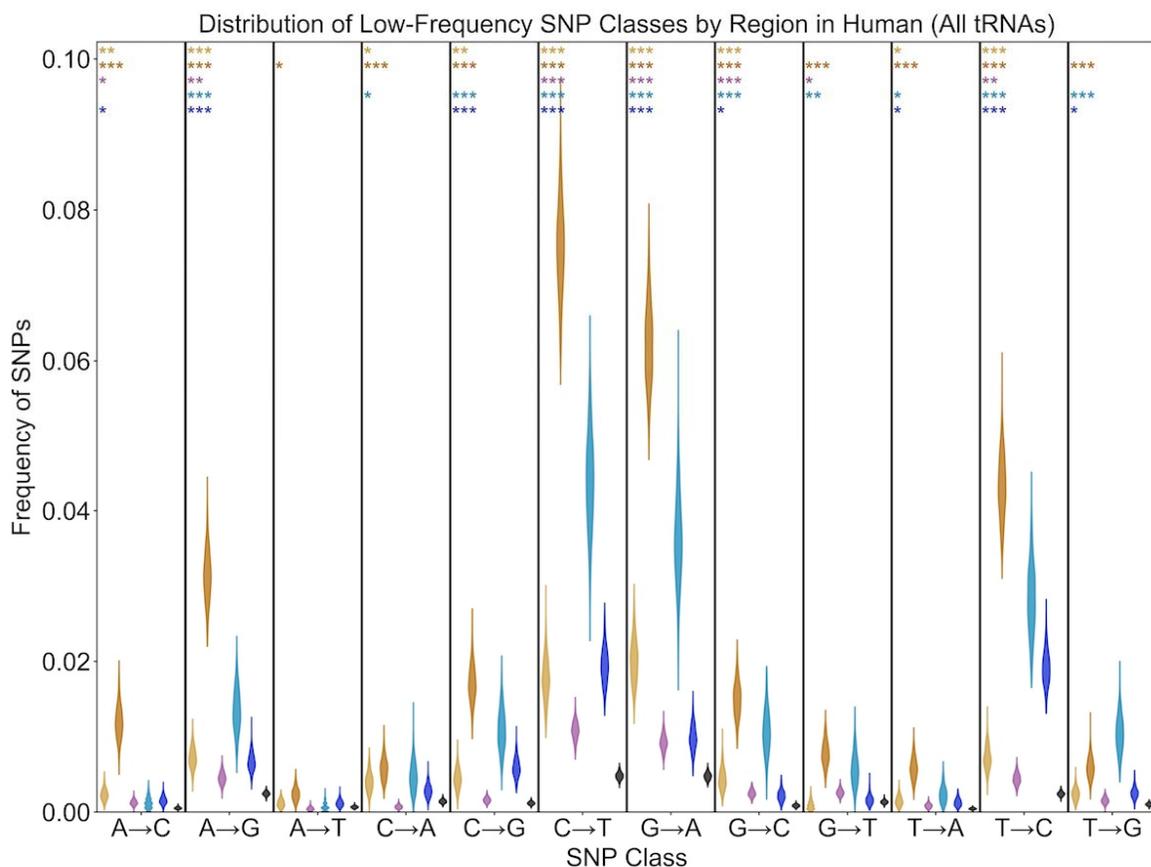
509 **Fig. 1. Strong pattern of variation in regions flanking human tRNA genes relative to vertebrates, upon comparison to**
510 **Rhesus macaque, and within the human population. A:** The negative of the average phyloP score (comparing humans to 100
511 vertebrate species) is plotted for each position within the tRNA and flank, across all human tRNAs. For consistency in plotting,
512 we multiplied the average score at each position by -1, so that more highly divergent regions would have higher, positive
513 scores. **B:** Divergence at non-gap alignments between the hg19 and rheMac2 genomes at each position within tRNAs and their
514 flanking regions. **C:** The frequency at which each position within tRNAs and flanks have a low-frequency SNP (minor allele
515 frequency less than or equal to 0.05) across all human tRNAs. The black dotted line in each plot represents the average value
516 across the untranscribed reference regions used in this study. The acceptor stem (gray), D-stem (red), C-stem (green) and V-stem
517 (blue) are highlighted within the tRNA, both in the plots and in the legend to the right, which shows the secondary structure of
518 the tRNA (19). The black vertical lines separate the inner and outer flanking regions. The 20 bases upstream and 10 bases
519 downstream of each tRNA are considered the inner 5' and inner 3' flanking regions, respectively, as these regions tended to show
520 a marked increase in variation relative to the outer flanking regions (see Methods). The dotted lines surrounding the plots depict
521 95% confidence intervals, calculated by bootstrapping by tRNA loci.

522



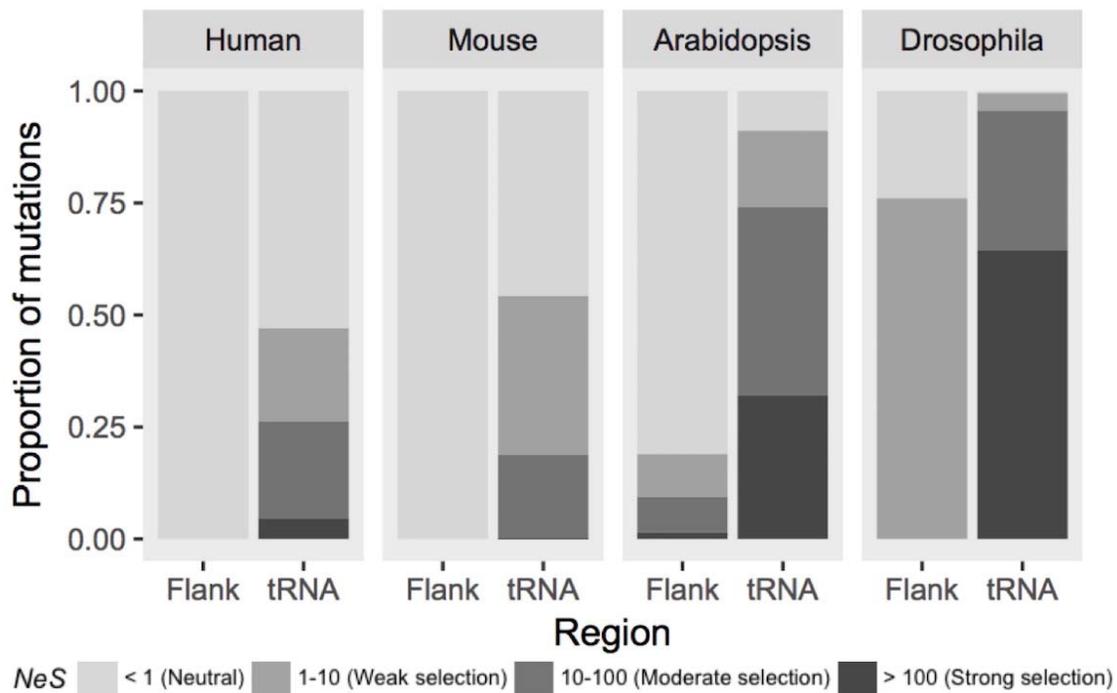
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538

Fig. 2. Three measures of tRNA expression are significantly correlated to both tRNA conservation and flanking region divergence. **A:** The negative of the average phyloP score (comparing humans to 100 vertebrate species) is plotted for each position within the tRNA and flank, across all active human tRNA loci. For consistency with Figure 1A, we used the negative of the phyloP score. **B:** The negative of the average phyloP score is plotted for each position within the tRNA and flank, across all inactive human tRNA loci, following the same format as **A**. **C:** Each tRNA's average phyloP score across its mature sequence is plotted against the value of the TBP peak corresponding to that tRNA. **D:** Each tRNA's average phyloP score across its inner 5' flanking region (20 nucleotides upstream of each tRNA gene) is plotted against the value of the TBP peak corresponding to that tRNA. **E:** The average phyloP score for each tRNA gene encoding a unique mature tRNA sequence is plotted against the log of the HEK293T cell read count for that tRNA (28). **F:** The average phyloP score across the inner 5' flanking region for each tRNA gene encoding a unique mature tRNA sequence is plotted against the log of the HEK293T cell read count for that tRNA (28). Several tRNAs are encoded by multiple sequentially identical genes. Because these would be expected to produce more tRNAs, and therefore have inflated read counts, we excluded these tRNAs from plots E and F. These tRNAs are included in Figure S3.



539
540
541
542
543
544
545
546
547

Fig. 3. SNP classes most common in regions affected by TAM are also most common at tRNA loci. The distribution of each class of low-frequency polymorphisms, here defined as a SNP with a minor allele frequency less than or equal to 0.05, is shown by region across all human tRNAs. At the top, the significance levels of Fisher's exact tests comparing the SNP distribution within each region of the tRNA and flank (outer 5' flank is yellow, inner 5' flank is orange, tRNA is purple, inner 3' flank is cyan, outer 3' flank is blue) to that of the untranscribed reference region (black) are represented by stars. One star represents a p value ≤ 0.05 , two stars represents a p value ≤ 0.005 , and three stars represents a p value ≤ 0.0005 .



548
549
550
551
552
553
554
555
556

Fig. 4. Estimated Distribution of Fitness Effects (DFE) indicates that tRNAs show high proportion of deleterious mutations are under strong selection. Estimated DFE of new deleterious mutations for tRNA genes and inner 3' flanking regions are shown in human, mouse, *A. thaliana* and *D. melanogaster*. The proportions of deleterious mutations are shown for each bin of purifying selection strength, estimated on a scale of NeS, where Ne is a measure of the recent effective population size and S is the strength of selection. The species are arranged by increasing Ne.

557 References

- 558 1. Higgs PG, Jameson D, Jow H, Rattray M (2003) The evolution of trna-leu genes in animal mitochondrial genomes. *Journal of Molecular Evolution* 57(4):435–445.
559
- 560 2. Kirchner S, Ignatova Z (2015) Emerging roles of tRNA in adaptive translation, signalling dynamics and disease. *Nat. Rev. Genet.* 16(2):98–
561 112.
- 562 3. Molla-Herman A, Valles AM, Ganem-Elbaz C, Antoniewski C, Huynh JR (2015) tRNA processing defects induce replication stress and Chk2-
563 dependent disruption of piRNA transcription. *EMBO J.* 34(24):3009–3027.
- 564 4. Jinks-Robertson S, Bhagwat AS (2014) Transcription-associated mutagenesis. *Annu. Rev. Genet.* 48:341–359.
- 565 5. Gnatt AL, Cramer P, Fu J, Bushnell DA, Kornberg RD (2001) Structural basis of transcription: an RNA polymerase II elongation complex at
566 3.3 A resolution. *Science* 292(5523):1876–1882.
- 567 6. Helmrich A, Ballarino M, Tora L (2011) Collisions between replication and transcription complexes cause common fragile site instability at
568 the longest human genes. *Mol. Cell* 44(6):966–977.
- 569 7. Timakov B, Liu X, Turgut I, Zhang P (2002) Timing and targeting of P-element local transposition in the male germline cells of *Drosophila*
570 *melanogaster*. *Genetics* 160(3):1011–1022.
- 571 8. Gomez-Gonzalez B, Aguilera A (2007) Activation-induced cytidine deaminase action is strongly stimulated by mutations of the THO complex.
572 *Proc. Natl. Acad. Sci. U.S.A.* 104(20):8409–8414.

- 573 9. Green P, Ewing B, Miller W, Thomas PJ, Green ED (2003) Transcription-associated mutational asymmetry in mammalian evolution. *Nat.*
574 *Genet.* 33(4):514–517.
- 575 10. White RJ (2011) Transcription by RNA polymerase III: more complex than we thought. *Nat. Rev. Genet.* 12(7):459–463.
- 576 11. Zhang J, Ferre-D'Amare AR (2016) The tRNA Elbow in Structure, Recognition and Evolution. *Life (Basel)* 6(1).
- 577 12. Sun C, et al. (2017) Roles of tRNA-derived fragments in human cancers. *Cancer Lett.* 414:16– 25.
- 578 13. Ziehler WA, Day JJ, Fierke CA, Engelke DR (2000) Effects of 5' leader and 3' trailer structures on pre-tRNA processing by nuclear RNase P.
579 *Biochemistry* 39(32):9909–9916.
- 580 14. Hopper AK (2013) Transfer RNA post-transcriptional processing, turnover, and subcellular dynamics in the yeast *Saccharomyces cerevisiae*.
581 *Genetics* 194(1):43–67.
- 582 15. Hasler D, et al. (2016) The Lupus Autoantigen La Prevents Mis-channeling of tRNA Fragments into the Human MicroRNA Pathway. *Mol.*
583 *Cell* 63(1):110–124.
- 584 16. Lee YS, et al. (2009) A novel class of small rnas: trna-derived rna fragments (trfs). *Cell* 23(23):2639–2649.
- 585 17. Maraia RJ, Lamichhane TN (2011) 3' processing of eukaryotic precursor tRNAs. *Wiley Interdiscip Rev RNA* 2(3):362–375.
- 586 18. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A (2010) Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.*
587 20(1):110–121.
- 588 19. Chan PP, Lowe TM (2009) GtRNAdb: a database of transfer RNA genes detected in genomic sequence. *Nucleic Acids Res.* 37(Database
589 issue):D93–97.
- 590 20. Sprinzl M, Horn C, Brown M, Ioudovitch A, Steinberg S (1998) Compilation of tRNA sequences and sequences of tRNA genes. *Nucleic*
591 *Acids Res.* 26(1):148–153.
- 592 21. Kundaje A, et al. (2015) Integrative analysis of 111 reference human epigenomes. *Nature* 518(7539):317–330.
- 593 22. Bloom-Ackermann Z, et al. (2014) A comprehensive tRNA deletion library unravels the genetic architecture of the tRNA pool. *PLoS Genet.*
594 10(1):e1004084.
- 595 23. Doran JL, Bingle WH, Roy KL (1988) Two human genes encoding tRNA(GCCGly). *Gene* 65(2):329–336.
- 596 24. Dieci G, Sentenac A (1996) Facilitated recycling pathway for RNA polymerase III. *Cell* 84(2):245–252.
- 597 25. Ciesla M, Boguta M (2008) Regulation of RNA polymerase III transcription by MafI protein. *Acta Biochim. Pol.* 55(2):215–225.
- 598 26. Orioli A, et al. (2011) Widespread occurrence of non-canonical transcription termination by human RNA polymerase III. *Nucleic Acids Res.*
599 39(13):5499–5512.
- 600 27. Kent WJ, et al. (2002) The human genome browser at UCSC. *Genome Res.* 12(6):996–1006.
- 601 28. Juo ZS, et al. (1996) How proteins recognize the TATA box. *J. Mol. Biol.* 261(2):239–254.
- 602 29. Dunham I, et al. (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature* 489(7414):57–74.
- 603 30. Myers RM, et al. (2011) A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biol.* 9(4):e1001046.
- 604 31. Kharchenko PV, Tolstorukov MY, Park PJ (2008) Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat. Biotechnol.*
605 26(12):1351–1359.
- 606 32. Zheng G, et al. (2015) Efficient and quantitative high-throughput tRNA sequencing. *Nat. Methods* 12(9):835–837.
- 607 33. Thul PJ, et al. (2017) A subcellular map of the human proteome. *Science* 356(6340).
- 608 34. Karolchik D, et al. (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* 32(Database issue):D493–496.

- 609 35. Schmidt S, et al. (2008) Hypermutable non-synonymous sites are under stronger negative selection. *PLoS Genet.* 4(11):e1000281.
- 610 36. Saini N, et al. (2017) APOBEC3B cytidine deaminase targets the non-transcribed strand of tRNA genes in yeast. *DNA Repair (Amst.)* 53:4–
611 14.
- 612 37. Tenesa A, et al. (2007) Recent human effective population size estimated from linkage disequilibrium. *Genome research* 17(4):520–526.
- 613 38. Phifer-Rixey M, et al. (2012) Adaptive evolution and effective population size in wild house mice. *Molecular biology and evolution*
614 29(10):2949–2955.
- 615 39. Cao J, et al. (2011) Whole-genome sequencing of multiple Arabidopsis thaliana populations. *Nat. Genet.* 43(10):956–963.
- 616 40. Shapiro JA, et al. (2007) Adaptive genic evolution in the *Drosophila* genomes. *Proceedings of the National Academy of Sciences*
617 104(7):2271–2276.
- 618 41. Lowe TM, Chan PP (2016) tRNAscan-SE On-line: integrating search and context for analysis of transfer RNA genes. *Nucleic Acids Res.*
619 44(W1):W54–57.
- 620 42. Narasimhan VM, et al. (2017) Estimating the human mutation rate from autozygous segments reveals population differences in human
621 mutational processes. *Nat Commun* 8(1):303.
- 622 43. Keightley PD (2012) Rates and fitness consequences of new mutations in humans. *Genetics* 190(2):295–304.
- 623 44. Haldane JBS (1937) The effect of variation of fitness. *The American Naturalist* 71(735):337–349.
- 624 45. Agrawal AF, Whitlock MC (2012) Mutation load: The fitness of individuals in populations where deleterious alleles are abundant. *Annual*
625 *Review of Ecology, Evolution, and Systematics* 43(1):115–135.
- 626 46. Lack JB, et al. (2015) The *Drosophila* genome nexus: a population genomic resource of 623 *Drosophila melanogaster* genomes, including 197
627 from a single ancestral range population. *Genetics* 199(4):1229–1241.
- 628 47. Lack JB, Lange JD, Tang AD, Corbett-Detig RB, Pool JE (2016) A Thousand Fly Genomes: An Expanded *Drosophila* Genome Nexus. *Mol.*
629 *Biol. Evol.* 33(12):3308–3313.
- 630 48. Fredrik Pettersson BM, Ardell DH, Kirsebom LA (2005) The length of the 5' leader of *Escherichia coli* tRNA precursors influences bacterial
631 growth. *J. Mol. Biol.* 351(1):9–15.
- 632 49. Bogu GK, et al. (2015) Chromatin and RNA Maps Reveal Regulatory Long Noncoding RNAs in Mouse. *Mol. Cell. Biol.* 36(5):809–819.
- 633 50. Consortium TGP (2015) A global reference for human genetic variation. *Nature* 526(7571):68–74.
- 634 51. Waterston RH, et al. (2002) Initial sequencing and comparative analysis of the mouse genome. *Nature* 420(6915):520–562.
- 635 52. authors listed N (2000) Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature* 408(6814):796–815.
- 636 53. Corbett-Detig RB, Hartl DL, Sackton TB (2015) Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol.*
637 13(4):e1002112.
- 638 54. Goddard KA, Hopkins PJ, Hall JM, Witte JS (2000) Linkage disequilibrium and allele-frequency distributions for 114 single-nucleotide
639 polymorphisms in five populations. *Am. J. Hum. Genet.* 66(1):216–234.
- 640 55. Afgan E, et al. (2016) The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids*
641 *Res.* 44(W1):W3–W10.
- 642 56. Gibbs RA, et al. (2007) Evolutionary and biomedical insights from the rhesus macaque genome. *Science* 316(5822):222–234.
- 643 57. Kersey PJ, et al. (2016) Ensembl genomes 2016: more genomes, more complexity. *Nucleic Acids Research* 44(D1):D574–D580.
- 644 58. Rosenbloom KR, et al. (2015) The UCSC genome browser database: 2015 update. *Nucleic Acids Research* 43(D1):D670–D681.
- 645 59. Clark AG, et al. (2007) Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* 450(7167):203–218.

- 646 60. Jukes TH, Cantor CR (1969) Evolution of protein molecules. *Mammalian Protein Metabolism* pp. 21–132.
- 647 61. Cozen AE, et al. (2015) ARM-seq: AlkB-facilitated RNA methylation sequencing reveals a complex landscape of modified tRNA fragments.
648 *Nat. Methods* 12(9):879–884.
- 649 62. Keightley PD, Eyre-Walker A (2007) Joint inference of the distribution of fitness effects of deleterious mutations and population demography
650 based on nucleotide polymorphism frequencies. *Genetics* 177(4):2251–2261.
- 651
- 652
- 653
- 654
- 655
- 656