






# How Does Transcription-Associated Mutagenesis Shape tRNA Microevolution?

Hector Baños <sup>1,\*</sup>, Ling Wang <sup>2</sup>, Corinne Simonti <sup>2</sup>, Annalise B. Paaby <sup>2</sup>,  
Christine Heitsch <sup>3</sup>

<sup>1</sup>Department of Mathematics, California State University, San Bernardino, CA 92407, USA

<sup>2</sup>School of Biological Sciences, Georgia Institute of Technology, Atlanta, GA 30332, USA

<sup>3</sup>School of Mathematics, Georgia Institute of Technology, Atlanta, GA 30332, USA

\*Corresponding author: E-mail: hector.banos@csusb.edu.

Accepted: February 23, 2026

## Abstract

Transfer RNAs (tRNAs) are among the most highly conserved and frequently transcribed genes. Recent studies have demonstrated that tRNAs experience exceptionally high rates of transcription-associated mutagenesis (TAM) as well as strong purifying selection. How the mutational input of TAM, which induces a nonuniform distribution of nucleotide substitutions, affects the fitness of tRNA molecules is unclear. Secondary structure in tRNAs is strongly conserved over macroevolutionary time, suggesting that mutations that disrupt paired sites may be especially deleterious, and TAM-induced mutations primarily involve nucleotide transitions, which tend to preserve base-pairing stability. To examine how TAM affects tRNA molecule fitness and shapes tRNA evolution over short timescales, we analyzed tRNA allelic variation in contemporary *Caenorhabditis elegans* strains. We propose a model of tRNA microevolution driven by TAM and demonstrate that the observed secondary structure characteristics align with our predicted TAM-biased patterns. Furthermore, we developed a continuous Markov substitution model that incorporates TAM-specific mutational biases. This TAM-biased model fits the *C. elegans* tRNA data more effectively than standard models, such as the general time-reversible model. Based on these results, we conclude that tRNAs in natural populations carry substantial levels of structure-destabilizing mutations, which may be tolerated but nevertheless likely induce meaningful fitness costs. Our findings are consistent with recent experimental studies on tRNA fitness in yeast, but challenge prior theoretical and computational analyses that emphasize RNA base-pairing as a primary determinant in genotype-phenotype systems.

**Key words:** Transcription-associated mutagenesis, tRNA, Microevolution, Markov substitution model.

## Significance

Transfer RNAs (tRNAs) are ancient molecules, encoded as genes in all living systems. tRNA genes are known to experience exceptionally high rates of both mutation and purifying selection, but how these opposing evolutionary forces shape tRNA evolution is unclear. We developed a sequence substitution model specific to tRNA mutagenesis and applied it to standing variation in a natural population in order to infer how mutation and selection affect the structural stability of tRNA molecules.

© The Author(s) 2026. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

## Introduction

Transfer RNAs (tRNAs) are ancient and indispensable genes essential for protein synthesis, renowned for their evolutionary conservation (Tang et al. 2009; Zhang and Ferré-D'Amaré 2016). This conservation is widely attributed to the structural integrity necessary for their function (Westhof et al. 2022), maintained by strong selection for base-pair preservation. As a result, the preservation of secondary structure is frequently used as a proxy for tRNA functionality, making tRNAs a classic model for exploring genotype–phenotype relationships.

However, the evolutionary forces shaping tRNA sequences remain poorly understood, particularly at the population level (Ishimura et al. 2014). The preservation of secondary structure over macroevolutionary time emphasizes its importance to molecular function, but tRNA sequence confers essential functionality beyond the context of folding. For example, sequence identity elements specify interactions with aminoacyl-tRNA synthetases and other factors necessary for appropriate amino acid charging (Giegé and Eriani 2023). Recent yeast experiments show that single mutations in loop regions can significantly affect fitness, while single mutations in the acceptor stem are surprisingly well-tolerated (Li et al. 2016). These findings provide the first empirical estimates of tRNA fitness effects and challenge the prioritization of paired sites over loop regions implied by the thermodynamic model (Fontana and Schuster 1998; Reidys et al. 2001; Aguirre et al. 2011). Alterations to sequence identity elements may disrupt coadaptations that arise over relatively short timescales (Meiklejohn et al. 2013; Adrion et al. 2015), but the fitness dynamics of tRNAs over microevolutionary time, including how selection shapes the distribution of standing variants in populations, remains unknown.

Transcription-associated mutagenesis (TAM) (Jinks-Robertson and Bhagwat 2014) has been identified as a major contributor to mutation rates in highly transcribed regions, and tRNAs are among the most highly transcribed genes in the genome (Palazzo and Lee 2015; Boivin et al. 2018). During transcription, deamination of the coding strand, along with DNA repair mechanisms responding to this deamination, is associated with an increase in C → T and consequently G → A substitutions in the coding strand (Green et al. 2003; Jinks-Robertson and Bhagwat 2014). tRNAs within and across species exhibit mutational variation with signatures of historical TAM, including in flanking regions just upstream and downstream of tRNA gene bodies (Thornlow et al. 2018). Both the gene and its flanking regions are vulnerable to TAM during transcription and are presumed to experience the same rate of mutation. However, strong purifying selection on the

gene itself purges some mutations, while flanking regions can accumulate exceedingly high numbers of mutations. Thus, these regions provide a record of the historical expression level and mutation rate of that gene (Thornlow et al. 2018, 2020).

The exceedingly high rates of mutation and purifying selection at tRNAs has led to the hypothesis that eukaryotic genomes carry substantial mutational load at tRNA loci (Thornlow et al. 2018). However, how this translates to a disease burden, including how the fitness of individual tRNA molecules are compromised, remains unexplored. A key question is how the nonuniform nucleotide substitutions associated with TAM might alleviate or exacerbate destabilization of tRNA secondary structure. For example, selection might tolerate transitions over transversions, if they are less likely to disrupt nucleotide pairs in stems; under this scenario, the C → T and G → A transitions induced most frequently by TAM may not be as deleterious as they could be. Also under this scenario, the coincidence of the highest frequency and least deleterious mutations would mean that selection would not erode the TAM signature from standing variation even as most mutations are purged (Thornlow et al. 2018). Our focus is to better assess the signature of TAM in tRNA allelic variation and determine how TAM affects overall tRNA fitness.

It is important to note that beyond TAM, additional mechanisms govern the occurrence of mutations that comprise standing variation in tRNA genes, including transcription-coupled repair (Svejstrup 2002). However, a strong and growing body of literature indicates that tRNA genes experience unusually high rates of mutation as a consequence of their exceptional rate of transcription (Murillo-Recio et al. 2025), and that TAM dominates the mutational landscape of tRNA genes (Park et al. 2012; Jinks-Robertson and Bhagwat 2014; Thornlow et al. 2018; Liu and Zhang 2020; Thornlow et al. 2020).

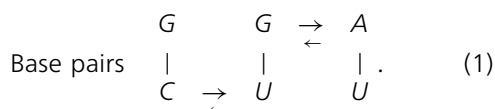
In this study, we assess how TAM affects tRNA fitness and microevolution by analyzing standing genetic variation in tRNA genes in *Caenorhabditis elegans*. Using a novel continuous Markov substitution model, we evaluate TAM's mutational signature and the role of selection on base-pair preservation. Aligned with the experimental results in Li et al. (2016) and Domingo et al. (2018), the data, supported by our model, show that paired sites experience higher substitution frequencies than unpaired ones under TAM. Additionally, we observed that most substitutions, whether preserving base-pairing or not, decrease thermodynamic stability, suggesting that TAM compromises structural integrity rather than maintaining it. These results highlight how TAM shapes tRNA evolution in ways beyond structural constraints under the thermodynamic model.

## Background

Markov models are widely used to study evolutionary processes, particularly DNA substitutions (Felsenstein 2003; Yang 2006). While these models are traditionally applied to gene tree inference across species, our approach departs from this by using them to describe allele formation within a population.

These models use different parameters to capture the various mutational mechanisms and selection pressures shaping DNA substitutions. Hence, the Markov exchangeability matrix for the substitution from  $i$  to  $j$  (denoted  $i \rightarrow j$ ) for  $i, j \in \{A, C, G, T\}$  may differ depending on the combination of nucleotides. In particular, transitions are expected to occur much more frequently than transversions since they preserve the purine ( $A \leftrightarrow G$ ) versus pyrimidine ( $C \leftrightarrow T$ ) distinction.

Moreover, for noncoding RNA genes, it is likely that transitions are further enhanced as they allow Watson–Crick base pairs G–C or C–G (denoted GC) and A–U/U–A (AU)<sup>1</sup> to interconvert via the wobble G–U/U–G (GU) pairing:



Note that of the three types of canonical base pairings, GC is the most thermodynamically favorable, followed by AU, and then GU. Hence, based on the pairing conservation, one would expect that, in the coding strand,  $C \rightarrow T$  transitions to be favored in GC pairs (resulting in GU pairs),  $A \rightarrow G$  ones in AU (resulting in GU pairs), and both  $G \rightarrow A$  and  $T \rightarrow C$  in GU (resulting in AU and GC pairs, respectively) with the last as the most preferred, and the first as the least, due to thermodynamic stability. TAM is the dominant source of mutations at tRNAs (Thornlow et al. 2018), which is biased toward more frequent  $C \rightarrow T$  and  $G \rightarrow A$  transitions on the coding strand, denoted in red above.

## Secondary Structure Conservation in RNA Macroevolution

A noncoding RNA molecule like tRNA is composed of a single-stranded sequence which folds hierarchically (Tinoco and Bustamante 1999) into a functional structure. The intrasequence base pairings of the secondary structure are the critical scaffold which is then organized by tertiary interactions into the final 3D conformation. In some cases, it is the pattern of base pairing, and not necessarily the nucleotide content, that is evolutionarily conserved. This selection pressure can be so strong that the occurrence of compensatory mutations in base pairs, identified under multiple sequence alignment, is the gold-standard

for RNA structural inference (Eddy and Durbin 1994; Cannone et al. 2002; Griffiths-Jones et al. 2005).

More precisely, a *secondary structure* for an RNA sequence  $R$  of length  $n$  is a set of paired indices  $S(R) = \{(i, j) \mid 1 \leq i < j \leq n\}$  such that the nucleotides in positions  $i$  and  $j$  form a canonical base pair, ie either Watson–Crick (GC, AU) or wobble (GU). Suppose that  $R'$  is a *one-point mutant* of  $R$ , meaning that it differs from  $R$  by a single substitution. Suppose further that  $S(R)$  is also a valid secondary structure for  $R'$  (meaning that the nucleotide mutation in  $R'$  does not disrupt any pairing of  $S(R)$ ). Then, either the substitution changed an unpaired nucleotide or it changed a paired one as described in Equation (1). In these cases, we say that the substitution *preserves pairing potential*. Otherwise, it necessarily disrupts the pairing involving the mutated nucleotide.

It is critical to emphasize that a substitution being pairing potential preserving says nothing about the secondary structure of  $R'$ . The mutational constraint only requires that  $R'$  *can* assume the same set of pairings as  $S(R)$ , not that it *does*. It is, however, a necessary condition for  $R'$  to be a *neutral neighbor* of  $R$ .

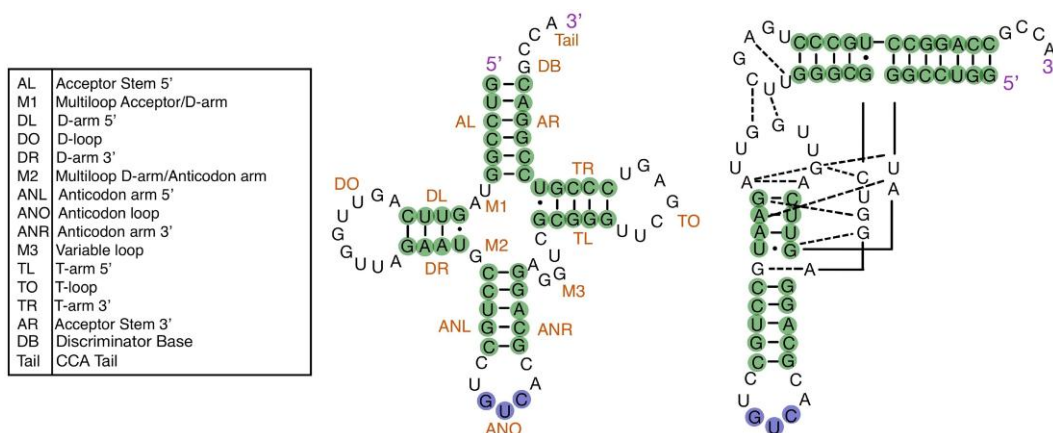
Using RNA secondary structure formation as a model genotype–phenotype system (Schuster et al. 1994; Cowperthwaite and Meyers 2007), the theory of neutral networks was developed as a framework for understanding structural fitness landscapes. In this context, evolution is modeled as a series of one-point mutations, and a substitution is considered “neutral” if it does not alter the secondary structure. The conservation of pairing (and not just its potential preservation) is most often determined by free energy minimization (Zuker and Stiegler 1981) under the nearest neighbor thermodynamic model (NNTM). Hence  $R'$  is a neutral neighbor of  $R$  if it differs by a single substitution and retains the same minimum free energy (MFE) structure, ie the one-point mutation must be *MFE-preserving*.

Under this NNTM optimization interpretation of fitness neutrality, disruption of canonical base pairs is highly deleterious. Hence, the macroevolution expectation of secondary structure conservation implies that substitutions should follow these patterns:

- (U) A bias against mutations at paired sites.
- (P1) A bias toward  $C \rightarrow T$  in GC pairs and  $A \rightarrow G$  in AU ones if a Watson–Crick pairing is mutated.
- (P2) A bias toward transitions if a wobble pairing is mutated, with  $T \rightarrow C$  clearly favored over  $G \rightarrow A$ .

## Experimental Insights into tRNA Microevolution

Transfer RNA has one of the most direct genotype–phenotype relationships possible due to the tight



**Fig. 1.** Representative tRNA secondary structure and tertiary interactions. From left to right: table lists the structural components, including the CCA Tail, which are then labeled on the secondary structure. The four-armed cloverleaf is closed by the acceptor (A) stem, and contains three hairpin stem-loop structures: the D, AN, and T arms. Paired nucleotides as well as the AN are highlighted. Most pairings are Watson–Crick (GC/AU, dash) but two wobble ones (GU, dot) are present. Various tertiary interactions (dashed lines) stabilize the 3D structure. The overall L-shape is critical to ribosome binding, and hence protein biosynthesis (Center and Right Figures are modified with permission from Liu (2013) and Klemm et al. (2016)).

coupling of sequence/structure/function. In order to read the mRNA codon and deliver the amino acid, a tRNA sequence must fold into a canonical 3D “L” shape. This essential structure is remarkably conserved across all three domains of life (Pak et al. 2017). Although variations exist, including mitochondrial tRNAs that adopt extreme, yet functional shapes (Jühling et al. 2018; Ozerova et al. 2024), the cloverleaf secondary structure remains the core scaffold of this critical arrangement. As illustrated<sup>2</sup> in Fig. 1, the intrasequence base pairing forms four runs of stacked base pairs, with three hairpin loops and the central multibranch loop. Three nucleotides in the anticodon (AN) loop target the complementary codon in the mRNA molecule while the corresponding amino acid, attached at the Tail, is loaded on the opposing acceptor stem. As pictured, the variable portion (M3) of the multiloop is unpaired, but it can sometimes be much longer (Berg and Brandl 2021) in which case it typically includes pairings.

As a model genotype–phenotype system, tRNA folding has attracted considerable theoretical attention (Fontana and Schuster 1998; Reidys et al. 2001; Aguirre et al. 2011) in macroevolution studies. These studies use secondary structure thermodynamic predictions, such as the `RNAfold` function of the ViennaRNA 2.0 package (Lorenz et al. 2011) which outputs a single MFE structure, as a proxy for evolutionary fitness. In this neutral-neighbor context, any change in the base pairing pattern, ie the positions which are paired, is deemed deleterious.

More recently, however, tRNA has been the focus of population growth experiments (Yona et al. 2013; Li et al. 2016; Domingo et al. 2018) that use sequencing

technology to explore fitness at the microevolution scale. All three studies investigated a particular *Saccharomyces cerevisiae* Arginine gene: tRNA<sup>Arg</sup><sub>CCU</sub>, which is distinctive as the only tRNA gene with the CCU AN.<sup>3</sup> This gene is not essential for survival, but its deletion compromises fitness as measured by population growth.

By deleting the gene entirely, Yona et al. (2013) demonstrated that within 200 generations the wild-type growth rate could be recovered. Investigating further they confirmed that this was achieved by mutating one of the 11 tRNA<sup>Arg</sup><sub>UCU</sub> copies available, and not always the same one. Hence, a single C → T substitution in the wobble position of the AN was sufficient to recover the wild-type fitness. Moreover, they then confirmed that such AN switching can be found across all three domains of life, leading to the conclusion that this is a wide-spread adaptation mechanism. Hence, while the AN is generally expected to be under very strong selection pressure, it can and will evolve rapidly to meet new translation demands.

In contrast, Li et al. (2016) comprehensively characterized all one-point mutants of tRNA<sup>Arg</sup><sub>UCU</sub>. They classified 1% as beneficial, 42% as deleterious, and the remaining majority (57%) as “(nearly) neutral.” Generally, the least mutable positions were in loops, and the most in stems. As expected, mutations in the AN itself significantly reduced fitness. Nonetheless, mutations at three positions in the T-loop (positions 53, 54, and 55), two in the T-stem (positions 52 and 60, which form a the base pair “closing” the T-loop), and one in the D-stem also had a substantial negative impact (position 18). Conversely, the acceptor stem seemed to tolerate all single substitutions, as did the variable loop.

This robustness to one-point mutants in key structural components was echoed by Domingo et al. (2018) who investigated all existing variants of tRNA<sup>Arg</sup><sub>UCU</sub> found across yeast species. They identified 14 mutations across 10 sites in six other genomes: 11 in the acceptor stem (five in AL, six in AR) and two in the variable loop (M3) with one from the anticodon stem (ANL). All but one of the AL/AR mutations came together in compensatory pairs, of which one was a GU, three AU, and one GC. The remaining AR mutation formed a GU pair as did the ANL one. They then generated all possible combinations of these substitutions and tested their fitness in growth competition. All 14 one-point mutants exhibited fitness comparable to the wild-type *S. cerevisiae* gene and the extant variants.

Both studies found that fitness decreased rapidly with subsequent mutations, with significant epistatic interactions also identified. The epistasis had a strong negative bias except—unsurprisingly—for compensatory mutations in paired positions. All of this suggests that, while conservation of secondary structure is critical at the macroevolutionary level, the situation may be much more nuanced and flexible across a population.

### Markov Models for DNA Macroevolution

DNA sequences evolving in time under nucleotide substitutions are typically modeled using a continuous-time Markov process. These models are parameterized by a nucleotide distribution vector  $\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$  and  $4 \times 4$  instantaneous rate matrix  $Q$ . The Markov matrix  $M$  of the model is defined as  $M = \exp(Qt)$ , where the entry  $m_{ij} \in M$  denotes the probability that nucleotide  $i$  mutates to nucleotide  $j$  after time  $t$ . Then, the joint distribution of nucleotides in the ancestral and descendant sequences over time  $t$  is given by  $P = \text{diag}(\pi) \exp(Qt)$  where the entry  $p_{ij}$  represents the probability of observing a substitution  $i \rightarrow j$  at a site.

In general, the models are distinguished by the constraints on  $\pi$  and  $Q$ . All forms require that  $\sum \pi_i = 1$  with  $\pi_i \geq 0$ , that  $q_{ij} > 0$  for  $i \neq j$ , and that  $q_{ii} = -\sum_{i \neq j} q_{ij}$ . In other words,  $\pi$  must be a probability distribution, all off-diagonal entries of  $Q$  are positive, and the diagonal ones make the row sums equal zero. In a context where the edge length is unknown, one off diagonal entry is set to 1, to avoid scaling redundancy. Additionally, the rate matrix used is typically normalized by multiplying all original entries of  $Q$  by

$$w = \left( -\sum_{i=1}^4 q_{ii} \pi_i \right)^{-1}.$$

This results in the expected number of substitutions being one per unit of time.

The earliest, and most mathematically tractable models, assume that  $\pi$  is uniform,  $t$  is unknown, and the 11 free entries of  $Q$  come from a set of either 0, 1, or 2 parameters (Jukes and Cantor 1969; Kimura 1980). When there are two, one is for transitions (which preserve the purine/pyrimidine division) and the other for transversions. The addition of a third parameter is used to distinguish transversions which are amino/keto preserving from the weak/strong symmetry.

More biologically realistic models allow nonuniform  $\pi$  and provide a larger number of exchangeability parameters. The most general form, associated with the *continuous general Markov model* (GM) over a single edge, has no additional constraints, and hence 12 free parameters (11 from the rate matrix, and 1 from the edge length). However, this great flexibility often results in over-fitting and is seldom needed to capture the substitution process accurately.

The most common substitution model used for phylogenetic inference is the *general time-reversible* (GTR) one (Tavaré 1986). Such model has five free parameters (six when considering a single edge length) satisfying  $\text{diag}(\pi)Q = Q^T \text{diag}(\pi)$ . This equality holds when the substitution process has the same parameters going forward in time as backward. In a GTR model, the rate matrix  $Q$  can be parameterized via a  $4 \times 4$  non-negative matrix known as an exchangeability matrix. As useful as the assumption of time-reversibility is for macroevolution analyses, it is not clear that it will be appropriate for the *C. elegans* tRNA population data. Although we considered GTR as a possibility, our TAM-biased Markov substitution model is built from the strand-symmetric rate matrix. We note that a model can be both strand-symmetric and time-reversible only if  $\pi_A = \pi_T$  and  $\pi_C = \pi_G$ .

The strand-symmetric (sym) model (Casanelles and Sullivant 2005) also has six free parameters. However, the equality now enforced is on the exchangeability of  $i \rightarrow j$  and its Watson–Crick complement, ie  $C \rightarrow T$  and  $G \rightarrow A$  must have the same parameter. The premise is that, for many mutational mechanisms, it is not possible to determine whether a substitution initially occurred in the coding strand or the template one. We consider such mutations as background noise, and take this matrix as the foundational one for our new TAM-biased model.

### Microevolution Signature of Mutational Bias

Transfer RNAs are not only highly conserved but also highly transcribed, exposing them to high levels of *transcription-associated mutagenesis* (TAM) (Thornlow et al. 2018). Critically, TAM's characteristic signature breaks the strand-symmetry assumption, as well as the time-reversibility one.

During RNA transcription, a hybrid DNA–RNA complex forms, leaving the nontemplate DNA strand exposed and vulnerable to mutagens (Jinks-Robertson and Bhagwat 2014). The formation of noncanonical structures by the unwound DNA further promotes mutagenesis by inhibiting repair pathways and interfering with replication processes (Gaillard and Aguilera 2016). Thus, the more frequently a gene is transcribed, the more vulnerable it is to mutation. TAM induces a nonuniform distribution of substitutions, biased toward a higher frequency of C → T and, secondarily, G → A mutations, on the coding (nontemplate) strand (Gómez-González and Aguilera 2007; Jinks-Robertson and Bhagwat 2014; Williams et al. 2023). This phenomenon arises from enzymatic and biochemical processes that have been most thoroughly characterized in yeast but appear universal across eukaryotes (Jinks-Robertson and Bhagwat 2014; Thornlow et al. 2018).

The discovery that tRNAs are subject to exceptionally high levels of TAM also demonstrated that tRNAs experience very strong purifying selection (Thornlow et al. 2018). This is evident by the reduced level of mutational variation within tRNA gene bodies relative to their flanking regions. Both tRNAs and their flanking regions presumably acquire a similar spectrum of mutations, but those deleterious changes to tRNA that are purged are not observed in population data. In this way, the mutational variation at tRNA flanking regions offers a relatively direct record of the historical TAM process on the gene itself. Here, we use this unique phenomenon to estimate the mutational input to tRNAs from flanking sequence data, and test hypotheses about how TAM affects tRNA fitness and how tRNAs respond to selection.

## Results and Discussion

### Model Validation

We developed a continuous Markov model that describes the distribution of substitutions arising from TAM within a population. The model begins with the concept of  $S$ , which is a concatenation of all tRNA genes into a single sequence;  $S$  represents the consensus sequence across the population, which is an approximation for  $S$ , the corresponding sequence for the most recent common ancestor of the population.  $S'$  represents the mutational variation sequence, which incorporates all observed substitutions in the population. Following several constraints (see Materials and Methods section), the model addresses the nucleotide substitution process from  $S$  to  $S'$ . To incorporate the nonuniform nucleotide substitutions produced by TAM, we estimated  $\tau$ , the TAM bias vector, from population data. Together with the nucleotide base-frequency

vector  $\pi$ , this allowed us to derive  $\pi^\tau$ , the distribution of nucleotides undergoing substitution. Our dataset includes sequences for 581 tRNA genes from 331 *C. elegans* strains collected from the wild; please see Materials and Methods section for complete details of our model, its implementation, and the experimental data.

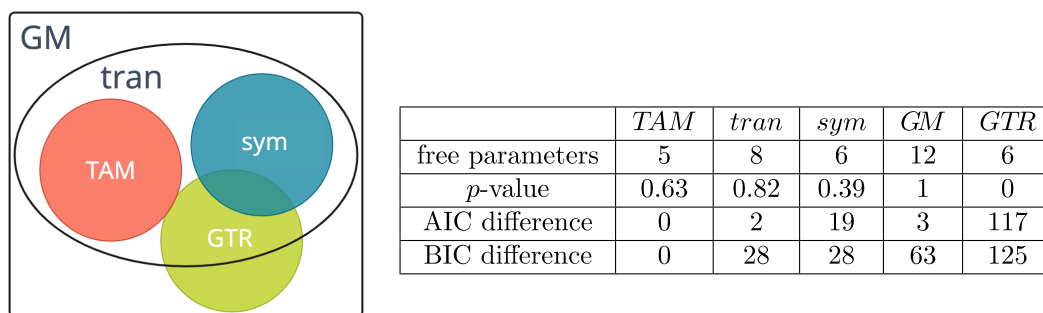
We find good agreement between model hypotheses and population data. Specifically, the new mutational variation sequence  $S'$  reasonably describes the observed tRNA distribution, and the consensus sequence  $S$  is a suitable approximation to the most recent common ancestor of those *C. elegans* alleles. Moreover, the flanking regions yield a robust estimate for the TAM bias vector  $\tau$ , and the tRNA distribution of substituted nucleotides in the tRNA (ie the proportion of nucleotides that mutated to another nucleotide) align with  $\pi^\tau$ , the proposed TAM biased distribution. Finally, of the Markov substitution models considered, the new Tam matrix best describes the data.

We observed 661 variants in our dataset of 581 tRNA genes. Of those alleles, 70 have an indel and cannot be described by a substitution model. Among the 591 remaining, no two differ at the same position, but 60 have two or more substitutions. Hence, there are 531 distinct-site, one-point mutant alleles, and the compound sequence  $S'$  encodes 80.3% of the total tRNA mutational variation in the *C. elegans* population.

We find that  $S$  is a good approximation to the hypothetical most recent common ancestor sequence  $S$  because the distribution of alleles is very sparse. In particular, nearly all consensus tRNA (484/531 = 91.1%) occur in nearly all strains (at least 298/331 so  $\geq 90.0\%$ ). The remaining 47 consensus sequences have more frequently occurring variants. However, of the 45 present in  $\geq 45\%$  of strains, only two have a variant which occurs nearly as often (<15% difference). Hence, in all but four tRNA genes, the consensus clearly dominates, allowing us to consider the evolution of  $S'$  from  $S$ .

The TAM bias vector  $\tau$  was estimated from the tRNA flanking regions, which contained 2070 mutations. Flanking sequence are subject to TAM but not to direct selection, so they provide a putatively unfiltered empirical estimate of the distribution of input mutations. This yielded  $\pi - \pi^\tau = (0.0454, -0.1111, -0.0456, 0.1110) = (\tau_R, -\tau_Y, -\tau_R, \tau_Y)$ . The close match in positive and negative estimates supports our TAM bias hypothesis; averaging yields  $\hat{\tau} = (0.0455, 0.1110)$ .

The nucleotides A/C/G/T in  $S$  are distributed as  $\pi = (0.181, 0.258, 0.323, 0.237)$ .<sup>4</sup> Under our TAM bias model, the expected distribution of substituted nucleotides is  $\pi^\tau = (0.135, 0.370, 0.369, 0.125)$  while  $\hat{\pi}^\tau = (0.117, 0.369, 0.375, 0.139)$  is observed. Under statistical testing, described in Materials and Methods section,



**Fig. 2.** Comparison of rate matrix parameterizations for Markov substitution model: GM, sym, tran, and Tam along with GTR. The AIC and BIC) are reported as the difference from Tam's score as this was the lowest, ie best.

we fail to reject the null hypothesis with a  $P$ -value of 0.510 whereas repeating the analysis with  $\tau = (0, 0)$  yields  $P < <0.001$ . Hence, we conclude that the tRNA substitutions in the data are well described by our proposed TAM bias distribution.

This distribution of substituted nucleotides is then used to parameterize a new rate matrix, denoted Tam, for the continuous Markov substitution process model. As illustrated in Fig. 2, Tam is a special case of tran, which was introduced as a generalization of the standard sym. We also consider the full GM as well as the frequently used GTR.

The calculated  $P$ -values, along with other model selection criteria, are also given in Fig. 2. With  $P < <0.001$ , it is clear that GTR does not fit the data. Having rejected the null hypothesis, and in view of the significantly lower AIC and BIC scores, we conclude that time-reversibility—a common macroevolution assumption—is not an appropriate expectation for this microevolution model.

Among the remaining four matrices with  $P > 0.05$ , the model selection criteria support Tam as the best fit for the substitution process which model the formation of  $S'$  from  $S$ . Tam has the lowest AIC, although the difference with tran is on the borderline of meaningful. (Recall that differences of  $<2$  units in AIC and 6 in BIC are generally considered negligible.) However, BIC is definitive; Tam had the lowest (best) score by a wide margin. Moreover, there is no reason to prefer tran over Tam as the latter provides more specific biological insights about the distribution of substituted nucleotides that the former does not.

Finally, we note that the GM  $P$ -value likely indicates overfitting, and yet the information criteria support Tam as a better fit to the *C. elegans* tRNA data.

### Assessing Preservation of Pairing Potential

As part of the model validation, we confirmed that the observed substitutions in the population follow a distribution consistent with the TAM bias expected by the

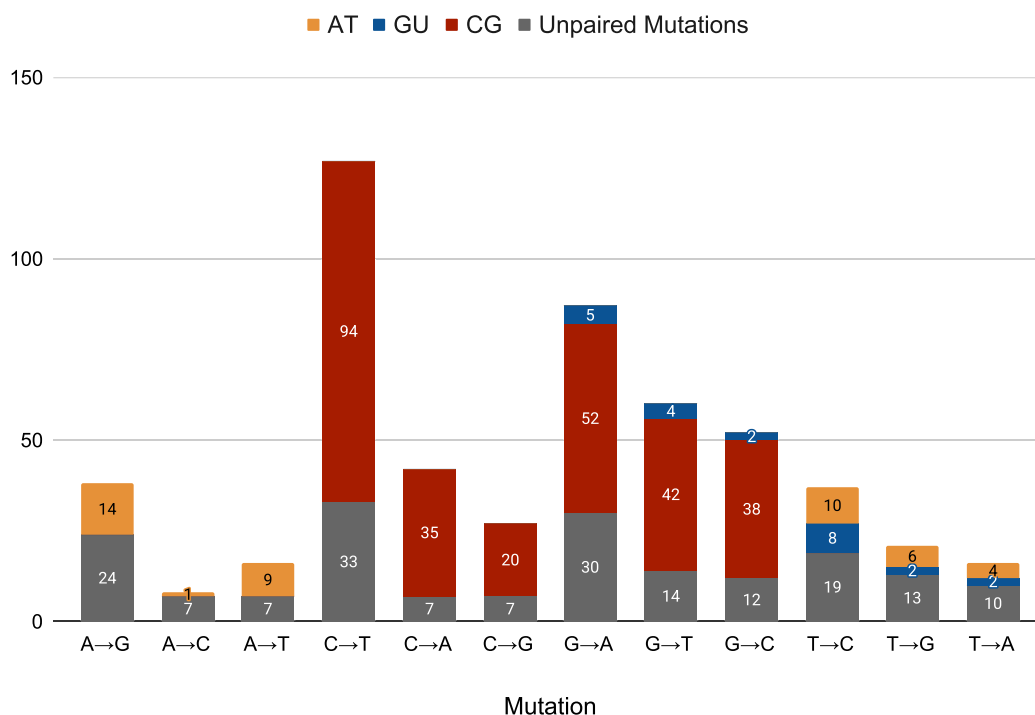
**Table 1.** Comparison of three expected pairing distributions under TAM bias model given initial ones from  $S$  and estimated  $\hat{\tau}$  with observed alleles encoded in  $S'$ .

	Nucleotides by pairedness				Pairedness	Base pairs
	$\Pi/\Pi^\tau/\hat{\Pi}^\tau$				$\Pi_{pu}/\Pi_{pu}^\tau/\hat{\Pi}_{pu}^\tau$	$\rho/\rho^\tau/\hat{\rho}^\tau$
Int	$\begin{pmatrix} 0.071 & 0.189 & 0.213 & 0.096 \\ 0.110 & 0.070 & 0.110 & 0.141 \end{pmatrix}$				(0.570, 0.430)	(0.662, 0.251, 0.087)
Exp	$\begin{pmatrix} 0.054 & 0.270 & 0.243 & 0.051 \\ 0.082 & 0.100 & 0.125 & 0.075 \end{pmatrix}$				(0.618, 0.382)	(0.785, 0.148, 0.067)
Obs	$\begin{pmatrix} 0.045 & 0.281 & 0.269 & 0.060 \\ 0.072 & 0.089 & 0.105 & 0.079 \end{pmatrix}$				(0.655, 0.345)	(0.807, 0.126, 0.066)
<i>P</i> -value	0.472				0.081	0.510

model. In particular, C's are overrepresented in  $\hat{\pi}^\tau$  by 43.0% while G's increased by 16.0% over the initial nucleotide distribution  $\pi$ . This demonstrates that the mutational variation in extant tRNAs exhibits a signature consistent with TAM, even though these sequences have also experienced strong selection that has removed many mutations over time.

We now assess whether three related distributions—nucleotides by pairedness ( $\Pi$ ), which separates  $\pi$  into paired and unpaired sites; pairedness ( $\Pi_{pu}$ ), the pairedness distribution achieved by column marginalization of  $\Pi$ ; and base pairs ( $\rho$ ), which indicates the proportion of GC, AU, and GU pairings in the secondary structure—also follow TAM-biased distributions. Details are given in Table 1. In all three cases, we fail to reject the null hypothesis for the expected and observed distributions of substituted nucleotides (ie paired sites also show signature of TAM).

In fact, as seen in the first and third columns, the proportion of C and G nucleotides in base pairs which are mutated is consistently *higher* than expected. That is, even though the TAM bias leads us to expect more C



**Fig. 3.** Number of observed substitutions, categorized by pairedness with pairings further subdivided by type. Substitutions are ordered first by the nucleotide that was substituted and then the resulting substitution. Bar segments are colored to show the base pair in the secondary structure when the substituted nucleotide was paired, and gray when the mutated nucleotide was unpaired.

and G substitutions, we observed even more than the model predicted. While this is not statistically significant, the trend, if representative, is surprising as it is inconsistent with expectations for the preservation of base pair formation.

One expectation for pairing preservation is that unpaired nucleotides should better tolerate mutations and therefore exhibit more substitutions in the dataset. This is clearly violated as paired positions are mutated almost twice as often (0.655 versus 0.345). Specifically, this represents an increase of 14.9% over the initial  $\Pi_{pu}$  one. As seen when  $\hat{\Pi}^r$  is compared to  $\Pi$ , this is entirely due to more substitutions in paired C's and G's which increased by 48.8% and 26.2%, respectively over the initial distribution. This is consistent with a TAM-biased substitution process. That is, paired nucleotides contain more C and G nucleotides, which makes these sites more likely to be targeted by TAM.

From the selective pressure to preserve secondary structure, one might then expect that mutated GC pairs are predominantly converted to GU pairings under the C → T transition as these would preserve structure in a wobble pair (Thornlow et al. 2018). As seen in Fig. 3, this is indeed the most common substitution. However, it happens only 33.5% of the time compared to the destabilizing G → A transition and the four

relevant transversions. Likewise, an AU pairing is potentially preserved as a GU one in only 31.8% of substitutions. However, these proportions are actually lower than when the corresponding unpaired substituted nucleotides combinations are considered: 47.1% for C,G and 45.3% for A,U. Hence, in this microevolution setting, we do not see strong support for selection favoring the formation of GU pairs from either GC or AU ones. Therefore, we see that C → T mutations are not preferentially biased to maintain structural conservation.

Similarly, we do not see enhanced conversion of GU pairings into GC/AU ones, as might also be expected. There are 974 wobble pairings in  $S$  (8.7% of 11,208), but GU is only 6.6% of the 348 pairings mutated to form  $S'$ . This is surprising as both transitions preserve pairing potential, and moreover would yield a base pair which is more thermodynamically stable.

Hence, we see a bias toward mutations in paired positions, particularly GC pairings, rather than away. Of the Watson–Crick pairings that are mutated, only 1/3 could form a wobble pairing. And there is no enhanced conversion of wobble pairings under mutation.

Thus, the observed substitutions follow the mutational pattern of TAM and remarkably, these are not selectively removed to preserve structural integrity. This is in line with the empirical estimates of fitness effects by

Li et al. (2016), which showed that for the single-copy of arginine-CCU tRNA in yeast, overall fitness decreased more at unpaired nucleotides than at paired ones.

### Comparison with Neutral Neighbors

When considering the distribution of alleles encoded in  $S'$ , we found that mutations arising from TAM are likely to substantially compromise tRNA fitness by destabilizing molecule structure. In particular, 42.7% of observed one-point mutations disrupt pairings. Yet, it is well-established that secondary structure conservation is a critical factor in tRNA macroevolution—so much so that MFE-preservation is considered a reliable proxy for fitness (Fontana and Schuster 1998; Reidys et al. 2001; Aguirre et al. 2011). Here, we show that even when our dataset is restricted to the neutral neighbors, clear deviations from macroevolution expectations are seen.

Neutrality under thermodynamic optimization is a strong constraint. Although 57.3% of substitutions from  $S$  to  $S'$  preserve pairing potential, only 37.3% are MFE-preserving. Interestingly, nearly all the difference is from the unpaired positions. The net result is that mutations still occur in paired nucleotides slightly more often than unpaired: 57.1% versus 42.9%. This is less than the substituted nucleotides for the full dataset (65.5% paired) but considerably more than those for the related neutral neighbors (32.6% paired). Thus, under this setting, all these observations agree with those in Assessing preservation of pairing potential section.

Moreover, under this scenario, we still find that GC pairings are mutated more than expected from the (initial) nucleotide distribution of MFE-preserving genes. This is again consistent with our TAM-bias model, but not with thermodynamic stability. In contrast, the initial distribution, which closely resembles the full one, favors GC over AU and GU pairings roughly 8:3:1. Hence, we see a bias toward the thermodynamic stability (with GC being by far the most stable pair, followed by AU) expected by macroevolution (Forster et al. 2006; Waldispühl et al. 2008) in the initial distribution, which is then being diminished by TAM-biased microevolution.

To quantify this bias further, recall that  $M$  denotes the MFE-preserving alleles from  $S'$  and  $\alpha(m)$  the corresponding gene in  $S$  for each  $m \in M$ . The set of all related neutral neighbors for  $m$  is  $\mathcal{N}_m$ , and  $\mathcal{N} = \bigcup_{m \in M} \mathcal{N}_m$ . For each  $m \in M$ , we now consider the relative thermodynamic stability  $f_m$  for  $m$  and  $\alpha(m)$  in comparison to  $N_m$ . As described in Materials and Methods section, the  $f_m$  value ranges from 0 to 1 with higher being less favorable.

As seen in Fig. 1 in Appendix, the values for the genes  $\alpha(m)$  are fairly narrowly clustered with a mean of 0.415 and standard deviation of 0.084. In contrast, those for

the alleles  $m$  are distributed over the full range with a (mean, SD) = (0.563, 0.258).

Hence, although GC pairs are clearly favored in the consensus genes, they are not the most thermodynamically stable in their neutral neighborhood. Interestingly, this may be indicative of mutational robustness (Gabzi et al. 2022) and/or evolvability (Wagner 2023). In contrast, the observed neutral neighbor alleles tend to be among the less energetically favorable possibilities, once again demonstrating a departure from the macroevolutionary expectation for favoring stability. Thus, even under this context, observed mutations are not preferentially biased toward maintaining structural stability; instead, structural stability is negatively affected.

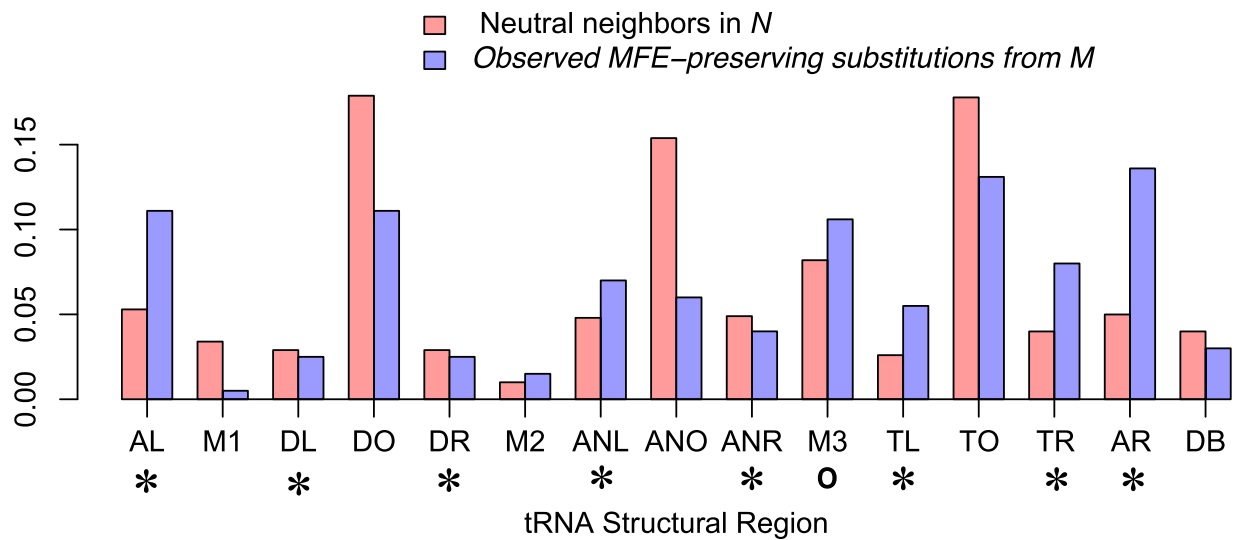
Finally, and to further reiterate such claim, we consider the positional distributions for  $M$  and  $\mathcal{N}$  across the different tRNA substructures. As illustrated in Fig. 4, the latter closely resembles findings for neutral neighbors from macroevolution studies (Fontana and Schuster 1998; Reidys et al. 2001). However, there are key differences for the former microevolution distribution of substituted nucleotides.

Observed MFE-preserving mutations occur more frequently in paired nucleotides, most notably in the acceptor stem (AR/AL) but also the 3' side of the T-arm (TR). The variable arm portion of the multibranch loop (M3) is also enhanced, but conclusions are confounded since it can contain paired as well as unpaired sites depending on the particular tRNA. In contrast, substitutions in all three hairpin loops (DO, ANO, and TO) are depressed, most notably for the AN.

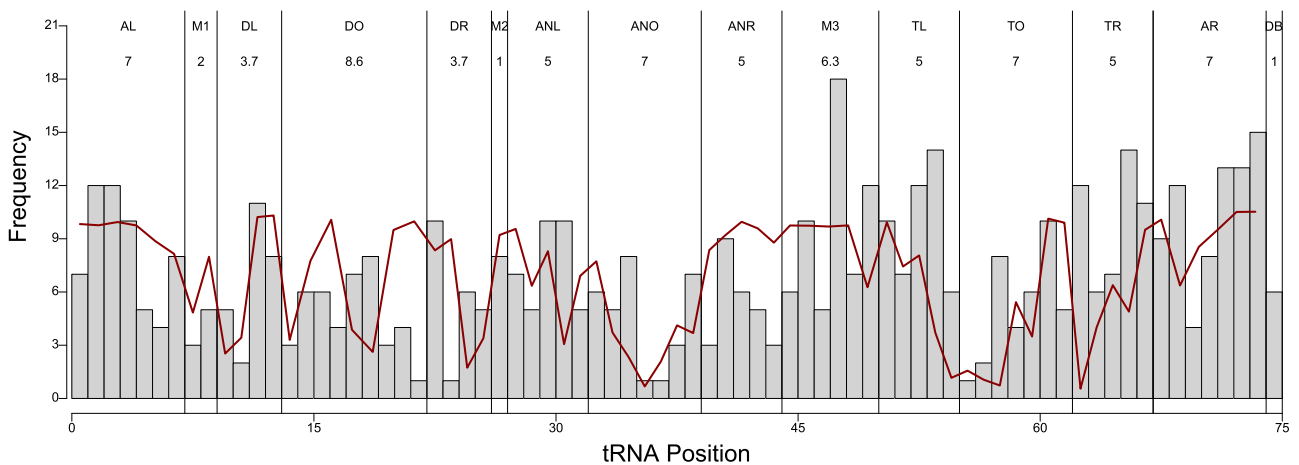
As a common identity element shared among most tRNAs, the three AN nucleotides are an essential component of a functional tRNA molecule, so it is hardly surprising if they are less mutable than might otherwise be expected under thermodynamic predictions. However, ANs are also known to evolve rapidly to meet new demands (Yona et al. 2013), and hence were included in our TAM-biased model. Although the AN itself represents <1/2 of its hairpin loop, it was not a priori clear that there would be a discernible difference at this level of structural resolution.

### Agreement with Experimental Fitness Assays

We have demonstrated that the distribution of alleles in *C. elegans* tRNA is described by our new TAM-biased Markov substitution model but deviates from the expectation that disruptions to structure are the least tolerated mutations, even when restricted to neutral neighbors. We now show that this distribution, and by extension our conclusions, accords with fitness tRNA studies (Yona et al. 2013; Li et al. 2016; Domingo et al. 2018).



**Fig. 4.** Distribution of observed MFE-preserving substitutions from  $M$  (blue, right bar per tRNA structural region) compared to all related neutral neighbors in  $N$  (red, left bar per tRNA structural region). Components of the tRNA secondary structure are labeled as in Fig. 1 with paired regions marked by \*. Note the enhancement of paired nucleotides that were substituted in  $M$ , particularly the acceptor stem, while mutations in hairpin loops, especially the AN one, are depressed. Conclusions about the variable arm (M3) are confounded since it can contain both kinds of sites depending on the specific tRNA.



**Fig. 5.** Comparison of substitution frequency over 531 *C. elegans* tRNA alleles (represented with gray bins) with mean site-fitness (Li et al. 2016) across all one-point mutants for the *S. cerevisiae* arginine-CCU tRNA (represented with the red line segments). Substructure components are labeled as in Fig. 1 along with their average length, which was used for the site normalization.

Recall that their experiments use the single-copy arginine-CCU gene from *S. cerevisiae*, so are insulated from the confounding effects when multiple copies of the same isoacceptor are present. Nonetheless, as seen in Fig. 5 and detailed below, we find good agreement between the site-fitness distribution for all one-point mutants obtained by Li et al. (2016) and the location of substitutions in *C. elegans* alleles.

To begin, we see that all *C. elegans* sites are mutated, but that lower frequency ones are often at (or near) a *S. cerevisiae* fitness minimum. Specifically, we observed

a point-wise correlation of around  $\sim 0.4$  when comparing the two datasets (with a  $P$ -value  $\sim 0.001$ , under the null hypothesis of no correlation). We note that, for many reasons, a perfect correlation is not expected as, for example, this analysis includes different tRNA genes which have distinct identity elements, and thus different constraints.

In particular, the lowest frequency/fitness sites are seen in hairpin loops, specifically the AN and the 5' end of the T-loop. The former agrees with both selection to preserve the AN as well as experimental evidence that

it can evolve rapidly when needed (Yona et al. 2013), while the latter is conserved as part of the internal promoter B-box region (Li et al. 2016). There are also proximal frequency/fitness lows in the D-stem which are not (yet) explained by functional conservation. However, this conservation may be driven by the D-stem's involvement in the A-box internal promoter essential for tRNA transcription (Mitra et al. 2015).

Conversely, regions with higher average frequency, ie both sides of the acceptor stem as well as the variable portion of the multiloop, also have high average fitness. This also echoes results from Domingo et al. (2018), whose dataset is built from 14 existing mutations of which 13 are from these regions.

Broadly, then, these two very different datasets agree; the acceptor stem and certain other paired sites are more mutable than expected under the (thermodynamic) neutral neighbor hypothesis, while the hairpin loops are much less so. These findings—aligned with our TAM-biased model—underscore the surprising robustness of many paired positions to mutations while also highlighting the influence of biochemical and other constraints beyond the secondary structure on the loop composition.

## Conclusions

It is well-established that tRNAs are subject to strong purifying selection that maintains their structural integrity over macroevolutionary time scales. This is most apparent in the conservation of secondary structure, as evidenced by covarying base pairs. Purifying selection is also evident at the microevolution time scale, as tRNA gene bodies exhibit substantially reduced rates of substitution compared to their flanking regions.

However, our results reveal that while the distribution of alleles in extant *C. elegans* strains is consistent with our model of TAM bias, it is not consistent with the expectations of thermodynamic neutrality for secondary structures. This shows that tRNAs in natural populations carry substantial levels of structure-destabilizing mutations, which may be tolerated but nevertheless likely induce meaningful fitness costs. This finding supports the conclusions in Thornlow et al. (2020), which claims that eukaryotic genomes carry substantial mutation load at tRNA genes.

This highlights a disconnect between scales. At the microevolutionary level, tolerable mutations often disrupt secondary structure, as demonstrated by empirical observations of fitness effects in yeast experiments. Yet, at the macroevolutionary level, selection pressure preserves canonical pairings. The resolution of this micro/macro dichotomy is likely to involve a variety of factors including the role of compensatory mutations (Kern and Kondrashov 2004;

Meer et al. 2010) and the interaction of isoacceptors (Yona et al. 2013) as well as the multicopy nature of most tRNA genes (Bloom-Ackermann et al. 2014).

Beyond this, our new population model may be applicable in other contexts. The key assumptions are high gene conservation and short evolutionary time to the most recent common ancestor. These enable the use of a consensus sequence as the ancestral approximation and the encoding of the allelic distribution in a single compound mutational variation “descendant.”

It may also be possible to integrate other mutational biases into a Markov model through careful rate matrix parameterization. This naturally leads to questions of model identifiability, whose resolution may reveal new biological insights.

## Materials and Methods

In this work, we explore how TAM's mutational patterns affect the secondary structure of tRNA molecules by analyzing tRNAs from a *C. elegans* population. To do this, we introduce a continuous Markov model that describes the distribution of alleles under mutational bias across a population. We then apply this new model to analyze the distribution of tRNA genes across extant *C. elegans* strains.

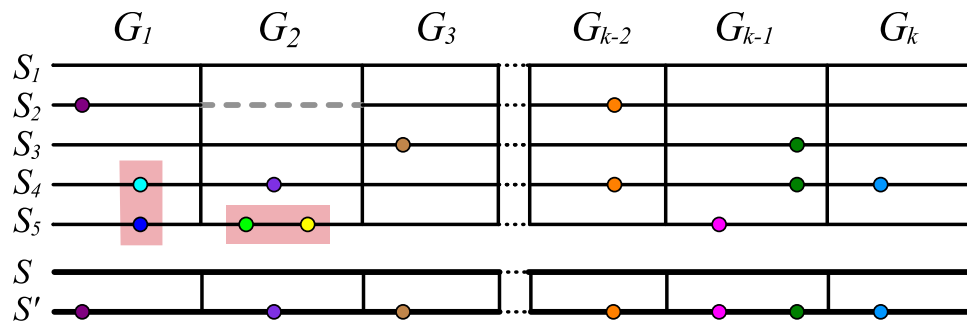
Typically (Felsenstein 2003; Yang 2006), Markov models are used for phylogenetic inference across species. This is effective since the populations have diverged enough that it is possible to reconstruct their evolutionary relationships in a gene tree. In contrast, we seek not to characterize the relationship among the *C. elegans* alleles but only to describe the distribution of mutations. In doing so, we show that they are consistent with our microevolution TAM-bias substitution model but not with macroevolution expectations of thermodynamic neutrality for secondary structure conservation.

### Model Development

Our model formulation is possible under a strong purifying selection and high mutation rate over a short evolutionary timescale—exactly the situation for tRNA microevolution. As we show below, these conditions allow us to capture nearly all of the observed allelic variation via a continuous Markov substitution process.

### Modeling Allele Formation within a Population

Let  $N$  be a population of  $n$  individuals and, for each individual  $i \in \{1, 2, \dots, n\}$ , let  $S_i$  be the concatenated collection of  $k$  genes under similar selection pressure, eg all tRNAs in extant *C. elegans* strains. Let  $\mathcal{S}$  be the corresponding sequence for the most recent common ancestor of  $N$ , and denote the consensus sequence for the  $S_i$  as  $S$ .



**Fig. 6.** Schematic of mutational variation sequence construction. From top to bottom: Each line  $S_i$  represents the  $k$  concatenated genes  $G_j$  for an individual  $i$  in the population  $N = \{1, \dots, 5\}$ . Every gene is not necessarily present in each individual; dashed gray lines denote a missing gene. Colored dots indicate nucleotide substitutions with respect to the consensus sequence  $S$  (shown below  $S_5$ ). Highlighted in red, we see that  $G_2$  in  $S_5$  violates condition (i) and  $G_1$  violates (ii) in  $S_4$  and  $S_5$ . By treating alleles violating these two conditions as missing genes, the remaining allelic variation can be encoded in the new compound sequence  $S'$  by substituting into  $S$  each distinct-site one-point mutation.

We assume that the  $k$  genes are highly conserved and the evolutionary time from  $S$  is relatively short. Thus, the sequences  $S_i$  are all close to each other, and also to  $S$ . In this scenario, although not necessarily others (Trudeau et al. 2016; Czech et al. 2018), it is reasonable to approximate  $S$  by  $S$ .

To formulate our microevolution population model, we impose the following two requirements. Moreover, we show in Results section that, for the tRNA microevolution situation being considered, these conditions are actually fairly mild.

- i. An allele differs from  $S$  by one nucleotide substitution.
- ii. Two alleles in distinct individuals cannot differ in the same site.

Alleles meeting (i) are, by definition, one-point mutants, and are often the individual steps in macroevolution models (Saks et al. 1998). We refer to those satisfying (ii) as *distinct-site* mutants. Under these two conditions, we can encode all mutational variation from the compliant alleles in a single compound sequence.

Specifically, we introduce the *mutational variation sequence*  $S'$  which is constructed by incorporating all observed substitutions in alleles meeting conditions (i) and (ii) into the consensus sequence  $S$ . As illustrated in Fig. 6, this guarantees a one-to-one correspondence between sites where  $S'$  differs from  $S$  and all distinct-site, one-point mutations observed in the population.

By integrating all those mutations into a single sequence, we can model allele formation without requiring a phylogenetic tree. This is particularly useful as inferring the true evolutionary relationship between individuals within a population can be very challenging (De Maio et al. 2015). Instead, we describe it simply through a nucleotide substitution process from  $S$  to  $S'$ .

Moreover, since multiple mutations within the same gene occur in different alleles, we strengthen the site independence for the Markov model. Hence, we parameterize our model with a  $4 \times 4$  rate matrix for nucleotide substitutions, rather than the more complicated dinucleotide ones (Allen and Whelan 2014) which model RNA pairing covariation for macroevolution analyses.

### Modeling the Distribution of Substituted Nucleotides under TAM

As described in the Background section, TAM's characteristic signature results in an asymmetric increase in  $C \rightarrow T$  and  $G \rightarrow A$  mutations on the coding strand. We hypothesize that the exchangeability asymmetry caused by TAM results in a bias of the distribution of substituted nucleotides, which we call the *TAM bias*.

To formulate this mathematically, let  $\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$  be the nucleotide distribution of the consensus sequence  $S$  (an approximation of  $S$ , the ancestral sequence). With no mutational bias, one would expect that the distribution of substituted nucleotides follow  $\pi$ .

To adjust for TAM, we introduce the *TAM bias vector*  $\tau = (\tau_R, \tau_Y) \in [0, 1]^2$ . We hypothesize that the nucleotides in  $S$  that mutated to form  $S'$  follow the *TAM bias distribution* given by:

$$\pi^\tau = (\pi_A^\tau, \pi_C^\tau, \pi_G^\tau, \pi_T^\tau), \quad (2)$$

where

$$\begin{aligned} \pi_A^\tau &= \pi_A - \tau_R, & \pi_C^\tau &= \pi_C + \tau_Y, & \pi_G^\tau &= \pi_G + \tau_R, & \text{and} \\ \pi_T^\tau &= \pi_T - \tau_Y. \end{aligned}$$

In other words, the TAM bias vector  $\tau$  adjusts the distribution of substituted nucleotides to account for differences in transition rates within purines ( $R$ ) and within

pyrimidines (Y). Hence,  $\pi_A^*$  is the proportion of A's in  $S$  that mutated to another nucleotide to form  $S'$  relative to all observed distinct-site, one-point substitutions.

### Modeling a TAM-biased Substitution Process

Our proposed TAM-biased distribution describes how the distribution of substituted nucleotides is altered from the root nucleotide distribution in  $S$ . However, to model allele formation within a population under TAM, we need to account for the mutation process itself. For this, we present a continuous Markov substitution model that incorporates the TAM bias to describe the production of  $S'$  from  $S$ .

Recall that for a Markov process over an edge of length  $t$ , the joint distribution of nucleotides in the ancestral and descendant sequences is given by  $P = \text{diag}(\pi) \exp(Qt)$ . In our context, the entry  $p_{ij}$  of  $P$  represents the probability of a site in  $S'$  differing from  $S$  by a substitution from  $i$  to  $j$ . To parameterize this new model  $P$ , we must propose a TAM-biased instantaneous rate matrix  $Q$ . Recall also that, typically, one entry of the matrix is set to 1 to avoid scaling redundancy with the edge length. We note that the rate matrices presented here are not-normalized form for simplicity, but all computations are done under normalization as described in the Background section.

We begin with a strand-symmetric model (sym) which assumes that a substitution is indistinguishable from its Watson–Crick complement. This accounts for the background noise, and its matrix has five free parameters with:

$$Q_{sym} = \begin{array}{c|cccc} & \text{A} & \text{C} & \text{G} & \text{T} \\ \hline \text{A} & * & \alpha & \beta & \gamma \\ \text{C} & \delta & * & \epsilon & \lambda \\ \text{G} & \lambda & \epsilon & * & \delta \\ \text{T} & \gamma & \beta & \alpha & * \end{array}$$

where  $\alpha, \beta, \gamma, \epsilon, \lambda \geq 0$ , and  $\delta = 1$ . Note that by accounting for the edge length, the model has 6 free parameters.

However, a characteristic of TAM is the asymmetry in transitions observed in the coding strand. Hence, we break the strand symmetry between substitutions  $A \rightarrow G$  with  $T \rightarrow C$  (circled in  $Q_{sym}$ ), and  $G \rightarrow A$  with  $C \rightarrow T$  (squared in  $Q_{sym}$ ), yielding a transition-asymmetric rate matrix

$$Q_{tran} = \begin{array}{c|cccc} & \text{A} & \text{C} & \text{G} & \text{T} \\ \hline \text{A} & * & \alpha & \beta & \gamma \\ \text{C} & \delta & * & \epsilon & \lambda \\ \text{G} & \theta & \epsilon & * & \delta \\ \text{T} & \gamma & \mu & \alpha & * \end{array}$$

with seven free parameters. Hence,  $Q_{sym}$  is a special case of  $Q_{tran}$  with  $\theta = \lambda$  and  $\mu = \beta$ , and this model has eight free parameters when accounting for the edge length.

Next, we incorporate the TAM bias model for the distribution of substituted nucleotides by adding restrictions to the  $Q_{tran}$  rate matrix entries. More specifically, we require that

$$\pi_A - \tau_R = \frac{s(1)}{Z}, \quad \pi_C + \tau_Y = \frac{s(2)}{Z}, \quad \pi_G + \tau_R = \frac{s(3)}{Z}, \quad (5)$$

and  $\pi_T - \tau_Y = \frac{s(4)}{Z},$

where  $s(i) = \sum_{j \neq i} p_{ij}$  is the sum of the off-diagonal entries in each row of  $P$  with  $Z = \sum_{i \neq j} p_{ij}$  as the normalizing factor. Thus, the left-hand sides come from Equation (2), whereas the right ones follow from the continuous Markov model.

We remark that these equations do not increase the number of model parameters. The entries of  $\pi$  are needed to compute  $P$  for every rate matrix, and hence do not count against any particular model. Additionally, the TAM bias vector  $\tau = (\tau_R, \tau_Y)$  is estimated from the flanking regions, which are subject to TAM but not strong selection pressure. Hence, since  $\tau_R$  and  $\tau_Y$  are independent from the tRNA alleles, they are not counted as model parameters either.

As a consequence, we improve the model parameterization under the constraints in Equation (5) by identifying linear dependencies among the  $Q_{tran}$  rate matrix entries using a computer algebra package. More precisely, given the high conservation of tRNAs, we expect a very slow substitution process, implying that the entries of  $Qt$  are quite small. Thus, we approximate  $\exp(Qt) \approx I + Qt$ , where  $I$  is the  $4 \times 4$  identity matrix. Under this approximation, we represent Equation (5) as a homogeneous linear system which can be solved using Macaulay2 (Grayson and Stillman 2009).

As a result, we obtain the linear dependencies given below for the transversions  $\delta$  and  $\gamma$  as well as the transition  $\theta$ . We then define  $Q_{TAM}$  to be the rate matrix obtained from  $Q_{tran}$  using these relationships. This yields a total of four free parameters for the new TAM-biased rate matrix: the

three transitions  $\beta$ ,  $\lambda$ , and  $\mu$  along with the two transversions  $\alpha$  and  $\epsilon$ —yielding five free parameters for the model.

We note that the choice of dependent variables among the matrix parameters was determined by Macaulay2 (version 1.25.06); any alternate formulation of the relationships would yield an equivalent model parameterization although a different form for the  $Q_{\text{tam}}$  rate matrix. The linear dependencies used here (which were manually derived from Macaulay2's initial output shown in the Appendix) are:

$$\begin{aligned}\delta &= -\epsilon - \lambda + (\mu - \beta)\Omega \\ \theta &= \lambda + (\mu - \beta)\Theta \\ \gamma &= -\alpha + \mu\Gamma + \beta\Psi,\end{aligned}$$

where

$$\begin{aligned}\xi &= \Lambda_4 - \Lambda_1 \quad \text{and} \quad \Omega = \frac{-1 - \Lambda_2}{\xi}, \quad \Theta = \frac{\Lambda_2 - \Lambda_3}{\xi}, \\ \Gamma &= \frac{\Lambda_1 - 1}{\xi}, \quad \Psi = \frac{1 - \Lambda_4}{\xi}\end{aligned}$$

for

$$\Lambda_i = \frac{\tau(i)}{\pi_i} \quad \text{with} \quad \tau(i) = \begin{cases} \tau_R & \text{if } i \text{ is odd} \\ \tau_Y & \text{if } i \text{ is even} \end{cases}.$$

Observe that the coefficients  $\Omega$ ,  $\Theta$ ,  $\Gamma$ , and  $\Psi$  are real numbers defined solely by the TAM bias vector  $\tau$  and the root nucleotide distribution  $\pi$ . As a function of  $\pi$  and  $\tau$ , they are unbounded and defined almost everywhere (except for a set of measure zero, ie a negligible subset of the parameter space).

We remark that  $\Omega \neq 0$ . Hence, the inequality  $\mu \neq \beta$  must hold; otherwise  $\delta < 0$ , which violates the definition. Likewise  $\Theta \neq 0$  for almost all  $\pi$  and  $\tau$ . But since we have already concluded that  $\mu \neq \beta$  when  $\Omega$  is defined, then generically  $\theta \neq \lambda$  as well. This shows that satisfying Equation (5) requires breaking the transition strand symmetries.

We show in Results sections that among the five rate matrices considered—strand-symmetric (sym), transition-asymmetric (tran), and TAM biased (Tam) along with general Markov (GM) and GTR—the new TAM-biased Markov substitution model best fits the allelic distribution in our *C. elegans* tRNA dataset.

## Data Analysis

Our dataset was generated from two publicly available databases: the Genomic tRNA Database (GtRNAdb, v2.0) (Chan and Lowe 2008, 2015) and the *C. elegans* Natural Diversity Resource (CeNDR, release 20180527) (Cook et al. 2016; Crombie et al. 2023). The starting point

was the reference genome (release WBcel235/ce11) (Harris et al. 2019) for the universal laboratory strain N2.

GtRNAdb identifies 581 functional tRNAs in the N2 nuclear genome, and CeNDR provided variant data for 330 *C. elegans* wild isolates collected from around the globe. We then retrieved 677 single nucleotide variants (SNVs) for the tRNA genes from the hard-filter variant call format (vcf) file using custom Python scripts.

Secondary structures for all tRNA sequences were obtained from tRNAscan-SE (version 2.0.11 Chan and Lowe 2019) using the default parameters. As illustrated in Fig. 1, nucleotides were classified as “paired” or “unpaired” according to these structures. We note that, although pictured, the Tail was not included in the analysis.

There were 16 SNVs which were found by tRNAscan-SE to occur in an intron; these were excluded from further analysis. Hence, our *C. elegans* microevolution dataset consists of 661 variants for 581 tRNA genes across a population of 331 strains.

We also collected data for the 40 nucleotides upstream and 40 downstream of each SNV and N2 gene. Since these tRNA flanking regions are unwound during transcription, they are subject to TAM (Thornlow et al. 2018). However, they are not under strong selection pressure, making them ideal for estimating the TAM bias vector  $\tau = (\tau_R, \tau_Y)$ . Under our TAM bias hypothesis, subtracting the observed distribution substituted nucleotides from the initial nucleotide distribution should yield  $\pi - \pi^\tau = (\tau_R, -\tau_Y, -\tau_R, \tau_Y)$ . As reported in Results section, the positive and negative values are in good agreement, yielding a robust estimate for  $\tau$ .

## Testing TAM-biased Distributions

Our model for TAM bias from Equation (2) proposes that the substituted nucleotides are distributed as  $\pi^\tau$ , where  $\pi_A^\tau = \pi_A - \tau_R$ ,  $\pi_C^\tau = \pi_C + \tau_Y$ , etc. To validate this hypothesis for the TAM-biased rate matrix, we compute a *P*-value using the `xmonte` function from the XNomial R package (Engels 2015). This function assesses whether the observed substitutions conform to the estimated multinomial distribution. Specifically, it examines 3,000 random outcomes to estimate the probability that a sample deviates from the expected distribution by at least as much as the observed data.

We later apply the same approach to test for the preservation of pairing potential. To do this, we classify the nucleotides of each tRNA gene as either paired (*p*) or unpaired (*u*) according to the tRNAscan-SE secondary structure. This separates  $\pi$ , the nucleotide distribution for *S*, by pairedness into

$$\Pi = \begin{pmatrix} A_p & C_p & G_p & T_p \\ A_u & C_u & G_u & T_u \end{pmatrix}, \quad (6)$$

where  $X_p$  and  $X_u$  denote the proportion of paired and unpaired nucleotides respectively for  $X \in \{A, C, G, T\}$ . Marginalizing the rows of  $\Pi$  reconstitutes  $\pi$ , while column marginalization yields the pairedness distribution, denoted  $\Pi_{pu} = (\Pi_p, \Pi_u)$ .

We also consider the distribution of base pairs, denoted  $\rho = (\rho_{GC}, \rho_{AT}, \rho_{GT})$ , where  $\rho_{GT}$  is the proportion of GU pairings in the tRNA secondary structures. It follows that

$$\rho_{GC} = \frac{C_p}{G_p + A_p}, \quad \rho_{AT} = \frac{A_p}{G_p + A_p}, \quad \text{and}$$

$$\rho_{GT} = \frac{G_p - C_p}{G_p + A_p} = \frac{T_p - A_p}{G_p + A_p}.$$

Under our TAM bias hypothesis, the distribution of substituted nucleotides by pairedness would be

$$\Pi^\tau = \begin{pmatrix} A_p^\tau & C_p^\tau & G_p^\tau & T_p^\tau \\ A_u^\tau & C_u^\tau & G_u^\tau & T_u^\tau \end{pmatrix}$$

$$= \begin{pmatrix} \pi_A^\tau \left(\frac{A_p}{\pi_A}\right) & \pi_C^\tau \left(\frac{C_p}{\pi_C}\right) & \pi_G^\tau \left(\frac{G_p}{\pi_G}\right) & \pi_T^\tau \left(\frac{T_p}{\pi_T}\right) \\ \pi_A^\tau \left(\frac{A_u}{\pi_A}\right) & \pi_C^\tau \left(\frac{C_u}{\pi_C}\right) & \pi_G^\tau \left(\frac{G_u}{\pi_G}\right) & \pi_T^\tau \left(\frac{T_u}{\pi_T}\right) \end{pmatrix}, \quad (7)$$

where, for instance,  $A_p^\tau$  is the proportion of one-point mutants in  $S'$  that differ from  $S$  by a substitution which changed a paired A. Row and column marginalization of  $\Pi^\tau$  yield  $\pi^\tau$  and  $\Pi_{pu}^\tau = (\Pi_p^\tau, \Pi_u^\tau)$ , respectively.

Analogously,  $\rho_{mGT}^\tau$  represents the proportion of substitutions that changed either the G or the T in what would have been a wobble GU pairing. Consequently, the components of  $\rho^\tau$  are

$$\rho_{GC}^\tau = C_p^\tau + G_p^\tau \left(\frac{\rho_{GC}}{\rho_{GC} + \rho_{GT}}\right), \quad \rho_{AT}^\tau = A_p^\tau + T_p^\tau \left(\frac{\rho_{AT}}{\rho_{AT} + \rho_{GT}}\right),$$

and

$$\rho_{GT}^\tau = G_p^\tau \left(\frac{\rho_{GT}}{\rho_{GC} + \rho_{GT}}\right) + T_p^\tau \left(\frac{\rho_{GT}}{\rho_{AT} + \rho_{GT}}\right). \quad (8)$$

### Assessing Markov Substitution Rate Matrices

To compare the five possible exchangeability matrices considered, we computed for each the three most common information criteria for model selection (Kalyaanamoorthy et al. 2017): the Akaike (AIC) (Akaike 1974), the corrected Akaike (AICc) (Burnham and Anderson 2002), and the Bayesian (BIC) (Schwarz 1978). We found that the differences between the first two were always negligible due to the length of the consensus sequence relative to the number of mutations. Hence, we focus on the AIC and BIC.

Loosely speaking, the lower the estimate, the “better” the model. Each criterion depends on the log-likelihood (LH) of the model being “penalized” by the number of parameters. The LH function is given by

$$L(Q, t) = \sum_{i=1}^n \log P(x_i | \pi, Q, t),$$

where  $\pi$  is the nucleotide distribution of  $S$ ,  $n$  is the sequence length of  $S$ ,  $x_i$  is the site pattern of  $S$  and  $S'$  at site  $i$ ,  $Q$  is the rate matrix,  $t$  are the expected number of substitutions between  $S$  and  $S'$ , and  $P(x_i | \pi, Q, t)$  is the conditional probability of  $x_i$  given  $\pi$ ,  $Q$ , and  $t$ .

To compute the LH function and estimate the maximum likelihood estimator (MLE) for all parameters, we again use the approximation  $\exp(Qt) \approx I + Qt$ . We then used MATLAB’s (version 25.2, R2025b MathWorks 2024) `fmincon` function to estimate  $Q$  and  $t$  MLEs using 1000 different random starting points.

We also computed  $\hat{P} = \text{diag}(\pi) \exp(\hat{Q}t)$  using the MLE estimates for each parameterization and compared the results to the data. The corresponding  $P$ -values for such comparisons were then obtained using the `xmonte` function.

### Analyzing Thermodynamic Neutrality

We also compared observed distributions against a background of neutral neighbors. Recall that these are alleles which are both one-point mutants and MFE-preserving. For consistency, we consider only observed alleles which are also neutral neighbors. The necessary thermodynamic predictions, ie the MFE score and its associated secondary structure, are computed with `RNAfold` (version 2.6.2 Lorenz et al. 2011) under default parameters.

More precisely, let  $P$  be the set of all distinct-site, one-point mutant alleles used to generate  $S'$ , and  $R$  the set of tRNA genes in  $S$ . For  $p \in P$ , let  $\alpha(p)$  be its corresponding  $r \in R$ . Then  $p \in P$  is MFE-preserving, and hence is a neutral neighbor, when it has the same `RNAfold` predicted secondary structure as  $\alpha(p)$ .

Let  $M \subseteq P$  be MFE-preserving, ie all observed neutral neighbors. For  $m \in M$ , let  $N_m$  be the set of all neutral neighbors of  $\alpha(m)$ , and  $\mathcal{N} = \bigcup_{m \in M} N_m$ . Since  $m \in N_m$ ,  $M$  is a subset of  $\mathcal{N}$ . We determine  $M$  and  $\mathcal{N}$  computationally, and consider how they differ.

First, we consider whether the substituted nucleotides were paired or not, according to the classification from the `tRNAscan-SE` secondary structure for  $\alpha(m)$ . Next we consider the thermodynamic stability of  $m$  and  $\alpha(m)$  relative to  $N_m$  over all of  $M$ .

This is done by ranking the set  $N_m \cup \{\alpha(m)\}$  by increasing MFE score with the lowest ranked as

1. The ranking is dense (García-Lapresta and Martínez-Panero 2024); ties share the same rank, and the next rank is assigned the next integer. For instance, if three sequences tie for the 5th lowest MFE score, all are assigned rank 5 with the next one(s) ranked 6. Rankings are then normalized by dividing by the total number of ranks for  $\alpha(m)$ . This defines a function  $f_m : N_m \cup \{\alpha(m)\} \rightarrow (0, 1]$  for each  $m \in M$ , where values near 0 are the most thermodynamically stable while those near 1 are the least so. We will compare  $f_m(m)$  and  $f_m(\alpha(m))$  over all of  $M$ .

Finally, we consider the positional distribution of substitutions in  $M$  and  $\mathcal{N}$ . The tRNA structural regions are defined in Fig. 1 and determined for  $\alpha(m)$  based on the tRNAscan-SE secondary structure.

### Comparing with Experimental Fitness

We consider the mutability of each site across the 531 alleles in our *C. elegans* dataset. Since not all tRNAs have the same length, the sites in each allele were normalized to the average length of the corresponding substructure component from the 15 (excluding the Tail) listed in the table in Fig. 1. The number of mutations for each normalized site can then be presented in a frequency histogram.

We compared this distribution to the mean site-fitness scores obtained by Li et al. (2016) for the *S. cerevisiae* single-copy arginine-CCU tRNA. The sites in the tRNA<sub>Arg</sub><sup>CCU</sup> were also normalized as above before reporting the average of their 3 one-point mutant fitness scores. For display purposes, this mean value was linearly rescaled from (0.685, 1.017) to (1, 10). Note that normalization of the data was performed only for the comparison against the experimental fitness data; for all other results, the data were not normalized. A correlation between the *C. elegans* dataset and the data analyzed in this work, along with their corresponding *P*-value, were computed using R's `cor.test()` function under Pearson's product-moment correlation.

### Notes

- 1 Uracil (U) will be used when discussing RNA structures and thymine (T) for their DNA genes.
- 2 Sequence pictured is *Bacillus subtilis* tRNA<sup>Asp</sup> with structures adapted with permission from Liu (2013).
- 3 Although there are 20 standard amino acids and 64 distinct codons, most eukaryotic nuclear genomes encode hundreds of tRNAs.
- 4 All data were analyzed at 10-digit precision. Due to rounding, reported values may not add up exactly.

### Acknowledgments

The authors used a large language model solely for spelling and grammar corrections; it was not used to generate original content.

### Funding

This work was supported by the NSF-Simons Southeast Center for Mathematics and Biology (SCMB) through the grants National Science Foundation DMS-1764406 and Simons Foundation/SFARI 594594. H.B. was also partially supported by National Science Foundation grant DMS-2331660. A.P. was also partially supported by National Science Foundation grant IOS-2319796. C.H. was also partially supported by the National Institutes of Health grant R01GM126554.

### Data Availability

The data used in this work were obtained from and are available at the Genomic tRNA Database (GtRNAdb, v2.0) (Chan and Lowe 2008, 2015) and the *C. elegans* Natural Diversity Resource (CeNDR, release 20180527) (Cook et al. 2016; Crombie et al. 2023). The repository <https://github.com/HectorBanos/TAM> contains the concatenated gene tRNA sequence  $S$  (composed of all 531 tRNAs), and the population variation sequence  $S'$ , as well as code to reproduce the results in this article.

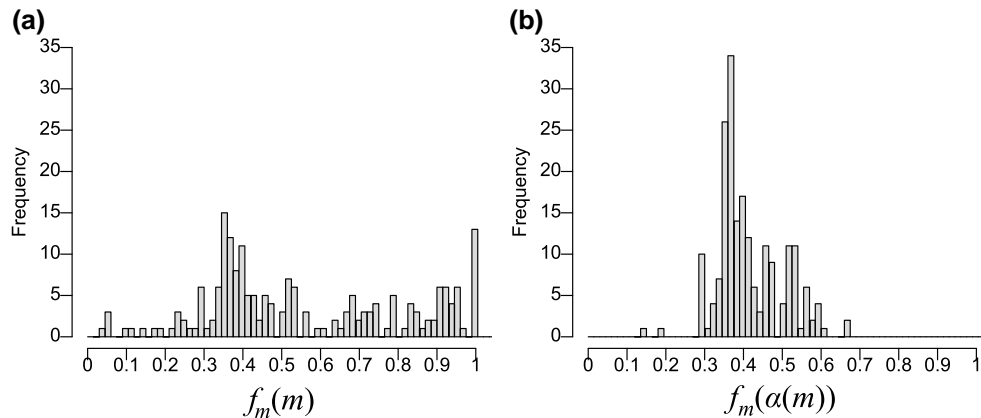
## Appendix

The output obtained from Macaulay2 by solving the linear system described in Modeling a TAM-biased Substitution Process section.

$$\lambda = \frac{1}{\pi_T \pi_G \pi_C \pi_A} ((\theta - \epsilon + \alpha) \pi_G + \alpha \tau_R) \pi_C^4 + ((3\theta - 2\epsilon + 2\alpha + \beta) \pi_G^2 + ((3\theta + \gamma - 2\epsilon + 3\alpha) \pi_T + (\theta - \epsilon + 2\alpha + \beta) \tau_R - 3\theta + 2\epsilon - 2\alpha) \pi_G + \tau_R (\gamma + 2\alpha) \pi_T - 2\alpha) \pi_C^3 + ((3\theta - \epsilon + \alpha + 2\beta) \pi_G^2 + ((6\theta + 2\gamma - 3\epsilon + 3\alpha + 3\beta) \pi_T + (2\theta - \epsilon + \alpha + 2\beta) \tau_R - 6\theta + 2\epsilon - 2\alpha - 2\beta) \pi_G^2 + ((3\theta + 3\gamma - \epsilon + 2\alpha) \pi_T^2 + ((2\theta + 2\gamma - \epsilon + 2\alpha + 2\beta) \tau_R + 2\tau_R (\pi_T - 1) (\gamma + \frac{\alpha}{2}) \pi_T - \frac{\alpha}{2})) \pi_C^2 + ((\theta + \beta) \pi_G^4 + ((3\theta + \gamma - \epsilon + 3\beta) \pi_T + (\theta + \beta) \tau_R - 3\theta - 2\beta) \pi_G^3 + ((3\theta + 3\gamma - \epsilon + 2\beta) \pi_T^2 + ((2\theta + \gamma - \epsilon + 2\beta) \tau_R - 6\theta - 2\gamma + \alpha \tau_Y + (\tau_Y - 3) \beta + \epsilon) \pi_T + (-2\theta - 2\beta) \tau_R + 3\theta + \beta) \pi_G^2 + ((\theta + 2\gamma) \pi_T^3 + ((\theta + 2\gamma + \beta) \tau_R - 3\theta + (\tau_Y - 3) \gamma + \alpha \tau_Y) \pi_T^2 + ((-2\theta - 2\gamma - 2\beta) \tau_R - \alpha \tau_Y + 3\theta + \gamma) \pi_T + (\theta + \beta) \tau_R - \theta) \pi_G + \pi_T \gamma \tau_R (\pi_T - 1)^2) \pi_C + \pi_T \pi_G \tau_Y (\pi_T + \pi_G - 1) (\gamma \pi_T + \beta \pi_G)$$

$$\mu = \frac{-1}{\pi_G \pi_T \pi_A} (2\pi_T^3 \gamma + ((3\gamma - \alpha + \beta) \pi_G + (3\pi_C + \tau_R - \tau_Y - 3) \gamma + \alpha \pi_C) \pi_T^2 + ((\gamma - \alpha + \beta) \pi_G^2 + ((2\pi_C + \tau_R - \tau_Y - 2) \gamma + \beta \pi_C - \beta \tau_Y + (1 - \tau_R) \alpha - \beta) \pi_G + (\pi_C - 1) (\pi_C + \tau_R - \tau_Y - 1) \gamma + \alpha \pi_C (\pi_C - \tau_Y - 1)) \pi_T - \tau_Y (\pi_C + \pi_G - 1) (\alpha \pi_C + \beta \pi_G))$$

$$\delta = \frac{-1}{\pi_G \pi_T \pi_A} ((\theta + \beta) \pi_G^3 + ((2\theta - \epsilon + \alpha + \beta) \pi_C + (2\theta + \gamma + \beta) \pi_T + (\theta + \beta) \tau_R - 2\theta - \beta) \pi_G^2 + ((\theta - \epsilon + \alpha) \pi_C^2 + ((2\theta + \gamma - \epsilon + \alpha) \pi_T + (\theta - \epsilon + \alpha + \beta) \tau_R - 2\theta + \epsilon - \alpha) \pi_C + (\theta + \gamma) \pi_T^2 + ((\theta + \gamma + \beta) \tau_R - 2\theta - \gamma) \pi_T + (-\theta - \beta) \tau_R + \theta) \pi_G + \tau_R (\pi_C + \pi_T - 1) (\gamma \pi_T + \alpha \pi_C))$$



**Fig. A1.** Relative thermodynamic stability  $f_m$  for the 198 observed neutral neighbors (left) and corresponding tRNA gene (right). Approximately half of the observed substitutions in the *C. elegans* population are disproportionately less stable than other related neutral neighbors, but no such skew is seen for the consensus tRNA genes approximating the ancestral sequence.

The dependencies of the  $Q_{TAM}$  matrix entries shown in Modeling a TAM-biased Substitution Process section where derived manually from these.

## Literature cited

Adrion JR, White PS, Montooth KL. The roles of compensatory evolution and constraint in aminoacyl tRNA synthetase evolution.

- Mol Biol Evol. 2015;33:152–161. <https://doi.org/10.1093/molbev/msv206>.
- Aguirre J, Buldú JM, Stich M, Manrubia SC. Topological structure of the space of phenotypes: the case of RNA neutral networks. *PLoS One*. 2011;6:e26324. <https://doi.org/10.1371/journal.pone.0026324>.
- Akaike H. A new look at the statistical model identification. *IEEE Trans Automat Contr*. 1974;19:716–723. <https://doi.org/10.1109/TAC.1974.1100705>.
- Allen JE, Whelan S. Assessing the state of substitution models describing noncoding RNA evolution. *Genome Biol Evol*. 2014;6:65–75. <https://doi.org/10.1093/gbe/evt206>.
- Berg MD, Brandl CJ. Transfer RNAs: diversity in form and function. *RNA Biol*. 2021;18:316–339. PMID: 32900285. <https://doi.org/10.1080/15476286.2020.1809197>.
- Bloom-Ackermann Z et al. A comprehensive tRNA deletion library unravels the genetic architecture of the tRNA pool. *PLoS Genet*. 2014;10:1–16. <https://doi.org/10.1371/journal.pgen.1004084>.
- Boivin V et al. Simultaneous sequencing of coding and noncoding RNA reveals a human transcriptome dominated by a small number of highly expressed noncoding genes. *RNA*. 2018;24:950–965. Epub 2018 Apr 27. <https://doi.org/10.1261/ma.064493.117>.
- Burnham K, Anderson D. Model selection and multimodel inference: a practical information-theoretic approach. Springer; 2002.
- Cannone JJ et al. The comparative RNA web (CRW) site: an online database of comparative sequence and structure information for ribosomal, intron, and other RNAs. *BMC Bioinformatics*. 2002;3:2. <https://doi.org/10.1186/1471-2105-3-2>.
- Casanellas M, Sullivant S. 16 - The Strand Symmetric Model. In: Edited by Pachter L, Sturmfels B, editors. *Algebraic statistics for computational biology*. Cambridge University Press; 2005. p. 305–321.
- Chan PP, Lowe TM. GtRNAdb: a database of transfer RNA genes detected in genomic sequence. *Nucleic Acids Res*. 2008;37:D93–D97. <https://doi.org/10.1093/nar/gkn787>.
- Chan PP, Lowe TM. GtRNAdb 2.0: an expanded database of transfer RNA genes identified in complete and draft genomes. *Nucleic Acids Res*. 2015;44:D184–D189. <https://doi.org/10.1093/nar/gkv1309>.
- Chan PP, Lowe TM. tRNAscan-SE: searching for tRNA genes in genomic sequences. In: *Gene prediction: methods and protocols*. Springer New York; 2019. p. 1–14.
- Cook DE, Zdraljevic S, Roberts JP, Andersen EC. CeNDR, the *Caenorhabditis elegans* natural diversity resource. *Nucleic Acids Res*. 2016;45:D650–D657. <https://doi.org/10.1093/nar/gkw893>.
- Cowperthwaite MC, Meyers LA. How mutational networks shape evolution: lessons from RNA models. *Annu Rev Ecol Evol Syst*. 2007;38:203–230. <https://doi.org/10.1146/ecolsys.2007.38.issue-1>.
- Crombie TA et al. CeNDR, the *Caenorhabditis* natural diversity resource. *Nucleic Acids Res*. 2023;52:D850–D858. <https://doi.org/10.1093/nar/gkad887>.
- Czech L, Barbera P, Stamatakis A. Methods for automatic reference trees and multilevel phylogenetic placement. *Bioinformatics*. 2018;35:1151–1158. <https://doi.org/10.1093/bioinformatics/bty767>.
- De Maio N, Schrepf D, Kosiol C. PoMo: an allele frequency-based approach for species tree estimation. *Syst Biol*. 2015;64:1018–1031. <https://doi.org/10.1093/sysbio/syv048>.
- Domingo J, Diss G, Lehner B. Pairwise and higher-order genetic interactions during the evolution of a tRNA. *Nature*. 2018;558:117–121. <https://doi.org/10.1038/s41586-018-0170-7>.
- Eddy SR, Durbin R. RNA sequence analysis using covariance models. *Nucleic Acids Res*. 1994;22:2079–2088. <https://doi.org/10.1093/nar/22.11.2079>.
- Engels B. *XNomial: Exact Goodness-of-Fit Test for Multinomial Data with Fixed Probabilities*. R package version 1.0.4, 2015.
- Felsenstein J. *Inferring phylogenies*. Sinauer Associates; 2003.
- Fontana W, Schuster P. Shaping space: the possible and the attainable in RNA genotype-phenotype mapping. *J Theor Biol*. 1998;194:491–515. <https://doi.org/10.1006/jtbi.1998.0771>.
- Forster R, Adami C, Wilke CO. Selection for mutational robustness in finite populations. *J Theor Biol*. 2006;243:181–190. <https://doi.org/10.1016/j.jtbi.2006.06.020>.
- Gabzi T, Pilpel Y, Friedlander T. Fitness landscape analysis of a tRNA gene reveals that the wild type allele is sub-optimal, yet mutationally robust. *Mol Biol Evol*. 2022;39:msac178. <https://doi.org/10.1093/molbev/msac178>.
- Gaillard H, Aguilera A. Transcription as a threat to genome integrity. *Annu Rev Biochem*. 2016;85:291–317. <https://doi.org/10.1146/annurev-biochem-060815-014908>.
- García-Lapresta JL, Martínez-Panero M. Two characterizations of the dense rank. *J Math Econ*. 2024;111:102963. <https://doi.org/10.1016/j.jmateco.2024.102963>.
- Giegé R, Eriani G. The tRNA identity landscape for aminoacylation and beyond. *Nucleic Acids Res*. 2023;51:1528–1570. <https://doi.org/10.1093/nar/gkad007>.
- Gómez-González B, Aguilera A. Activation-induced cytidine deaminase action is strongly stimulated by mutations of the THO complex. *Proc Natl Acad Sci U S A*. 2007;104:8409–8414. <https://doi.org/10.1073/pnas.0702836104>.
- Grayson DR, Stillman ME. Macaulay2, a software system for research in algebraic geometry. 2009. Available at <http://www2.macaulay2.com>.
- Green P et al. Transcription-associated mutational asymmetry in mammalian evolution. *Nat Genet*. 2003;33:514–517. <https://doi.org/10.1038/ng1103>.
- Griffiths-Jones S et al. Rfam: annotating non-coding RNAs in complete genomes. *Nucleic Acids Res*. 2005;33:D121–D124. <https://doi.org/10.1093/nar/gki081>.
- Harris TW et al. WormBase: a modern model organism information resource. *Nucleic Acids Res*. 2019;48:D762–D767. <https://doi.org/10.1093/nar/gkz920>.
- Ishimura R et al. Ribosome stalling induced by mutation of a CNS-specific tRNA causes neurodegeneration. *Science*. 2014;345:455–459. <https://doi.org/10.1126/science.1249749>.
- Jinks-Robertson S, Bhagwat AS. Transcription-associated mutagenesis. *Annu Rev Genet*. 2014;48:341–359. <https://doi.org/10.1146/annurev-genet-120213-092015>.
- Jühling T et al. Small but large enough: structural properties of armless mitochondrial tRNAs from the nematode *Romanomermis culicivorax*. *Nucleic Acids Res*. 2018;46:9170–9180. <https://doi.org/10.1093/nar/gky593>.
- Jukes TH, Cantor CR. Evolution of protein molecules. In: Munro HN, editor. *Mammalian protein metabolism*. Academic Press; 1969. p. 21–132.
- Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermini LS. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods*. 2017;14:587–589. <https://doi.org/10.1038/nmeth.4285>.
- Kern AD, Kondrashov FA. Mechanisms and convergence of compensatory evolution in mammalian mitochondrial tRNAs. *Nat Genet*. 2004;36:1207–1212. <https://doi.org/10.1038/ng1451>.
- Kimura M. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotides

- sequences. *J Mol Evol.* 1980;16:111–120. <https://doi.org/10.1007/BF01731581>.
- Klemm BP et al. The Diversity of Ribonuclease P: Protein and RNA Catalysts with Analogous Biological Functions. *Biomolecules.* 2016;6:27. <https://doi.org/10.3390/biom6020027>.
- Li C, Qian W, Maclean CJ, Zhang J. The fitness landscape of a tRNA gene. *Science.* 2016;352:837–840. <https://doi.org/10.1126/science.aae0568>.
- Liu H, Zhang J. Higher germline mutagenesis of genes with stronger testis expressions refutes the transcriptional scanning hypothesis. *Mol Biol Evol.* 2020;37:3225–3231. <https://doi.org/10.1093/molbev/msaa168>.
- Liu X. 2013. Molecular recognition of inhibitors, metal ions and substrates by Ribonuclease P [Phd thesis]. University of Michigan, Ann Arbor, MI. Available at <https://deepblue.lib.umich.edu/handle/2027.42/100013>.
- Lorenz R et al. ViennaRNA package 2.0. *Algorithms Mol Biol.* 2011;6:26. <https://doi.org/10.1186/1748-7188-6-26>.
- MathWorks T. MATLAB version: 24.2.0.2871072 (r2024b), 2024.
- Meer MV, Kondrashov AS, Artzy-Randrup Y, Kondrashov FA. Compensatory evolution in mitochondrial tRNAs navigates valleys of low fitness. *Nature.* 2010;464:279–282. <https://doi.org/10.1038/nature08691>.
- Meiklejohn CD et al. An incompatibility between a mitochondrial tRNA and its nuclear-encoded tRNA synthetase compromises development and fitness in *Drosophila*. *PLoS Genet.* 2013;9:1–12. <https://doi.org/10.1371/journal.pgen.1003238>.
- Mitra S et al. Eukaryotic tRNAs fingerprint invertebrates vis-à-vis vertebrates. *J Biomol Struct Dyn.* 2015;33:2104–2120. PMID: 25581620. <https://doi.org/10.1080/07391102.2014.990925>.
- Murillo-Rocio M et al. 2025. Characterization of somatic mutations in human tRNA genes reveals tumor-specific mutational loads, and the generation of tRNA variants that alter the genetic code [preprint]. *bioRxiv.* <https://doi.org/10.1101/2025.06.09.658696>
- Ozerova I et al. Aberrant mitochondrial tRNA genes appear frequently in animal evolution. *Genome Biol Evol.* 2024;16:evae232. <https://doi.org/10.1093/gbe/evae232>.
- Pak D, Root-Bernstein R, Burton ZF. tRNA structure and evolution and standardization to the three nucleotide genetic code. *Transcription.* 2017;8:205–219. PMID: 28632998. <https://doi.org/10.1080/21541264.2017.1318811>.
- Palazzo AF, Lee ES. Non-coding RNA: what is functional and what is junk? *Front Genet.* 2015;6:2. <https://doi.org/10.3389/fgene.2015.00002>.
- Park C, Qian W, Zhang J. Genomic evidence for elevated mutation rates in highly expressed genes. *EMBO Rep.* 2012;13:1123–1129. <https://doi.org/10.1038/embor.2012.165>.
- Reids C, Forst CV, Schuster P. Replication and mutation on neutral networks. *Bull Math Biol.* 2001;63:57–94. <https://doi.org/10.1006/bulm.2000.0206>.
- Saks ME, Sampson JR, Abelson J. Evolution of a transfer RNA gene through a point mutation in the anticodon. *Science.* 1998;279:1665–1670. <https://doi.org/10.1126/science.279.5357.1665>.
- Schuster P, Fontana W, Stadler PF, Hofacker IL. From sequences to shapes and back: a case study in RNA secondary structures. *Proc R Soc Lond B Biol Sci.* 1994;255:279–284. <https://doi.org/10.1098/rspb.1994.0040>.
- Schwarz G. Estimating the dimension of a model. *Ann Stat.* 1978;6:461–464. <https://doi.org/10.1214/aos/1176344136>.
- Svejstrup JQ. Mechanisms of transcription-coupled DNA repair. *Nat Rev Mol Cell Biol.* 2002;3:21–29. <https://doi.org/10.1038/nrm703>.
- Tang DT, Glazov EA, McWilliam SM, Barris WC, Dalrymple BP. Analysis of the complement and molecular evolution of tRNA genes in cow. *BMC Genomics.* 2009;10:188. <https://doi.org/10.1186/1471-2164-10-188>.
- Tavaré S. Lectures on mathematics in the life sciences. Vol. 17. American Mathematical Society; 1986. 57–86. <https://api.semanticscholar.org/CorpusID:82212051>.
- Thornlow BP et al. Transfer RNA genes experience exceptionally elevated mutation rates. *Proc Natl Acad Sci U S A.* 2018;115:8996–9001. <https://doi.org/10.1073/pnas.1801240115>.
- Thornlow BP et al. Predicting transfer RNA gene activity from sequence and genome context. *Genome Res.* 2020;30:85–94. <https://doi.org/10.1101/gr.256164.119>.
- Tinoco I, Bustamante C. How RNA folds. *J Mol Biol.* 1999;293:271–281. <https://doi.org/10.1006/jmbi.1999.3001>.
- Trudeau DL, Kaltenbach M, Tawfik DS. On the potential origins of the high stability of reconstructed ancestral proteins. *Mol Biol Evol.* 2016;33:2633–2641. <https://doi.org/10.1093/molbev/msw138>.
- Wagner A. Evolvability-enhancing mutations in the fitness landscapes of an RNA and a protein. *Nat Commun.* 2023;14:3624. <https://doi.org/10.1038/s41467-023-39321-8>.
- Waldispühl J, Devadas S, Berger B, Clote P. Efficient algorithms for probing the RNA mutation landscape. *PLoS Comput Biol.* 2008;4:1–22. <https://doi.org/10.1371/journal.pcbi.1000124>.
- Westhof E, Thornlow B, Chan PP, Lowe TM. Eukaryotic tRNA sequences present conserved and amino acid-specific structural signatures. *Nucleic Acids Res.* 2022;50:4100–4112. <https://doi.org/10.1093/nar/gkac222>.
- Williams JD et al. Spontaneous deamination of cytosine to uracil is biased to the non-transcribed DNA strand in yeast. *DNA Repair (Amst).* 2023;126:103489. <https://doi.org/10.1016/j.dnarep.2023.103489>.
- Yang Z. Computational molecular evolution. Oxford University Press; 2006.
- Yona AH et al. tRNA genes rapidly change in evolution to meet novel translational demands. *Elife.* 2013;2:e01339. <https://doi.org/10.7554/eLife.01339>.
- Zhang J, Ferré-D'Amaré AR. The tRNA elbow in structure, recognition and evolution. *Life.* 2016;6:3. <https://doi.org/10.3390/life6010003>.
- Zuker M, Stiegler P. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Res.* 1981;9:133–148. <https://doi.org/10.1093/nar/9.1.133>.

Associate editor: Barbara Holland