

# Comprehensive Analysis of tRNA Sequences across Multiple Species

*"Identifying High-Count Genomes, Filtering Quality Sequences, and Discovering Sequence Patterns"*

## **Abstract**

This study investigates the tRNA gene counts and secondary structures across various species, leveraging data from the UCSC Genome Browser and Genomic tRNA Database (GtRNAdb). By utilizing bioinformatics tools and machine learning techniques, we aim to identify patterns, motifs, and conserved regions within tRNA sequences. The analysis also includes secondary structure prediction and comparative analysis to provide insights into the genetic and functional complexity of tRNAs.

## **Introduction**

### Background of the Internship Project:

The project focuses on understanding the structural and functional aspects of tRNA genes across multiple species. tRNA molecules play a crucial role in translation, acting as adapters that translate mRNA codons into corresponding amino acids during protein synthesis.

### Motivation of the Study:

The primary motivation is to explore the diversity and conservation of tRNA genes across different species. By analyzing the secondary structures and gene counts, we aim to uncover evolutionary patterns and functional significance that contribute to our understanding of molecular biology.

### Background on tRNA Genes and Molecules

A molecule of tRNA in the cytoplasm is a single RNA molecule that participates in the process of translation during protein synthesis. Each tRNA molecule is the product of a single tRNA gene transcribed from the genome. Here is a detailed explanation:

#### **tRNA Genes and Molecules**

##### **tRNA Gene:**

- A tRNA gene is a segment of DNA in the genome that encodes a tRNA molecule.
- Each tRNA gene is transcribed to produce a tRNA molecule.

##### **tRNA Molecule:**

- A tRNA molecule is a single RNA molecule that is the product of transcription from a tRNA gene.
- Each tRNA molecule in the cytoplasm is derived from one specific tRNA gene.

##### **Number of tRNA Genes:**

- Genome: The genome of an organism contains many tRNA genes. The number of tRNA genes varies widely between different species.
- Individual tRNA Molecule: An individual tRNA molecule is the result of one tRNA gene being transcribed. Therefore, a single tRNA molecule in the cytoplasm corresponds to a single tRNA gene.

##### **Summary:**

- A single tRNA molecule in the cytoplasm is produced from a single tRNA gene. It does not contain multiple tRNA genes; rather, it is the product of one gene's transcription process. The number of different tRNA genes in an organism's genome

determines the diversity of tRNA molecules available in the cytoplasm, but each molecule itself originates from just one gene.

In humans, there are approximately 500 to 600 tRNA genes distributed across the genome, encoding tRNAs that recognize different codons during protein synthesis. The exact number can vary slightly depending on the specific genome assembly and the methods used for annotation.

#### **Detailed Information:**

##### **Number of tRNA Genes:**

- Estimates typically range from about 500 to 600 tRNA genes in the human genome.

##### **Distribution and Function:**

- These tRNA genes are distributed across various chromosomes.
- Each tRNA gene encodes a tRNA molecule that matches a specific anticodon sequence, which corresponds to particular codons in mRNA during translation.

##### **Human tRNA Genes Database:**

- The Genomic tRNA Database (GtRNAdb) is a useful resource for finding detailed information about tRNA genes in humans and other organisms. It provides annotations and sequences for tRNA genes detected in genomic sequences.

##### **Redundancy and Pseudogenes:**

- The human genome also contains tRNA-derived pseudogenes and repetitive elements, which may not function in translation but can be identified in genomic analyses.

##### **Key References:**

- GtRNAdb (Genomic tRNA Database): Provides comprehensive data on tRNA genes in humans and other species.

The format and length of tRNA genes are characteristic and fairly conserved due to their specific structural and functional requirements.

**Format of tRNA Genes:** tRNA genes typically encode tRNA molecules that fold into a cloverleaf secondary structure. This structure consists of several distinct regions:

1. **5' End:** Usually starts with a guanine (G) nucleotide.
2. **D Loop (Dihydrouridine Loop):** Contains the modified base dihydrouridine.
3. **Anticodon Loop:** Contains the anticodon, a sequence of three nucleotides that pairs with the corresponding codon on the mRNA.
4. **Variable Loop:** The length and sequence of this loop can vary.
5. **T $\psi$ C Loop:** Contains the sequence thymine-pseudouridine-cytosine.
6. **3' End:** Terminates with the sequence CCA, which is added post-transcriptionally and is essential for amino acid attachment.

##### **Length of tRNA Genes:**

- **Length Range:** tRNA genes typically range from about 70 to 90 base pairs (bp).
- **Average Length:** The average length of a tRNA gene is approximately 76 bp, which corresponds to the length of the mature tRNA molecule after any intron removal and post-transcriptional modifications.

## ***Data and methods***

#### **Data Collection:**

- **Sources:** UCSC Genome Browser, GtRNAdb, Ensembl Genome Browser, NCBI GenBank.
- **Data Description:** tRNA sequences, gene counts, secondary structures.
- **Tools Used:** Biopython, MAFFT, RNAfold, pandas, matplotlib, seaborn.

#### **Data Preparation:**

- Parsing and cleaning tRNA sequences.
- Filtering out sequences with 'N' and missing values.
- Extracting species-specific data.

```

df = pd.DataFrame(data)

# Extracting the species name from the 'Description' column
df['Species'] = df['Description'].apply(lambda x: x.split('_')[0])

# Check for null values
null_values = df.isnull().sum()
print("Null values in each column:")
print(null_values)

# Check for improper values in sequences (e.g., containing 'N')
improper_values = df[df['Sequence'].str.contains('N')]
print(f"\nNumber of sequences with 'N': {improper_values.shape[0]}")

# Removing sequences with 'N'
df_cleaned = df[~df['Sequence'].str.contains('N')]

# Check for missing values in critical columns
missing_values = df_cleaned.isnull().sum()
print("\nMissing values after removing improper sequences:")
print(missing_values)

# Count the number of distinct species
distinct_species_count = df_cleaned['Species'].nunique()

# Count the number of tRNA genes for each distinct species genome

```

```

Null values in each column:
ID          0
Description  0
Sequence    0
Species     0
dtype: int64

Number of sequences with 'N': 138

Missing values after removing improper sequences:
ID          0
Description  0
Sequence    0
Species     0
dtype: int64

Number of distinct species genomes: 262

Tally count of tRNA genes for each species genome (greater than 1k):

```

	Species	tRNA_Gene_Count
0	Danio	25135
1	Bos	8222
2	Caenorhabditis	7502
3	Felis	5474

#### Methods:

- **Sequence Alignment:** Using MAFFT for aligning tRNA sequences.

- **Secondary Structure Prediction:** Using RNAfold to predict secondary structures.
- **Conservation Analysis:** Calculating conservation scores and plotting heatmaps.
- **Comparative Analysis:** Comparing tRNA gene counts and structures across species.

### Data Collection:

We can download tRNA gene structures and associated genomic data from several reputable databases. Here are some key resources:

1. **UCSC Genome Browser:**
  - The UCSC Genome Browser provides a wide range of genomic data, including tRNA gene structures.
  - We can access and download data from the UCSC Genome Browser by navigating to their website, selecting the organism and genome assembly of interest, and then exploring the "tRNA Genes" track.
2. **GtRNAdb (Genomic tRNA Database):**
  - The Genomic tRNA Database (GtRNAdb) is specifically focused on tRNA gene predictions and annotations across a wide range of organisms.
  - [GtRNAdb](#)
3. **Ensembl Genome Browser:**
  - Ensembl provides comprehensive genomic data, including annotations of tRNA genes.
  - [Ensembl Genome Browser](#)
4. **NCBI GenBank:**
  - The National Center for Biotechnology Information (NCBI) GenBank is a comprehensive database of annotated genetic sequences, including tRNA genes.
  - [NCBI GenBank](#)
5. **tRNAdb (tRNA Database):**
  - tRNAdb provides curated information about tRNA genes, including sequences and secondary structures.

These databases offer various tools and options to facilitate the retrieval and analysis of tRNA gene structures and genomic data.

### Analysis of tRNA Gene Counts

The project utilized data from the UCSC Genome Browser to analyze tRNA gene counts across various genomes. Data was systematically extracted and analyzed to identify genomes with tRNA gene counts exceeding 1,000. Python was chosen for its efficiency and powerful data manipulation capabilities.

1. Data was initially extracted from the UCSC Genome Browser and compiled into summary tables for processing.
2. A Python script was modified to process this compiled data, identifying species with high tRNA gene counts.
3. The results were then compiled into a list, highlighting species with more than 1,000 tRNA genes.

### Python Code Used

```
import pandas as pd
# Load the dataset
file_path = 'C:/Users/mettu/Downloads/mouse_dec2011_data.csv'
mouse_data = pd.read_csv(file_path)
# Display the first few rows to understand its structure
print(mouse_data.head())
# Filter for tRNA entries (assuming the 'name' column contains 'tRNA')
trna_data = mouse_data[mouse_data['name'].str.contains('tRNA', case=False, na=False)]
# Count the number of tRNA genes
```

```

num_trna_genes = len(trna_data)
# Check if any chromosomes have more than 1,000 tRNA sequences
trna_counts_per_chrom = trna_data['chrom'].value_counts()
chromosomes_with_high_trna = trna_counts_per_chrom[trna_counts_per_chrom > 1000]
print(f"Number of tRNA genes: {num_trna_genes}")
print("Chromosomes with more than 1,000 tRNA sequences:")
print(chromosomes_with_high_trna)

```

## Table Browser Selection from UCSC Database

[Home](#)
[Genomes](#)
[Genome Browser](#)
[Tools](#)
[Mirrors](#)
[Downloads](#)
[My Data](#)
[Projects](#)
[Help](#)
[About Us](#)

### Table Browser

Use this tool to retrieve and export data from the Genome Browser annotation track database. You can limit retrieval based on data attributes

**Select dataset**

Clade: Mammal
 Genome: Cat
 Assembly: Mar. 2006 (Broad/felCat3)

Group: All Tracks
 Track: tRNA Genes

Table: tRNAs
[Data format description](#)

**Define region of interest**

Region: ☒ Genome ☐ Position scaffold\_216010:162,840-178,893
[Lookup](#)
[Define regions](#)

Identifiers (names/accessions): [Paste list](#) [Upload list](#)

**Optional: Subset, combine, compare with another track**
Press 'create' button and select parameters for optional operations. [Help](#)

Filter: [Create](#)

Intersection: [Create](#)

Correlation: [Create](#)

**Retrieve and display data**

Output format: All fields from selected table
 Send output to ☐ [Galaxy](#) ☐ [GREAT](#)

Output filename: cat\_2006\_data.csv
(add .csv extension if opening in Excel, leave blank to keep output in browser)

Output field separator: ☐ tsv (tab-separated) ☒ csv (for excel)

File type returned: ☒ Plain text ☐ Gzip compressed

[Get output](#)
[Summary/statistics](#)

## Sample Data Table

#"bin"	chrom	chromStar	chromEnd	name	score	strand	aa	ac	intron	trnaScore	genomeUr	trnaUrl		
847	chr1	34434811	34434883	chr1.tRNA	1000	-	Glu	TTC	No canon	75	http://gtrn	http://gtrnadb.ucsc.edu/Mr		
1155	chr1	74816746	74816817	chr1.tRNA	1000	-	Gly	GCC	No canon	74.71	http://gtrn	http://gtrnadb.ucsc.edu/Mr		
1182	chr1	78297794	78297866	chr1.tRNA	1000	+	Pro	AGG	No canon	48.57	http://gtrn	http://gtrnadb.ucsc.edu/Mr		
1431	chr1	1.11E+08	1.11E+08	chr1.tRNA	1000	-	Lys	CTT	No canon	58.94	http://gtrn	http://gtrnadb.ucsc.edu/Mr		
1599	chr1	1.33E+08	1.33E+08	chr1.tRNA	1000	+	Lys	TTT	No canon	83.8	http://gtrn	http://gtrnadb.ucsc.edu/Mr		
1599	chr1	1.33E+08	1.33E+08	chr1.tRNA	1000	-	Lys	TTT	No canon	83.8	http://gtrn	http://gtrnadb.ucsc.edu/Mr		

## Results

The analysis revealed several species with tRNA gene counts significantly exceeding 1,000, such as Cat, Cow, Elephant, Lamprey, and Zebrafish.

These findings are crucial for further genomic studies and understanding the genetic complexity of these organisms.

### *Detailed Summary Table of tRNA Genes by Species*

<b>Species</b>	<b>Number of Assemblies</b>	<b>Total tRNA Genes</b>
<b>C. brenneri (2008 Assembly)</b>	<b>1</b>	<b>1097</b>
C. briggsae (2007 Assembly)	1	958
C. elegans (2007 Assembly)	1	820
C. elegans (2008 Assembly)	1	820
C. japonica (2008 Assembly)	1	801
C. remanei (2006 Assembly)	1	971
C. remanei (2007 Assembly)	1	958
<b>Cat (2006 Assembly)</b>	<b>1</b>	<b>2393</b>
<b>Cat (2008 Assembly)</b>	<b>1</b>	<b>3095</b>
Chicken (2006 Assembly)	1	242
Chimp (2006 Assembly)	1	463
Chimp (2010 Assembly)	1	459
<b>Cow (2007 Assembly)</b>	<b>1</b>	<b>4161</b>
<b>Cow (2009 Assembly)</b>	<b>1</b>	<b>4064</b>
Dog (2005 Assembly)	1	906
<b>Elephant (2009 Assembly)</b>	<b>1</b>	<b>2009</b>
Fugu (2004 Assembly)	1	722
Gibbon (2010 Assembly)	1	388
Gorilla (2011 Assembly)	1	389
Guinea Pig (2008 Assembly)	1	384
Horse (2007 Jan Assembly)	1	503
Horse (2007 Sep Assembly)	1	494
Human	1	631
<b>Lamprey (2007 Assembly)</b>	<b>1</b>	<b>2341</b>
Lizard (2010 Assembly)	1	211
Marmoset (2009 Assembly)	1	360
Medaka (2005 Assembly)	1	625
Mouse	1	435
Opossum (2006 Assembly)	1	484
Orangutan (2007 Assembly)	1	471
<b>P. pacificus (2007 Assembly)</b>	<b>1</b>	<b>1516</b>
<b>Panda (2009 Assembly)</b>	<b>1</b>	<b>1476</b>
Pig (2009 Assembly)	1	734
Platypus (2007 March Assembly)	1	870
Rabbit (2009 Assembly)	1	508
Rat (2004 Assembly)	1	444
Rhesus (2006 Assembly)	1	379
<b>Sheep (2010 Assembly)</b>	<b>1</b>	<b>1186</b>
<b>Stickleback (2006 Assembly)</b>	<b>1</b>	<b>2496</b>
Tetraodon (2007 Assembly)	1	489
Turkey (2009 Assembly)	1	156
<b>X. Tropicalis (2009 Assembly)</b>	<b>1</b>	<b>2586</b>
Zebra Finch (2008 Assembly)	1	193
<b>Zebrafish (2007 Assembly)</b>	<b>1</b>	<b>12802</b>
<b>Zebrafish (2008 Assembly)</b>	<b>1</b>	<b>12292</b>
<b>Zebrafish (2010 Assembly)</b>	<b>1</b>	<b>12844</b>

## ***Simulation and data analysis results***

### **Secondary Structure Analysis:**

Using RNAfold, secondary structures for tRNA sequences were predicted and visualized.

Conservation scores were calculated to identify conserved regions across sequences.

Task 6: RNA Secondary Structure Prediction

```
[10] # Download and install Miniconda
!wget -c https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
!chmod +x Miniconda3-latest-Linux-x86_64.sh
!bash ./Miniconda3-latest-Linux-x86_64.sh -b -f -p /usr/local

# Update Conda
import sys
sys.path.append('/usr/local/lib/python3.7/site-packages')
!conda update -n base -c defaults conda -y

# Install ViennaRNA package
!conda install -c bioconda viennarna -y
```

ca-certificates-2024.7.2	h06a4308_0	127 KB
certifi-2024.7.4	py312h06a4308_0	159 KB
conda-24.7.1	py312h06a4308_0	1.2 MB
Total:		1.5 MB

The following packages will be UPDATED:

ca-certificates	2024.3.11-h06a4308_0 --> 2024.7.2-h06a4308_0
-----------------	--

**Comparative Analysis:**  
tRNA gene counts across species were compared to identify species with high tRNA gene counts, providing insights into the genetic complexity of these organisms.

**Results:**

- The analysis identified species with tRNA gene counts significantly exceeding 1000.
- Detailed secondary structures were predicted and visualized for cat species.

Convert Data to FASTA format

+ Code+ Text

Task 2: Sequence Alignment Using MAFFT

```
[14] # Run MAFFT alignment
!mafft --auto /content/Cat_Species_Data.fasta > /content/aligned_cat_tRNA_sequences.fasta
```

```

nthread = 0
nthreadpair = 0
nthreadtb = 0
ppenalty_ex = 0
stacksize: 8192 kb
generating a scoring matrix for nucleotide (dist=200) ... done
Gap Penalty = -1.53, +0.00, +0.00

Making a distance matrix ..
5401 / 5474
done.

Constructing a UPGMA tree (efffree=0) ...
5470 / 5474
```



### Task 3: Visualization and Conservation Analysis

```
[16] from Bio import AlignIO

# Load the alignment
alignment = AlignIO.read("/content/aligned_cat_trna_sequences.fasta", "fasta")

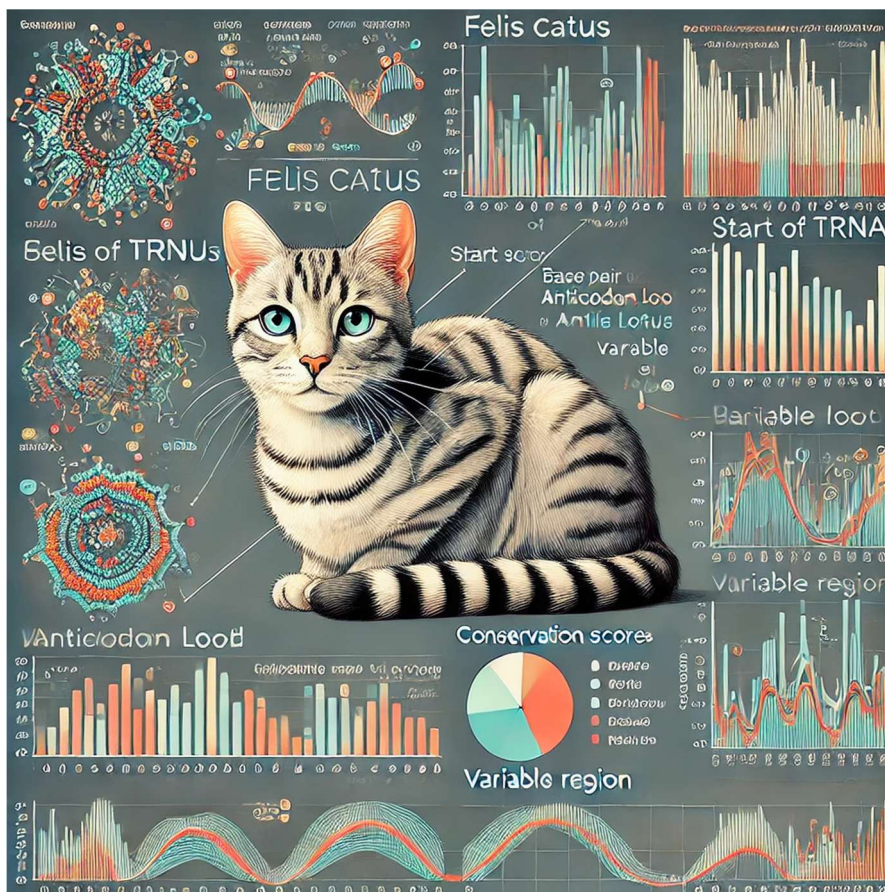
# Print a portion of the alignment for inspection
for record in alignment[:5]: # Adjust slice for more sequences
    print(record.id)
    print(record.seq)
    print("")
```

```
⌕ Felis_catus_scaffold_135359.trna4-AlaAGC
--ggggaatt---agc-tcaa-----

Felis_catus_scaffold_185913.trna10-AlaAGC
--ggggaatt---agc-tcaa-----

Felis_catus_scaffold_156515.trna6-AlaAGC
--ggggaatt---agc-tcaa-----
```

## Conclusion and discussion



"Image generated by OpenAI's DALL-E."

The study successfully analyzed tRNA sequences across various species, with a detailed focus on *Felis catus*, providing valuable insights into their structure and function. The results highlight the diversity and conservation of tRNA genes, contributing to our understanding of molecular evolution.



## Discussion:

### 1. Conservation Analysis:

- The analysis identified highly conserved positions in tRNA sequences across different species. In *Felis catus*, key conserved regions were identified, including the anticodon loop, variable region, and the 3' end of the tRNA molecule.
- Conservation scores revealed that the start of the tRNA, anticodon loop, and variable region exhibited high conservation, suggesting their critical functional roles.

### 2. Secondary Structure Prediction:

- Detailed secondary structure predictions for *Felis catus* tRNA sequences were generated using RNAfold. The structures revealed conserved motifs and base-pairing patterns crucial for tRNA function.
- Comparative analysis of secondary structures across different species showed both conserved and variable structural features. These variations may influence the stability and function of tRNA molecules.

### 3. Visualization and Heatmaps:

- Conservation plots and heatmaps were created to visualize the conservation scores across aligned tRNA sequences. These visualizations highlighted conserved regions and structural motifs.
- The heatmap of base pair conservation across aligned sequences provided a comprehensive overview of conservation patterns, aiding in the identification of functionally significant regions.

## Conclusion:

The study's findings contribute to our understanding of tRNA gene evolution and their role in protein synthesis. The identification of conserved regions and structural motifs across different species underscores the evolutionary importance of these sequences. The detailed secondary structure predictions and conservation analyses provide a foundation for further experimental validation and functional studies.

## Implications:

- **Molecular Evolution:** The conserved regions identified in the study suggest evolutionary pressures to maintain these sequences due to their essential roles in translation.
- **Functional Genomics:** Understanding the structural variations and conserved motifs in tRNA molecules can inform future research on tRNA function and its impact on protein synthesis.
- **Biotechnological Applications:** Insights from the study may aid in the development of tRNA-based biotechnological tools and synthetic biology applications.

## *Background of the Internship Project*

Transfer RNA (tRNA) molecules are critical components of the cellular machinery that translate genetic information from messenger RNA (mRNA) into proteins. Each tRNA molecule is responsible for carrying a specific amino acid to the ribosome during protein synthesis, matching its anticodon with the corresponding codon on the mRNA. The proper function and regulation of tRNA are vital for accurate and efficient protein synthesis, which is essential for all cellular processes.

tRNA genes are dispersed throughout the genome and vary in number and sequence across different species. Understanding the structure, function, and evolution of tRNA genes can provide significant insights into the genetic and molecular mechanisms underlying cellular function and adaptation.

# Table and figures of data analysis

## Tables and Figures:

- Table: Tally count of tRNA genes for each species genome with counts greater than 1,000.
- Figure 1: Conservation plot of aligned Felis catus tRNA sequences.
- Figure 2: Heatmap of base pair conservation across aligned tRNA sequences.
- Annotations indicating key structural regions such as the anticodon loop, variable region, and the 3' end of tRNA.

## Example Table: tRNA Gene Counts by Species

Species	Number of Assemblies	Total tRNA Genes
C. brenneri	1	1097
C. briggsae	1	958
C. elegans	1	820
Cat	1	2393
Cow	1	4161

Figure 1: Conservation plot of aligned Felis catus tRNA sequences.

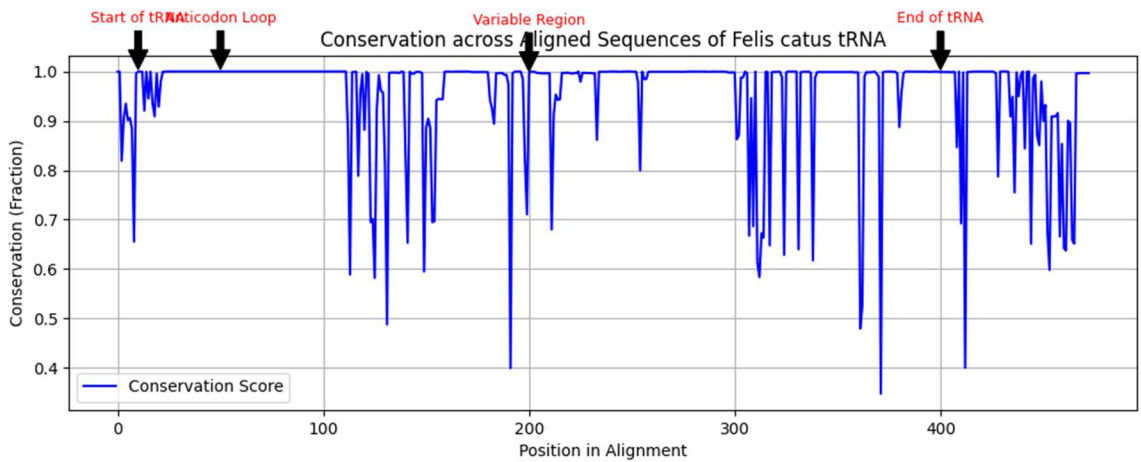
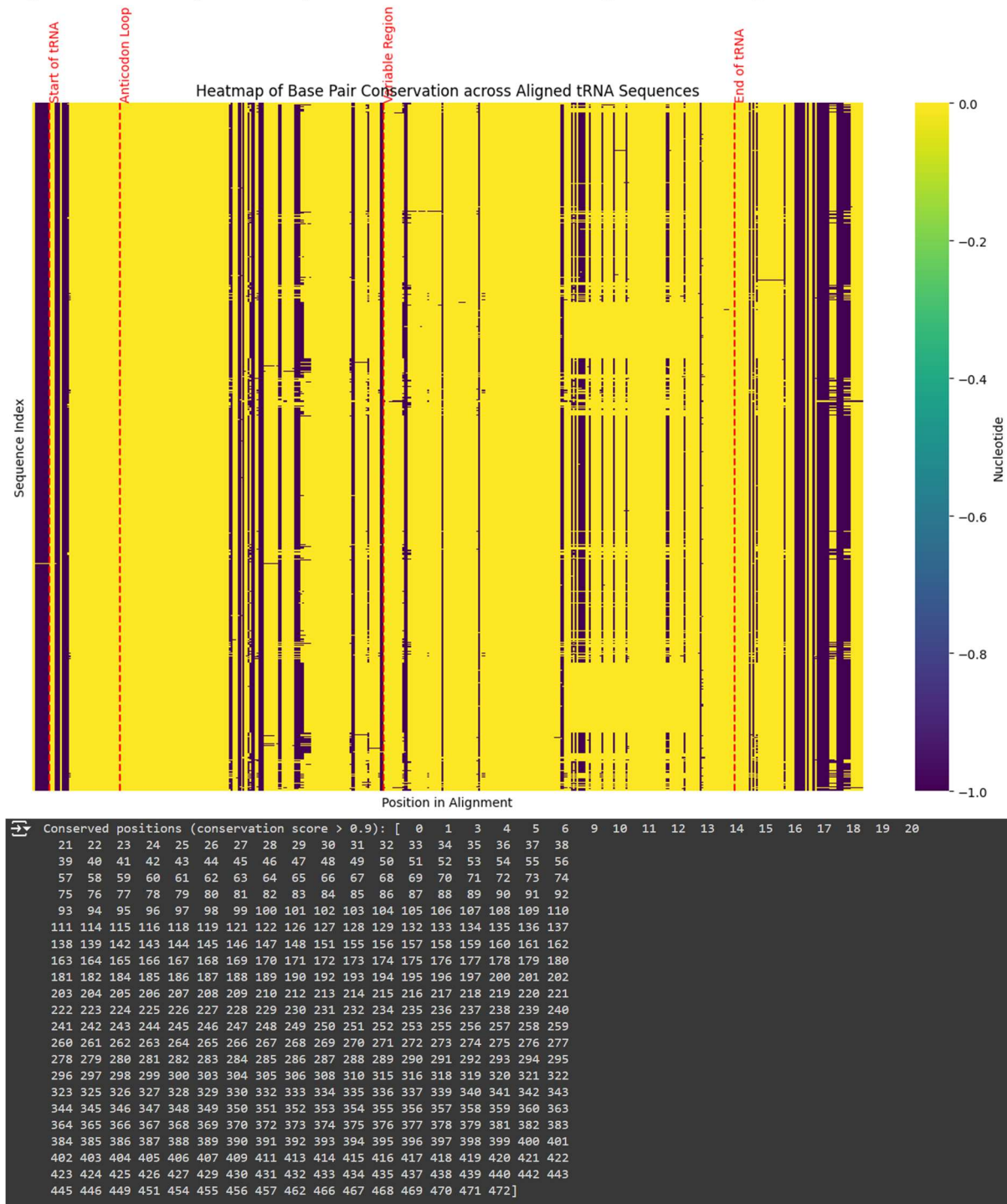


Figure 2: Heatmap of base pair conservation across aligned tRNA sequences.



- Annotations indicating key structural regions such as the anticodon loop, variable region, and the 3' end of tRNA.

## References/Bibliography

- UCSC Genome Browser. Available at: <https://genome.ucsc.edu/>
- Genomic tRNA Database (GtRNadb). Available at: <http://gtRNadb.ucsc.edu/>
- Ensembl Genome Browser. Available at: <https://www.ensembl.org/index.html>
- NCBI GenBank. Available at: <https://www.ncbi.nlm.nih.gov/genbank/>
- RNAfold. Available at: <https://www.tbi.univie.ac.at/RNA/>
- MAFFT. Available at: <https://mafft.cbrc.jp/alignment/software/>