



# Transfer RNA genes experience exceptionally elevated mutation rates

Bryan P. Thornlow<sup>a</sup>, Josh Hough<sup>a</sup>, Jacquelyn M. Roger<sup>a</sup>, Henry Gong<sup>a</sup>, Todd M. Lowe<sup>a,b,1</sup>, and Russell B. Corbett-Detig<sup>a,b,1</sup>

<sup>a</sup>Department of Biomolecular Engineering, University of California, Santa Cruz, CA 95064; and <sup>b</sup>Genomics Institute, University of California, Santa Cruz, CA 95064

Edited by Andrew G. Clark, Cornell University, Ithaca, NY, and approved July 30, 2018 (received for review February 5, 2018)

**Transfer RNAs (tRNAs) are a central component for the biological synthesis of proteins, and they are among the most highly conserved and frequently transcribed genes in all living things. Despite their clear significance for fundamental cellular processes, the forces governing tRNA evolution are poorly understood. We present evidence that transcription-associated mutagenesis and strong purifying selection are key determinants of patterns of sequence variation within and surrounding tRNA genes in humans and diverse model organisms. Remarkably, the mutation rate at broadly expressed cytosolic tRNA loci is likely between 7 and 10 times greater than the nuclear genome average. Furthermore, evolutionary analyses provide strong evidence that tRNA genes, but not their flanking sequences, experience strong purifying selection acting against this elevated mutation rate. We also find a strong correlation between tRNA expression levels and the mutation rates in their immediate flanking regions, suggesting a simple method for estimating individual tRNA gene activity. Collectively, this study illuminates the extreme competing forces in tRNA gene evolution and indicates that mutations at tRNA loci contribute disproportionately to mutational load and have unexplored fitness consequences in human populations.**

tRNA | transcription | mutagenesis | TAM | computational prediction

**T**ransfer RNAs (tRNAs) are essential to protein synthesis across all of life. Their primary function is in translation of the genetic code into the corresponding amino acid sequences that make up proteins. Thus, tRNA molecules are critical for virtually all cellular processes, and the genes encoding tRNA molecules have been highly conserved over evolutionary time (1, 2). Mitochondrial tRNAs have been the subject of many studies, as mutations in these genes lead to a large number of maternally inherited genetic diseases (3). However, eukaryotic genomes contain ~10- to 20-fold as many tRNA genes encoded in their nuclear chromosomes, which are required for cytosolic protein translation (2, 4). Despite their importance to the cell, there has been little study of evolutionary conservation or pathogenic mutations in cytosolic tRNA genes (5, 6). tRNAs are required in exceptionally large quantities, and therefore tRNA genes may experience greater levels of transcription than even the most highly transcribed protein-coding genes (7, 8). In turn, this may lead to high levels of transcription-associated mutagenesis (TAM). As the largest, most ubiquitous RNA gene family, cytosolic tRNAs constitute an ideal gene set for studying the interplay between natural selection and elevated mutation rates.

Transcription affects the mutation rates of transcribed genes (9) through the unwinding and separation of cDNA strands (10). During transcription, a nascent RNA strand forms a hybrid DNA–RNA complex with a template DNA strand. While the complementary tract of nontemplate DNA is temporarily isolated, it is chemically reactive and thus accessible by potential mutagens (10). Transcription can lead to the formation of noncanonical DNA structures, which can hinder repair pathways and promote errors by the polymerase (11). The RNA strand can also reanneal to the template DNA strand, prolonging isolation and increasing vulnerability to mutations (12, 13). Furthermore, if transcription and DNA replication occur concomitantly at a particular locus, collisions between

RNA polymerase and the DNA replication fork may also damage DNA (9, 11, 14). In human cancer cells, increased transcription and replication induce torsional stress and collisions (11).

Several cellular agents have also been shown to cause damage in highly expressed genes (15). Among the most notable sources of mutation associated with transcription is activation-induced cytidine deaminase (AID) (16). AID accompanies RNA polymerase II and deaminates cytosine nucleotides. To resolve the resulting base-pair mismatch, the opposing guanine is converted to adenine and uracil to thymine, resulting in excess C→T mutations on the nontemplate strand and excess G→A mutations on the template strand (9, 17). AID is a member of the APOBEC (apolipoprotein B mRNA editing catalytic polypeptide-like) gene family, many of which are involved in double-stranded break repair in transcription (9). Some members of the APOBEC family act strongly at short genes, suggesting increased activity at tRNA loci (18, 19). For example, APOBEC3B causes 1,000-fold more DNA damage at tRNA loci than at other genomic regions in yeast (19). AID also acts on highly transcribed genes in immune B cells, causing transition mutations and double-stranded breaks (9). Due to the strong association of the APOBEC family with transcription, relative excesses of C→T and G→A mutations are a signature of TAM (9).

To conserve mature tRNA sequence identity in the presence of an elevated mutation rate, tRNA genes should experience strong purifying selection. tRNA transcription requires sequence-specific binding of transcription factors to the internal box A and box B

## Significance

**While transcription-associated mutagenesis (TAM) has been demonstrated for protein-coding genes, its implications in shaping genome structure at transfer RNA (tRNA) loci in metazoans have not been fully appreciated. We show that cytosolic tRNAs are a striking example of TAM because of their variable rates of transcription, well-defined boundaries, and internal promoter sequences. tRNA loci have a mutation rate approximately 7- to 10-fold greater than the genome-wide average, and these mutations are consistent with signatures of TAM. These observations indicate that tRNA loci are disproportionately large contributors to mutational load in the human genome. Furthermore, the correlations between tRNA locus variation and transcription indicate that prediction of tRNA gene expression based on sequence variation data is possible.**

Author contributions: B.P.T., T.M.L. and R.B.C.-D. conceived the project; B.P.T., T.M.L., and R.B.C.-D. designed research; B.P.T. and J.M.R. performed research; B.P.T., J.H., J.M.R., and H.G. analyzed data; and B.P.T., T.M.L., and R.B.C.-D. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

This open access article is distributed under [Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 \(CC BY-NC-ND\)](https://creativecommons.org/licenses/by-nc-nd/4.0/).

<sup>1</sup>To whom correspondence may be addressed. Email: tmjlowe@ucsc.edu or rucorbet@ucsc.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801240115/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801240115/-DCSupplemental).

Published online August 20, 2018.

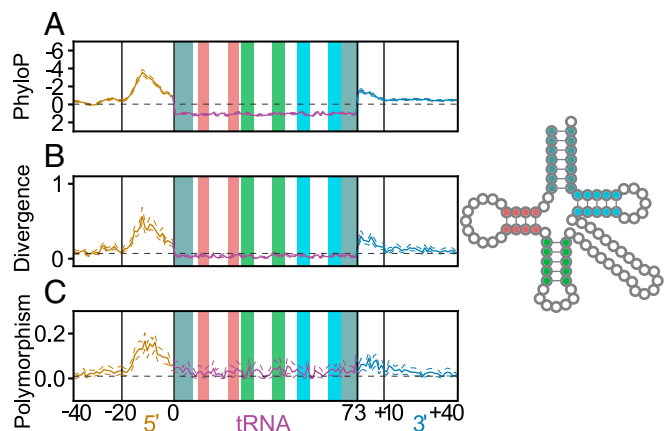
promoter elements (20). Once transcribed, precursor tRNAs must fold properly to undergo maturation, which can be disrupted by sequence-altering mutations. The unique structure of tRNAs dictates processing by RNases, addition of modifications, accurate recognition by aminoacyl tRNA synthetases, incorporation into the translating ribosome, and accurate positioning of the anticodon relative to mRNA codons (21, 22). Because of the need to maintain sequence specificity, DNAs encoding the mature portions of tRNAs are well conserved (21). Therefore, we expect that a large proportion of mutations arising in tRNA genes will be deleterious and will quickly be purged by natural selection.

While most human tRNA genes do not have external promoters (20, 21), tRNA transcripts include leader and trailer sequences extending roughly two to five nucleotides upstream and 5–15 nucleotides downstream of the annotated mature tRNA gene, based on the position of the genomically encoded poly(T) transcription termination sequence. Aside from the termination sequence, these flanking sequences appear to have limited sequence-specific functionality in most cases (23–26). Very early in maturation, all tRNA flanking sequences are removed by RNase P (22–24) and RNase Z (22, 27). Because these flanking genomic sequences are frequently unwound and therefore vulnerable to TAM, we expect that these regions will experience mutation rates similar to those of tRNAs. Whereas tRNA genes should experience purifying selection, the flanking regions should be neutral or under weak selection. Here we investigate the patterns of conservation, divergence, and within-species variation of cytosolic tRNAs in humans and other model organisms to elucidate the forces shaping the evolution of this essential RNA gene family.

## Results and Discussion

**Flanking Regions of tRNA Genes Are Highly Variable Despite Strong Conservation of Mature tRNA Sequences.** To estimate evolutionary conservation, we examined PhyloP, which measures the conservation of each human genomic position across 100 vertebrate species (28), by position within each tRNA locus (*Methods*). Positive PhyloP scores indicate strong conservation, and negative scores indicate accelerated evolution. To study the effects of evolution on a shorter timescale, we also estimated sequence divergence between human and *Macaca mulatta* at each tRNA locus. Mature tRNA sequences are highly conserved across all positions, based on both average PhyloP score (Fig. 1*A* and *Dataset S1*) (28) and *M. mulatta* alignment (Fig. 1*B*). However, the inner 5' flanking region (20 bases upstream of the tRNA; see *Methods*) is roughly four times more divergent than the untranscribed reference regions. We also find increased rates of divergence in the inner 3' flanking region, which is roughly three times more divergent than the reference regions (Fig. 1*B*). Both the outer 5' flank (21–40 bases upstream of the tRNA) and the outer 3' flank (11–40 bases downstream of the tRNA) are also roughly 1.5 times more divergent than the reference regions. For tRNAs that contain introns (2), we find that intronic variation correlates with flanking variation (*SI Appendix*, Fig. S1). Furthermore, intergenic regions within clusters of active tRNAs show similar patterns in their PhyloP scores (*SI Appendix*, Fig. S2).

We also studied population-level variation at low-frequency SNPs (minor allele frequency <0.05) for each tRNA locus. Low-frequency SNPs are evolutionarily young and are less affected by selection (29). Consistent with our divergence analyses, we find that low-frequency SNPs are more common across both the tRNA gene sequence and flanking regions than in untranscribed reference regions (Fig. 1*C*). Although the inner flanking regions are most polymorphic, the mature tRNA sequences have about twice as many low-frequency SNPs as reference regions. Overall, our results are consistent on multiple timescales, indicating that tRNAs and flanking sequences are prone to mutation. Indeed, of the 247 sites in the genome that have the lowest possible PhyloP scores, –20 (28, 30), 14 are 10–15 bases upstream of the start of an active tRNA gene,



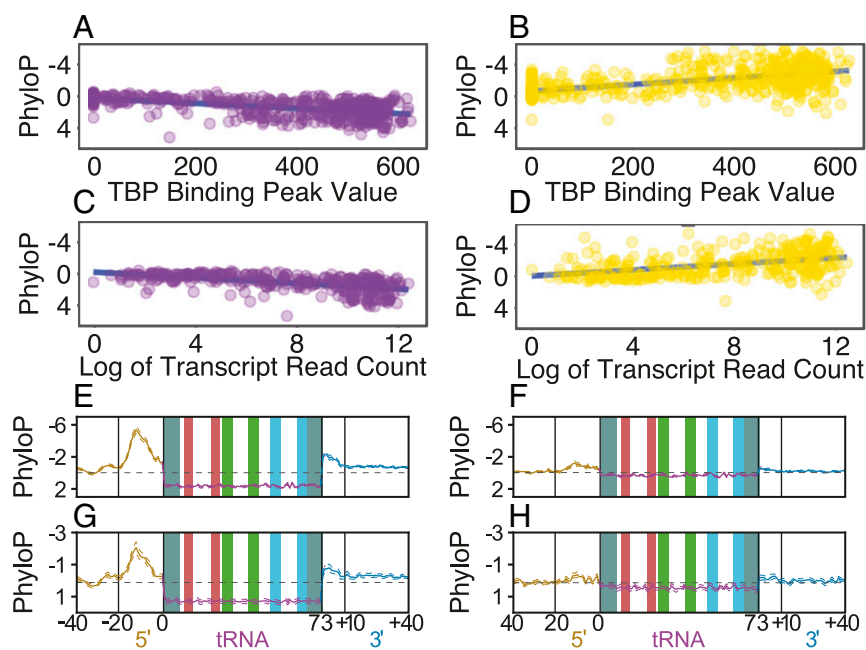
**Fig. 1.** There is a strong pattern of variation in regions flanking human tRNA genes by three measures: relative to vertebrates, by comparison with Rhesus macaque alone, and within the human population. (*A*) The average PhyloP score (comparing humans to 100 vertebrate species) is plotted for each position within the tRNA and flanking region across all human tRNAs. (*B*) Divergence between the human and *M. mulatta* tRNA genes and their flanking regions. (*C*) Frequency of low-frequency SNPs (minor allele frequency  $\leq 0.05$ ) across all human tRNAs. The acceptor stem (gray), D-stem (red), anticodon stem (green), and T-stem (blue) are highlighted within the tRNA both in the linear plots and in the 2D structure legend to the right (2, 61). Nucleotide numbering below the plots is relative to mature tRNA boundaries, with inner and outer flanks demarcated by a shift in mutation rate (*Methods*). Dotted lines surrounding plots depict 95% CIs calculated by nonparametric bootstrapping by tRNA loci.

indicating disproportionate enrichment (hypergeometric test,  $P < 1.65e-48$ ) and that tRNA flanking regions are among the least conserved in the genome. Nonetheless, mature tRNA gene sequences are strongly conserved by purifying selection, which purges mutations.

**Transcription Is Correlated with Variation in tRNA and Flanking Regions.** We hypothesized that, if transcription-associated mutagenesis drives variation among tRNA loci, highly active tRNA genes would show the greatest mutation rates. Because tRNA transcript abundance measures are often not attributable to individual loci due to identical gene copies and difficulty sequencing full-length tRNAs, we estimated relative transcriptional activity based on chromatin state data from the Epigenomic Roadmap Project (31). Based on these data, we classified human tRNA genes as “active” if they are located in expressed regions in several cell lines and otherwise as “inactive” (*Methods* and Fig. 2). In some cases, multiple cell lines correspond to a single tissue or organ, so tissue-specific tRNAs [e.g., the brain-specific arginine tRNA in mouse (6)] are considered active.

We find that active tRNA genes are significantly more conserved than inactive tRNA loci (Mann–Whitney  $U$  test,  $P < 8.40e-53$ ), and the flanking regions of active tRNAs are significantly more divergent than the flanking regions of inactive tRNAs ( $P < 7.98e-61$ ). The peak measure of divergence between human and *M. mulatta* tRNA genes in the inner 5' flanking regions is roughly five times greater in active tRNAs than in inactive tRNAs (Fig. 2*E* and *F*). Active tRNAs in human populations also have significantly more low-frequency SNPs per site than inactive tRNAs across the entire locus, including the tRNA and flanking regions ( $P < 3.72e-36$ ) (*SI Appendix*, Fig. S3). Inactive tRNAs are still significantly more conserved ( $P < 2.02e-12$ ) and polymorphic ( $P < 0.007$ ) than the untranscribed reference regions, and their flanks are significantly more divergent than the reference regions ( $P < 1.36e-16$ ).

That the peak in both divergence and polymorphism in all species is consistently 12–15 nucleotides upstream of the mature tRNA sequence is curious. At the most divergent position, 55% of all tRNA loci differ between human and *M. mulatta*, and 15% of human tRNA loci have a low-frequency SNP (Fig. 1). Furthermore,



**Fig. 2.** tRNA expression is significantly correlated to both tRNA conservation and flanking region divergence. (A and B) TBP peak value (expression) is plotted versus PhyloP score (conservation) for each mature tRNA (A) and adjacent inner 5' flanking region (B). (C and D) Log of the HEK293T cell DM-tRNA-seq read count (expression) (46) is plotted versus PhyloP score (conservation) for each gene encoding a unique mature tRNA sequence (C) and the corresponding inner 5' flanking region (D). Both TBP occupancy and transcript abundance are greater for highly conserved mature tRNA loci (A and C) and those with the most divergent flanks (B and D). (E and F) Plotted as in Fig. 1, human tRNA loci that are separated into active (E) versus inactive (F) groups show the characteristic differences seen in A–D. (G and H) Mouse tRNA loci split into active (G) versus inactive (H) groups show a pattern strikingly similar to that seen in human (A–F).

virtually all active tRNA loci differ at this nucleotide between human and *M. mulatta*, and 25% have a low-frequency SNP at this site (*SI Appendix, Fig. S3B*). This implies that this region either does not face uniform selective pressures or is not uniformly vulnerable to TAM. While distant flanking sequences can affect tRNA expression in yeast (32), few studies have shown that flanking regions affect expression in higher eukaryotes (33). Transcription initiation is long relative to elongation (34, 35), which may lead to prolonged isolation of the nontemplate DNA strand at the initiation site and increased vulnerability to TAM. A poised initiation complex might also increase the likelihood of collisions between Pol3 and the replication fork (14). Thus, frequent initiation at highly transcribed tRNA loci may contribute to the nonuniform pattern of variation.

This may also explain the increased variation in the outer 3' flank relative to the outer 5' flank, as positioning of downstream transcription termination sites varies among tRNA genes (2, 36), whereas transcription start site positions are more consistent. While most tRNAs do not have clear TATA boxes, the TATA-binding protein (TBP) still binds to the DNA duplex ~25 nucleotides upstream of the tRNA (37), which coincides with a decrease in variability. Furthermore, while both flanking regions for many other Pol3-transcribed genes are divergent, the 5' flanking regions are generally more divergent than the 3' flanking regions, suggesting that the underlying mechanism is not tRNA-specific (*Dataset S1*). However, additional studies are necessary to support the assertion that this pattern is due to transcription.

Two orthogonal analyses strengthen the observed correlations between gene expression and variation at tRNA loci. First, we find a significant correlation between the TBP intensity peaks (38–40) and conservation of the mature tRNA sequence (Spearman's  $\rho = 0.64$ ,  $P < 2.2e-16$ ) across all human tRNAs and the opposite relationship in the flanking regions (Spearman's  $\rho = -0.64$ ,  $P < 2.2e-16$ ) (Fig. 2). TBP ChIP-sequencing (ChIP-seq) data directly reflect transcriptional activity for each locus, as its occupancy is significantly correlated with and required for transcription (20, 41–45). Second, mature tRNA sequence read counts are strongly correlated with tRNA conservation (Spearman's  $\rho = 0.18$ ,  $P < 0.001$ ) and flanking region divergence (Spearman's  $\rho = -0.61$ ,  $P < 2.2e-16$ ) (Fig. 2 and *SI Appendix, Fig. S4*). These read counts were collected from a single HEK cell line by Zheng et al. (46)

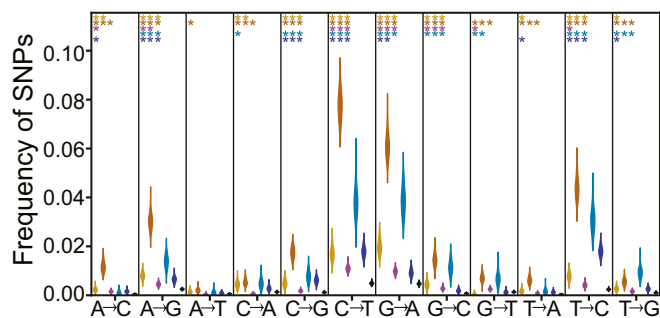
using DM-tRNA-seq, a specialized tRNA-sequencing method. These correlations are consistent with the idea that more highly transcribed tRNAs vary more in their flanking regions.

**Patterns of Divergence and Conservation Can Be Leveraged to Predict tRNA Gene Expression.** Regardless of how tRNA expression is measured, we find highly significant correlations between gene expression and tRNA sequence conservation. The consistency of these correlations across methods of measurement and across species indicates that it may be possible to predict relative tRNA with DNA sequence conservation patterns and other correlates of tRNA transcriptional activity (e.g., tRNAscan-SE bit scores). Indeed, active and inactive tRNAs are largely distinguishable using only flank and gene PhyloP data (*SI Appendix, Fig. S5*). As sequencing technology becomes more accessible, predicting tRNA gene-expression levels through analysis of DNA data is enticing. Such a model could make future tRNA gene annotation more detailed and cost-effective.

**Variation Patterns Observed at tRNAs Are Not Observed in Most Other Gene Families.** Applicability of this proposed tool is likely best suited for tRNAs, other Pol3 genes, and unique classes of highly expressed protein-coding genes such as histones. Among the histone protein-coding genes less than 1,000 nucleotides in length, the average PhyloP score per nucleotide across the coding sequence and flanking regions is 3.449 and  $-2.052$ , respectively, comparable to tRNA loci (*SI Appendix, Fig. S6*). In contrast, most genes transcribed by RNA Pol2 do not appear to be good targets (*Dataset S1*). For example, ribosomal proteins are very highly transcribed (47) and have well-conserved exons, but their introns and flanking regions are not as divergent as tRNA flanking regions (28, 48). tRNAs are likely ideal for studying TAM because they have predictable transcript start and end sites, internal promoters, and high transcription rates.

**Patterns of Low-Frequency SNPs Are Consistent with TAM.** In TAM, repair pathways activated in response to deaminations lead to excess conversions between guanine and adenine and between thymine and cytosine on the coding strand (9, 17). Across all tRNA loci, we found that the most common low-frequency SNPs are C→T and G→A and that these mutations are significantly more common in both tRNA genes and flanking regions than in untranscribed reference regions (Fisher's exact test,  $P < 0.05$  for all comparisons) (Fig. 3). Removal





**Fig. 3.** The SNP classes most common in regions affected by TAM are also most common at tRNA loci. Shown is the distribution of each class of low-frequency polymorphisms by region across all human tRNAs. Stars indicate the significance levels of Fisher's exact tests comparing the SNP distribution within each region of the tRNA and flank (outer 5' flank in yellow, inner 5' flank in orange, tRNA in purple, inner 3' flank in cyan, outer 3' flank in blue) with that of the untranscribed reference region (black): one star,  $P \leq 0.05$ ; two stars,  $P \leq 0.005$ ; and three stars,  $P \leq 0.0005$ .

of CpG sites (49) does not significantly affect these results. The relative excesses of these SNPs are much more pronounced in active tRNA loci than in inactive tRNA loci (*SI Appendix, Fig. S7 A and B*). These results suggest that deamination of the noncoding strand due to TAM and the DNA repair mechanisms acting in response to deamination is especially common at these loci (9, 17, 19).

It is difficult to discern whether this increased prevalence is due to TAM or selection to preserve the structural integrity of the tRNA. To preserve tRNA secondary structure, we expect transition mutations (e.g., A–U to G–U base pairs, C–G to U–G base pairs) to be more common than transversions, as they should disrupt stem helices less often. However, the mutational skew expected of regions affected by TAM is stronger in regions flanking tRNAs. Transcription initiation is relatively long compared with elongation (34, 35), which might contribute to increased mutagenesis by APOBEC enzymes or more collisions (14) or double-stranded breaks. However, divergence at tRNA flanking regions is correlated with divergence at introns in both human (Spearman's rank,  $\rho = 0.734$ ,  $P < 5.58e-6$ ) (*SI Appendix, Fig. S1B*) and mouse ( $\rho = 0.733$ ,  $P < 5.24e-4$ ) (*SI Appendix, Fig. S1D*), indicating similar mutation rates across tRNA loci. Our results therefore suggest that TAM drives the excess of transitions among low-frequency SNPs across tRNA loci.

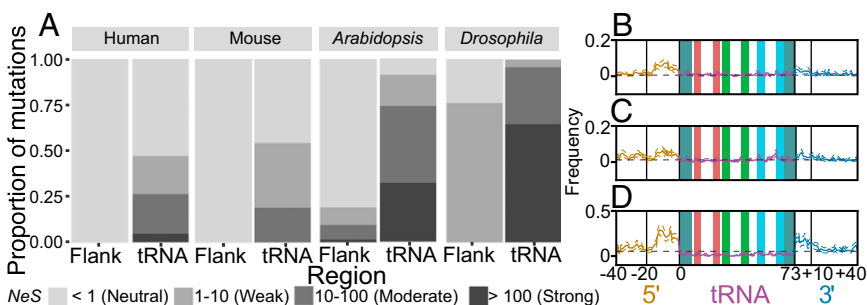
**tRNA Flanking Region Variation in Other Model Organisms Is Consistent with Variation Observed in Humans.** To confirm that our results are not restricted to humans, we also analyzed tRNAs in *Mus musculus*, *Drosophila melanogaster*, and *Arabidopsis thaliana*. We find similar patterns of sequence conservation of tRNA loci in each when measuring PhyloP or divergence to outgroups (*SI Appendix, Fig. S8*). The 5' flanks are consistently more divergent than the 3' flanks, and the most divergent sites are roughly 10–15 bases upstream of

the tRNA in all species. We also used ChIP data across nine mouse tissues to classify mouse tRNAs based on their expression (50). Active mouse tRNAs are more strongly conserved than their inactive counterparts (Mann–Whitney  $U$  test,  $P < 1.81e-19$ ), and their flanks are more divergent ( $P < 7.04e-22$ ) (Fig. 2 *G* and *H*), consistent with our results from the human data (Fig. 2 *E* and *F*). Active mouse tRNAs also have more low-frequency SNPs in their flanking regions than inactive mouse tRNAs ( $P < 2.23e-4$ ) (*SI Appendix, Fig. S9*). Such consistency suggests that a shared underlying molecular mechanism drives these patterns of sequence variation.

Low-frequency SNPs in the tRNA gene sequences also follow qualitative patterns similar to those in the human data. We observe excess transitions in all species studied (*SI Appendix, Fig. S10*), and active mouse tRNAs show a greater excess of low-frequency transitions than do inactive mouse tRNAs (*SI Appendix, Fig. S7 C and D*). However, these patterns vary across species (Fig. 4 *B–D*). For example, in mouse, tRNA genes have more low-frequency SNPs than the untranscribed reference regions (Fig. 4*B*), but the opposite is true in *D. melanogaster* (Fig. 4*D*). Low-frequency SNPs are thought not to be strongly affected by selection (29), but selection is more efficient in species with greater effective population sizes (Fig. 4*A*). Effective population size (51–54) and tRNA copy number vary across species, and because the sample sizes and data quality differ among population samples, these differences may be attributable to differences in the impact of selection or in ascertainment of low-frequency variation.

**Functional tRNA Sequences Experience Strong Purifying Selection in All Species Studied.** Our analysis of the distribution of fitness effects (DFE) of deleterious mutations demonstrates that tRNAs evolve under strong purifying selection in all analyzed species. In contrast, regions flanking tRNAs are inferred to be either neutral or subject to weak selection ( $N_e S < 10$ , where  $N_e$  is the effective population size and  $S$  is the strength of selection) (Fig. 4*A*). Our estimates of the proportions of new mutations falling into each  $N_e S$  range of the DFE for tRNAs indicate far fewer nearly neutral mutations ( $N_e S < 1$ ) and substantially more strongly deleterious mutations ( $N_e S > 100$ ) in *D. melanogaster* and *A. thaliana* than in the human or mouse populations (Fig. 4*A*). Given that estimates of effective population size in humans (7,000) (51) and mouse (25,000–120,000) (52) are substantially lower than in *A. thaliana* (300,000) (53) and *D. melanogaster* (>1,000,000) (54), this difference in strength of selection may partially reflect differences in effective population size and might explain the differences in low-frequency SNPs in tRNA loci across species (Fig. 4 *B–D*). In turn, this might indicate that the strength of purifying selection, independent of effective population size, at tRNA loci is consistent across diverse species.

The strength of selection across species may also reflect the number of unique tRNA gene sequences in each genome. For example, roughly half of all human tRNA genes have unique sequences, but the majority of *D. melanogaster* tRNAs have identical copies (2). tRNAs with the same anticodon but different sequences



**Fig. 4.** The estimated DFE indicates that high proportions of deleterious mutations in tRNAs are under strong selection. (A) Estimated DFE of new deleterious mutations for tRNA genes and inner 3' flanking regions shown in human, mouse, *A. thaliana*, and *D. melanogaster*. Proportions of deleterious mutations are shown for each bin of purifying selection strength, estimated on a scale of  $N_e S$ . Species are arranged by increasing effective population size. (B–D) Low-frequency SNPs plotted as in Fig. 1*C* for mouse (B), *A. thaliana* (C), and *D. melanogaster* (D).

may have different functions, and this may affect strength of selection at each locus as well. Indeed, a significantly greater proportion of sites are invariant (Fisher's exact test,  $P < 7.50e-5$ ) and fewer sites are divergent ( $P < 3.85e-8$ ) in active single-copy human tRNA genes than in active multicopy human tRNA genes. We observe the same patterns in the inner 5' ( $P < 5.87e-5$ ;  $P < 0.025$ ) and inner 3' ( $P < 8.90e-5$ ;  $P < 4.04e-4$ ) flanks of active tRNA genes, suggesting increased transcription of active multicopy tRNA genes. However, few SNP data are available for multicopy tRNAs compared with single-copy tRNAs, limiting our ability to identify consistent differences among tRNA subgroups.

**tRNA Loci Contribute Disproportionately to Mutational Load.** Our discovery of a highly elevated mutation rate at tRNA loci suggests that tRNA genes may contribute disproportionately to mutational load, the reduction in individual fitness due to deleterious mutations (55, 56). To estimate the relative mutation rates at active tRNA loci, we calculated the average ratios of  $\theta$  for the inner 3' and 5' flanking regions of active human tRNA genes to the untranscribed reference regions using the approach of Messer (*Methods*) (29). We estimate  $\theta$  in the flanking regions instead of the tRNAs because strong selection can cause underestimation of  $\theta$  (29), and our results indicate that active human tRNAs are subject to strong selection while the flanking regions are likely selectively neutral (Fig. 44). We therefore estimate that the mutation rate is between 7.24 (inner 3'; 95% CI 7.12–7.33) and 10.36 (inner 5'; 95% CI 10.16–10.41) times greater at tRNA loci than the genome-wide average. Given that there are 25,852 base pairs of active human tRNA sequence, and using  $1.45e-8$  as the genome-wide mutation rate (57), we estimate that U (the genome-wide rate of deleterious mutation per diploid genome) contributed by tRNAs is between 0.0054 and 0.0078. Since active tRNAs make up only 0.0009% of the human genome (2), this implies that mutations in tRNAs contribute disproportionately to mutational load. Our findings highlight that mutations at tRNA loci are likely an important source of fitness and disease variation in human populations.

## Conclusions

Our findings demonstrate that the exceptional transcription rates of tRNA genes cause a similarly substantial increase in mutation rates through TAM. Our results are consistent across a broad range of taxonomically diverse species, indicating that elevated mutation rates due to TAM and strong purifying selection are widespread and may be a good predictor of relative tRNA gene transcription levels. The conflict between extreme TAM and consequent strong purifying selection at tRNA loci is potentially an unappreciated source of genetic disease and may have a profound impact on human fitness that is yet to be fully addressed.

## Materials and Methods

**Defining tRNA Loci and Flanking Regions.** We used tRNA coordinates from GtRNAdb (2) for the human, *M. musculus*, *D. melanogaster*, and *A. thaliana* genomes. For each species, we defined untranscribed reference regions by searching 10 kilobases upstream of each tRNA and selecting a 200-nucleotide tract. If this tract was within a highly transcribed region of the genome [based on genome-wide ChIP data (31)], overlapped a conserved element [defined as a region with a phastCons log odds score greater than 0 (28)], was within 1,000 nucleotides of a known gene (48), or overlapped a reference region assigned to another tRNA, we selected a different tract 1,000 bases further upstream and repeated the selection until we found an acceptable region. For the mouse genome, we checked known genes, previously assigned reference regions, and conserved elements. For the *D. melanogaster* and *A. thaliana* genomes, we began our searches only 1,000 bases upstream of each tRNA and searched for 200-nucleotide tracts that were at least 100 nucleotides away from any annotated genetic element (58, 59) due to the high functional densities of these species' genomes.

For each tRNA in all species, we defined the inner 5' flank as the 20 bases immediately upstream of the 5' end of the tRNA gene on the coding strand and the outer 5' flank as the 20 bases directly upstream of the inner 5' flank. The

inner 3' flank refers to the 10 bases downstream of the tRNA gene, and the outer 3' flank refers to the 30 bases downstream of these 10 bases. We made these decisions based on inflection points in our data, as the flanking regions up to 20 bases upstream and 10 bases downstream of tRNA genes have less variation. Transcription usually ends about 10 bases downstream of tRNA genes (36).

**Classifying tRNAs Based on Breadth of Expression.** The Roadmap Epigenomics Consortium compiled genome-wide epigenomic data across 127 human tissues and cell lines to characterize the chromatin state across the genome (31). We analyzed the regions surrounding each tRNA in each epigenome sample and used clustering to classify each genomic region according to its most common epigenomic state. We classified all human tRNAs based on the epigenomic state annotation in the genome. In the corresponding model, regions in state 1 are likely to be transcribed. The 342 tRNAs in state 1 in at least 4 of the 127 tissues analyzed are active tRNAs, and we consider the remaining 254 tRNAs to be inactive. To classify mouse tRNAs, we used a 15-state Hidden Markov Model based on ChIP data in which states 5 and 7 corresponded to regions near active promoters (50). We considered the 272 tRNAs in genomic regions annotated as state 5 or 7 in at least 3% of tissues as active and the remaining 188 tRNAs as inactive.

**Aligning tRNAs.** We aligned all tRNAs across all species using covariance models (60) and assigned coordinates to each position in each tRNA and flank based on the Sprinzl numbering system (61). We averaged the PhyloP, divergence, and low-frequency SNP data for all sites assigned to the same Sprinzl coordinate for their respective tRNA loci. Because some tRNAs have variations in structure (2), this alignment was necessary for positionwise comparisons between tRNAs. We filtered tRNAs with fewer than 50 aligned bases from our analyses. If a conserved element (regions with a phastCons log odds score greater than 0; ref. 28) was present 4–10 bases up- or downstream of a tRNA, the tRNA was excluded from our analyses, as these regions might contribute to the secondary structure of mature tRNAs and be subject to anomalous levels of selection. We also excluded nuclear-encoded mitochondrial tRNA genes.

**Parsing Variation Data.** We analyzed human variation data from the African superpopulation of 661 humans from phase 3 of the 1000 Genomes Project (62). We acquired *D. melanogaster* variation data for the Siavonga, Zambia populations from the *Drosophila* Genome Nexus Database (58, 59). We obtained *M. musculus* and *A. thaliana* data from Waterston et al. (63) and the *Arabidopsis* Genome Initiative (64), respectively. All nonhuman data were aligned and genotypes curated as described in ref. 65.

Within each gene, flank, or reference region, we considered positions with minor allele frequencies between 0 and 0.05 to be low-frequency SNPs. We also determined the frequency each class of mutations (e.g., A→G) within each region of each tRNA locus where the identity of each base is defined according to the coding strand sequence. We found the frequency of divergences and low-frequency SNPs by position across all tRNAs and flanking regions. For conservation studies across multiple species, we used the PhyloP track (28) from the University of California, Santa Cruz (UCSC) Genome Browser (48) and calculated the average score for each position within the tRNAs and flanking regions. No PhyloP data were available for *A. thaliana* (28). For direct comparisons between the species of interest and an outgroup, we used the Multiz track from the UCSC Table Browser (66) and the Stitch MAFs tool from Galaxy (67) to create sequence alignments. Details are available in *SI Appendix*.

**Transcription Factor Binding.** The ENCODE Project Consortium used ChIP-seq data to identify binding regions for regulatory factors (38–40) including the TBP and Pol3 transcription factors in the human genome (20). These data were taken from the UCSC Genome Browser (48). The intensity of a given peak correlates with a greater frequency of transcription factor binding to that region. For each human tRNA, we found the strongest TBP peak in the 50 base pairs immediately upstream of the tRNA across the GM12878, H1-hESC, HeLa-S3, HepG2, and K562 cell lines. We also calculated the average PhyloP score across the flanking regions for each tRNA (28) and used Spearman's rank correlation test on these data.

**Correlating Variation to Cell-Line Read Counts.** Zheng et al. (46) used demethylation sequencing to detect tRNAs within HEK293T cells (46, 68). We used Spearman's rank correlation tests to correlate mature tRNA transcript read counts and tRNA and flanking region conservation. Because Zheng et al. (46) sequenced mature tRNAs, which are often encoded by multiple genes, we excluded identical genes to control for the correlation between gene copy number and overall expression (Fig. 2 C and D and refs. 32 and 46).

Separately, we summed the average PhyloP scores at these loci and correlated the summed scores to total tRNA read counts (*SI Appendix, Fig. S4*).

**Estimating the Distribution of Fitness Effects.** We estimated the DFE for each species using the method of Keightley et al. (69) and the DFE- $\alpha$  software. See *SI Appendix* for details.

**Estimating the Mutation Rate in Active tRNA Genes.** We used the equation  $\theta(k) = kG^k$  (defined in ref. 29) to estimate the mutation rate at active tRNA loci. We calculated the ratios of  $\theta$  in active tRNA flanking regions to  $\theta$  in the

reference regions for  $k = 1, 2, 3$  and bootstrapped by tRNA loci to calculate 95% CIs. See *SI Appendix* for more details.

**ACKNOWLEDGMENTS.** We thank Craig Mello for helpful discussions and input on this project, Brian Lin for the Fig. 1 legend, Andrew Holmes for assisting with mouse data analysis, and the R.B.C.-D. and T.M.L. laboratories for suggestions and feedback. This work was supported by National Human Genome Research Institute/NIH Grant 2R01HG006753-04A1 (to T.M.L.) and National Institute of General Medical Sciences/NIH Grant R35GM128932-01 (to R.B.C.-D.). B.P.T. was funded by NIH Training Grant T32 HG008345.

- Tang DT, Glazov EA, McWilliam SM, Barris WC, Dalrymple BP (2009) Analysis of the complement and molecular evolution of tRNA genes in cow. *BMC Genomics* 10:188.
- Chan PP, Lowe TM (2016) GTRNAdb 2.0: An expanded database of transfer RNA genes identified in complete and draft genomes. *Nucleic Acids Res* 44:D184–D189.
- Suzuki T, Nagao A, Suzuki T (2011) Human mitochondrial tRNAs: Biogenesis, function, structural aspects, and diseases. *Annu Rev Genet* 45:299–329.
- Schimmel P (2018) The emerging complexity of the tRNA world: Mammalian tRNAs beyond protein synthesis. *Nat Rev Mol Cell Biol* 19:45–58.
- Kutter C, et al. (2011) Pol III binding in six mammals shows conservation among amino acid isotypes despite divergence among tRNA genes. *Nat Genet* 43:948–955.
- Ishimura R, et al. (2014) RNA function. Ribosome stalling induced by mutation of a CNS-specific tRNA causes neurodegeneration. *Science* 345:455–459.
- Kirchner S, Ignatova Z (2015) Emerging roles of tRNA in adaptive translation, signalling dynamics and disease. *Nat Rev Genet* 16:98–112.
- Molla-Herman A, Vallés AM, Ganem-Elbaz C, Antoniewski C, Huynh JR (2015) tRNA processing defects induce replication stress and Chk2-dependent disruption of piRNA transcription. *EMBO J* 34:3009–3027.
- Jinks-Robertson S, Bhagwat AS (2014) Transcription-associated mutagenesis. *Annu Rev Genet* 48:341–359.
- Gnatt AL, Cramer P, Fu J, Bushnell DA, Kornberg RD (2001) Structural basis of transcription: An RNA polymerase II elongation complex at 3.3 Å resolution. *Science* 292:1876–1882.
- Gaillard H, Aguilera A (2016) Transcription as a threat to genome integrity. *Annu Rev Biochem* 85:291–317.
- Kim N, Jinks-Robertson S (2012) Transcription as a source of genome instability. *Nat Rev Genet* 13:204–214.
- Aguilera A, García-Muse T (2013) Causes of genome instability. *Annu Rev Genet* 47:1–32.
- Helmrich A, Ballarino M, Tora L (2011) Collisions between replication and transcription complexes cause common fragile site instability at the longest human genes. *Mol Cell* 44:966–977.
- Timakov B, Liu X, Turgut I, Zhang P (2002) Timing and targeting of P-element local transposition in the male germline cells of *Drosophila melanogaster*. *Genetics* 160:1011–1022.
- Gómez-González B, Aguilera A (2007) Activation-induced cytidine deaminase action is strongly stimulated by mutations of the THO complex. *Proc Natl Acad Sci USA* 104:8409–8414.
- Green P, Ewing B, Miller W, Thomas PJ, Green ED; NISC Comparative Sequencing Program (2003) Transcription-associated mutational asymmetry in mammalian evolution. *Nat Genet* 33:514–517.
- Taylor BJ, Wu YL, Rada C (2014) Active rnap pre-initiation sites are highly mutated by cytidine deaminases in yeast, with aid targeting small rna genes. *elife* 3:e03553.
- Saini N, et al. (2017) APOBEC3B cytidine deaminase targets the non-transcribed strand of tRNA genes in yeast. *DNA Repair (Amst)* 53:4–14.
- White RJ (2011) Transcription by RNA polymerase III: More complex than we thought. *Nat Rev Genet* 12:459–463.
- Zhang J, Ferré-D'Amaré AR (2016) The tRNA elbow in structure, recognition and evolution. *Life (Basel)* 6:E3.
- Sun C, et al. (2018) Roles of tRNA-derived fragments in human cancers. *Cancer Lett* 414:16–25.
- 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. *Biochemistry* 39:9909–9916.
- Hopper AK (2013) Transfer RNA post-transcriptional processing, turnover, and subcellular dynamics in the yeast *Saccharomyces cerevisiae*. *Genetics* 194:43–67.
- Hasler D, et al. (2016) The Lupus Autoantigen La prevents mis-channeling of tRNA fragments into the human MicroRNA pathway. *Mol Cell* 63:110–124.
- Lee YS, Shibata Y, Malhotra A, Dutta A (2009) A novel class of small RNAs: tRNA-derived RNA fragments (tRFs). *Genes Dev* 23:2639–2649.
- Maraia RJ, Lamichhane TN (2011) 3' processing of eukaryotic precursor tRNAs. *Wiley Interdiscip Rev RNA* 2:362–375.
- Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A (2010) Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res* 20:110–121.
- Messer PW (2009) Measuring the rates of spontaneous mutation from deep and large-scale polymorphism data. *Genetics* 182:1219–1232.
- Karolchik D, et al. (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids Res* 32:D493–D496.
- Kundaje A, et al.; Roadmap Epigenomics Consortium (2015) Integrative analysis of 111 reference human epigenomes. *Nature* 518:317–330.
- Bloom-Ackermann Z, et al. (2014) A comprehensive tRNA deletion library unravels the genetic architecture of the tRNA pool. *PLoS Genet* 10:e1004084.
- Doran JL, Bingle WH, Roy KL (1988) Two human genes encoding tRNA(GCCGly). *Gene* 65:329–336.
- Dieci G, Sentenac A (1996) Facilitated recycling pathway for RNA polymerase III. *Cell* 84:245–252.
- Cieśla M, Boguta M (2008) Regulation of RNA polymerase III transcription by Maf1 protein. *Acta Biochim Pol* 55:215–225.
- Orioli A, et al. (2011) Widespread occurrence of non-canonical transcription termination by human RNA polymerase III. *Nucleic Acids Res* 39:5499–5512.
- Juo ZS, et al. (1996) How proteins recognize the TATA box. *J Mol Biol* 261:239–254.
- Dunham I, et al.; ENCODE Project Consortium (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature* 489:57–74.
- Myers RM, et al.; ENCODE Project Consortium (2011) A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biol* 9:e1001046.
- Kharchenko PV, Tolstorukov MY, Park PJ (2008) Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat Biotechnol* 26:1351–1359.
- Roberts DN, Stewart AJ, Huff JT, Cairns BR (2003) The RNA polymerase III transcriptome revealed by genome-wide localization and activity-occupancy relationships. *Proc Natl Acad Sci USA* 100:14695–14700.
- Mason PB, Struhl K (2003) The FACT complex travels with elongating RNA polymerase II and is important for the fidelity of transcriptional initiation in vivo. *Mol Cell Biol* 23:8323–8333.
- Kuras L, Kosa P, Mencia M, Struhl K (2000) TAF-Containing and TAF-independent forms of transcriptionally active TBP in vivo. *Science* 288:1244–1248.
- Zanton SJ, Pugh BF (2004) Changes in genomewide occupancy of core transcriptional regulators during heat stress. *Proc Natl Acad Sci USA* 101:16843–16848.
- Li XY, Virbasius A, Zhu X, Green MR (1999) Enhancement of TBP binding by activators and general transcription factors. *Nature* 399:605–609.
- Zheng G, et al. (2015) Efficient and quantitative high-throughput tRNA sequencing. *Nat Methods* 12:835–837.
- Thul PJ, et al. (2017) A subcellular map of the human proteome. *Science* 356:eaal3321.
- Kent WJ, et al. (2002) The human genome browser at UCSC. *Genome Res* 12:996–1006.
- Schmidt S, et al. (2008) Hypermutable non-synonymous sites are under stronger negative selection. *PLoS Genet* 4:e1000281.
- Bogu GK, et al. (2015) Chromatin and RNA maps reveal regulatory long noncoding RNAs in mouse. *Mol Cell Biol* 36:809–819.
- Tenesa A, et al. (2007) Recent human effective population size estimated from linkage disequilibrium. *Genome Res* 17:520–526.
- Phifer-Rixey M, et al. (2012) Adaptive evolution and effective population size in wild house mice. *Mol Biol Evol* 29:2949–2955.
- Cao J, et al. (2011) Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat Genet* 43:956–963.
- Shapiro JA, et al. (2007) Adaptive genic evolution in the *Drosophila* genomes. *Proc Natl Acad Sci USA* 104:2271–2276.
- Haldane JBS (1937) The effect of variation of fitness. *Am Nat* 71:337–349.
- Agrawal AF, Whitlock MC (2012) Mutation load: The fitness of individuals in populations where deleterious alleles are abundant. *Annu Rev Ecol Syst* 43:115–135.
- Narasimhan VM, et al. (2017) Estimating the human mutation rate from autozygous segments reveals population differences in human mutational processes. *Nat Commun* 8:303.
- Lack JB, et al. (2015) The *Drosophila* genome nexus: A population genomic resource of 623 *Drosophila melanogaster* genomes, including 197 from a single ancestral range population. *Genetics* 199:1229–1241.
- Lack JB, Lange JD, Tang AD, Corbett-Detig RB, Pool JE (2016) A thousand fly genomes: An expanded *Drosophila* genome nexus. *Mol Biol Evol* 33:3308–3313.
- Lowe TM, Chan PP (2016) tRNAcanSE On-line: Integrating search and context for analysis of transfer RNA genes. *Nucleic Acids Res* 44:W54–7.
- Sprinzl M, Horn C, Brown M, Ioudovitch A, Steinberg S (1998) Compilation of tRNA sequences and sequences of tRNA genes. *Nucleic Acids Res* 26:148–153.
- Auton A, et al.; 1000 Genomes Project Consortium (2015) A global reference for human genetic variation. *Nature* 526:68–74.
- Waterston RH, et al.; Mouse Genome Sequencing Consortium (2002) Initial sequencing and comparative analysis of the mouse genome. *Nature* 420:520–562.
- Arabidopsis Genome Initiative (2000) Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature* 408:796–815.
- Corbett-Detig RB, Hartl DL, Sackton TB (2015) Neutral selection constrains neutral diversity across a wide range of species. *PLoS Biol* 13:e1002112.
- Blanchette M, et al. (2004) Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res* 14:708–715.
- Afgan E, et al. (2016) The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res* 44:W3–W10.
- Cozen AE, et al. (2015) ARM-seq: AlkB-facilitated RNA methylation sequencing reveals a complex landscape of modified tRNA fragments. *Nat Methods* 12:879–884.
- Keightley PD, Eyre-Walker A (2007) Joint inference of the distribution of fitness effects of deleterious mutations and population demography based on nucleotide polymorphism frequencies. *Genetics* 177:2251–2261.