# Creating new and updated codon usage tables in HIVE using species-specific genomic and tissue-specific transcriptomic information
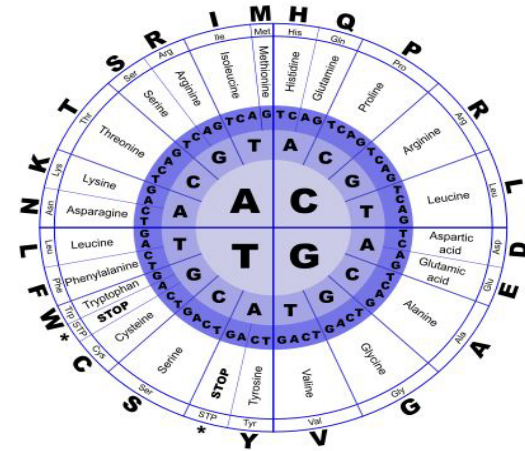
Luis Santana-Quintero, PhD
Office of Biostatistics and Pharmacovigilance
**Center for Biologics Evaluation and Research**
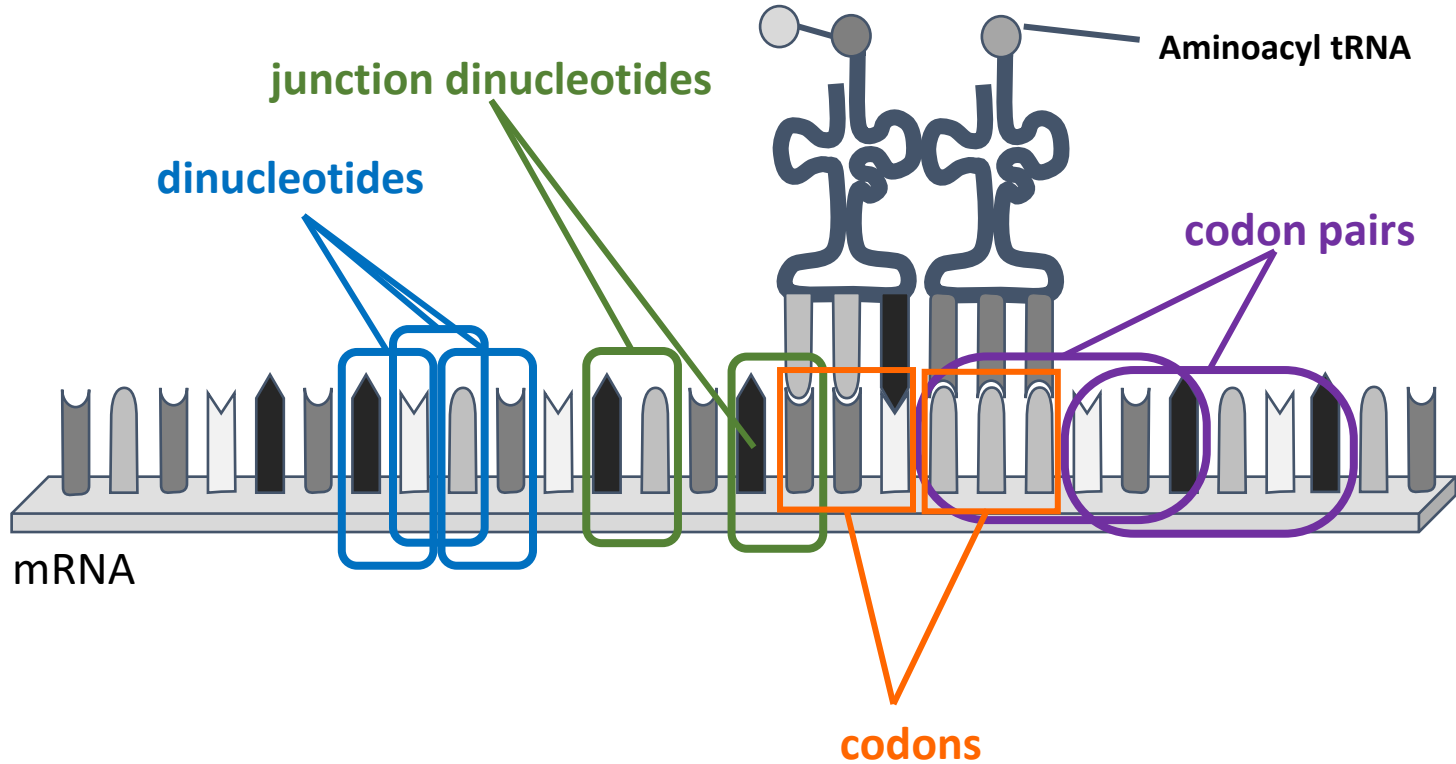
# Outline

- Background
  - Codon, Codon Pairs, dinucleotides and junction dinucleotides

- Introduction

- Methods
  - Codon Usage Tables (CUT)
  - Transcriptomic Weighted Usage

- HIVE Platform

- Results and Examples

- Conclusions

- Acknowledgements

# Background

- A codon is a DNA (or RNA) sequence of three nucleotides that encodes a particular **amino acid** or signaling the start/stop of protein synthesis.

- Genetic code is made up of codons and there are **64 different codons**.

- The codon usage tables are linked to a taxonomic reference, and they allow comparative analysis of the **codon usage frequencies**.

# Background: Codons, Codon Pairs and Dinucleotides

# Introduction

- We have created updated usage tables with the sequence information from GenBank and RefSeq.  Including information for the non-random distribution of occurrences in genes within a given species:
  - Codon
  - Codon pair
  - Dinucleotide
  - Junction Dinucleotide
  - GC content

In total, 288 million coding sequences (35 million from GenBank, 253 million from RefSeq) were included in the database, resulting in the creation of over 855,000 codon usage tables.

- Accounting for differential gene expression profiles in various human tissues, we have also created usage tables for normal **human tissues** and for **human primary tumors**.

# Methods

We created a publicly accessible repository of comprehensive, regularly updated Codon Usage Tables, HIVE-CUTs.

- CoCoPUTs – Codon and Codon Pair Usage Tables at species level (Genbank & RefSeq)

- TissueCoCoPUTs – 52 Human tissue-specific from GTEx Portal

- CancerCoCoPUTs – Tumor-specific contain transcriptome-derived data from 32 primary cancer types from TCGA

Athey *et. al*., *BMC Bioinformatics*, 2017
Alexaki *et. al,* Journal of Molecular Biology, 2019
Holcomb *et. al*, *Infection, Genetics and Evolution*, 2019

Kames *et. al*., *Journal of Molecular Biology*, 2019
Meyer *et. al*., *Genome Medicine*, 2021

# Codon Usage Table

Input: Listeria Genbank file (.gbff file)

Output: Listeria raw count and proportion of each codon that appears on the sequence

Protein sequence

CDS genomic sequence

Listeria (1637) Codon Usage Table

Table contains 12627246 CDSs (3847640570 codons), taken from RefSeq.

To select all data, click on the table and then press Ctrl+A.

Codon usage table

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| TTT | 30.91 | (118938555) | TCT | 12.69 | ( 48832316) | TAT | 23.73 | ( 91295446) | TGT | 4.11 | (15826247) |
| TTC | 14.38 | ( 55328760) | TCC | 6.36 | ( 24476826) | TAC | 10.77 | ( 41436606) | TGC | 1.99 | ( 7657220) |
| TTA | 36.97 | (142264519) | TCA | 10.23 | ( 39346546) | TAA | 2.29 | ( 8827467) | TGA | 0.59 | ( 2275879) |
| TTG | 12.99 | ( 49997367) | TCG | 6.32 | ( 24326534) | TAG | 0.39 | ( 1486335) | TGG | 9.32 | (35872854) |
| | | | | | | | | | | | |
| CTT | 21.10 | ( 81166995) | CCT | 7.90 | ( 30383779) | CA | | (7132305) |
| CTC | 5.72 | ( 21989291) | CCC | 1.72 | ( 6613497) | CA | | (8118617) |
| CTA | 12.93 | ( 49755746) | CCA | 17.65 | ( 67901523) | CAA | 29.12 | (112049842) | CGA | 5.91 | (22722832) |
| CTG | 5.29 | ( 20361616) | CCG | 7.23 | ( 27812164) | CAG | 5.43 | ( 20880675) | CGG | 3.08 | (11845696) |
| | | | | | | | | | | | |
| ATT | 50.25 | (193337084) | ACT | 15.66 | ( 60240444) | AAT | 31.88 | (122673962) | AGT | 13.76 | (52957091) |
| ATC | 18.33 | ( 70528562) | ACC | 7.00 | ( 26932588) | AAC | 14.46 | ( 55638648) | AGC | 8.42 | (32385364) |
| ATA | 9.47 | ( 36441806) | ACA | 25.01 | ( 96234848) | AAA | 60.47 | (232679709) | AGA | 6.84 | (26304995) |
| ATG | 26.70 | (102735509) | ACG | 13.21 | ( 50818597) | AAG | 10.89 | ( 41889161) | AGG | 1.32 | ( 5085993) |
| | | | | | | | | | | | |
| GTT | 26.43 | (101707555) | GCT | 23.00 | ( 88478148) | GAT | 39.87 | (153420147) | GGT | 23.54 | (90556103) |
| GTC | 9.27 | ( 35679738) | GCC | 8.80 | ( 33851641) | GAC | 14.30 | ( 55023297) | GGC | 14.54 | (55948592) |
| GTA | 21.00 | ( 80815623) | GCA | 27.87 | (107235318) | GAA | 60.58 | (233096896) | GGA | 19.14 | (73660319) |
| GTG | 13.65 | ( 52510924) | GCG | 17.29 | ( 66540277) | GAG | 13.50 | ( 51943665) | GGG | 9.08 | (34949864) |

7

# Cancer CoCoPUTs

RNA-seq files from NCI containing counts from 32 primary cancer types from TCGA

- Primary tumor
- Solid Tissue Normal

https://cancergenome.nih.gov/

THE CANCER GENOME ATLAS

| | 1 | 2 | 3 | | 19018 | 19019 |
|---|---|---|---|---|---|---|
| AAA | 1 | 3 | 6 | ... | 1 | 4 |
| AAT | 2 | 4 | 1 | ... | 3 | 1 |
| AAC | 6 | 34 | 1 | ... | 11 | 31 |
| | ... | ... | ... | ... | ... | ... |
| GGC | 2 | 19 | 9 | ... | 5 | 13 |
| GGG | 0 | 23 | 11 | ... | 3 | 7 |

| GENE | TPM-1 | TPM-2 |
|---|---|---|
| 1 | 54 | 700 |
| 2 | 438 | 22 |
| 3 | 2 | 1754 |
| ... | ... | ... |
| 19018 | 828 | 876 |
| 19019 | 9772 | 9821 |

| | Sample1 | Sample2 |
|---|---|---|
| AAA | 568,695 | 846,256 |
| AAT | 592,587 | 187,695 |
| AAC | 113,968 | 257,665 |
| | | |
| GGC | 12,896 | 385,411 |
| GGG | 1,898,987 | 75,850 |

Gene level usage
(codon and codon pair)

Primary transcript quantifications
(TPM for each gene)

Transcriptomic
Weighted Usage

# HIVE Platform

o A **cloud-based environment** that comprises both a storage library of data and a powerful computing capacity.

o Can **consume**, **digest**, **analyze**, **manage**, and **share all** this data.



High throughput infiniband multichannel network 40Gb/s

Metadata Database

**Distributed storage cloud**

**Cloud control server**

**Distributed computational cloud**

**Web portal drop-box**

**Instruments**

NCBI

NIH NATIONAL CANCER INSTITUTE

**External data providers**

**Local hard drives**

**Web-browser**

# Results: CoCoPUTs
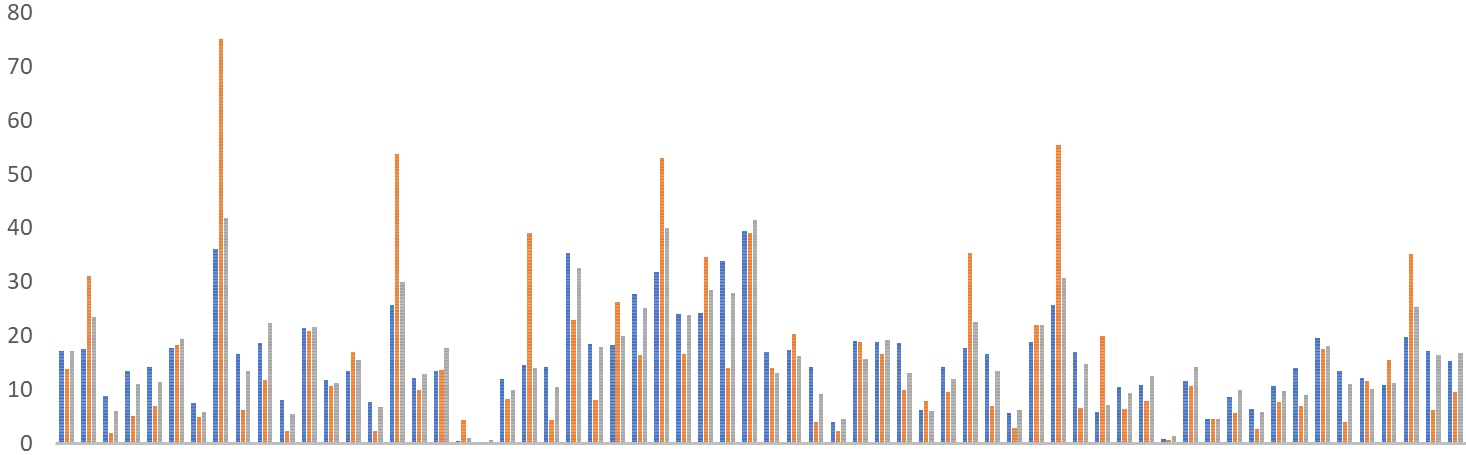
# Results: TissueCoCoPUTs

# Results: CancerCoCoPUTs

# Example: Tissue-Specific Codon Usage

- Calculated from tissue-specific transcriptome (GTEx Portal and Protein Atlas) and human codon usage per CDS
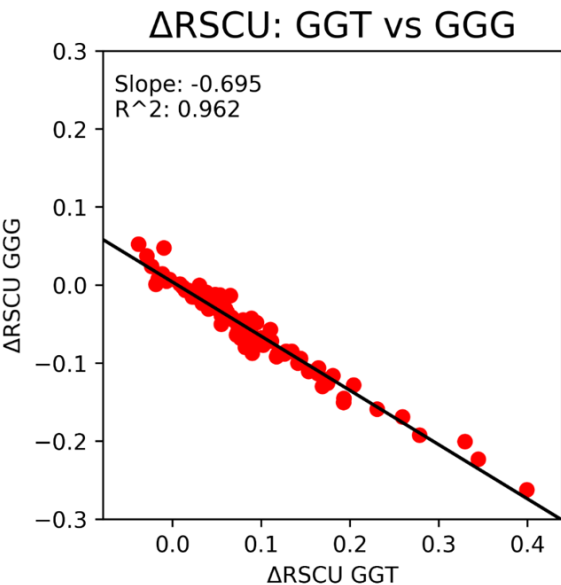- Results show large variability between tissues

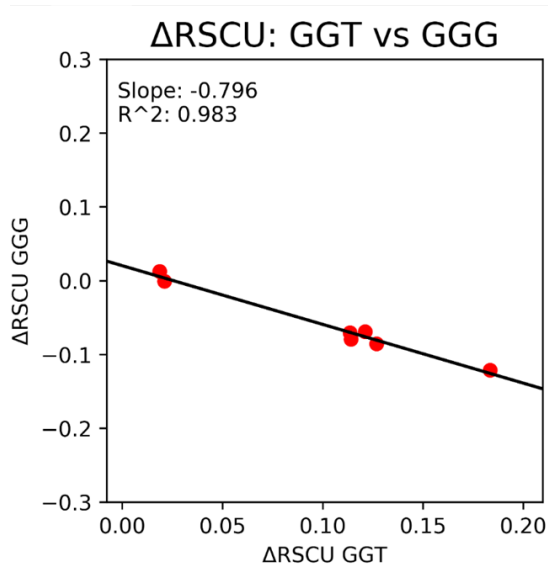**Codon Usage Frequencies between Genomic, Whole Blood and Liver**

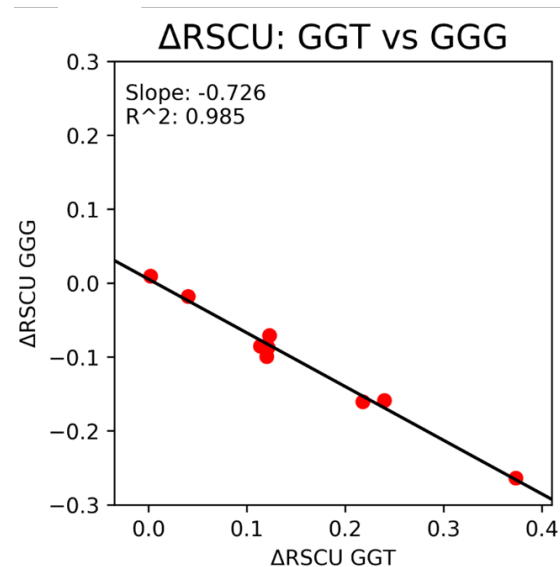# Example: Strong Negative Relationship Between GGT and GGG across Breast Cancer Subtypes



Meyer D., Kames J. *et. al.*, *Genome Medicine*, 2021
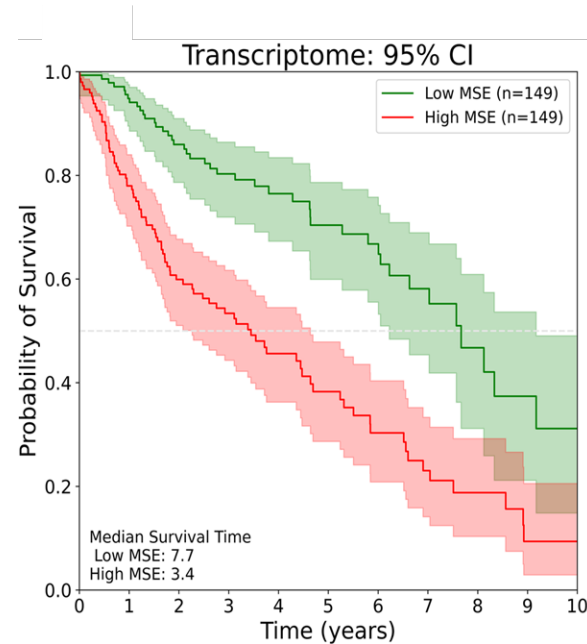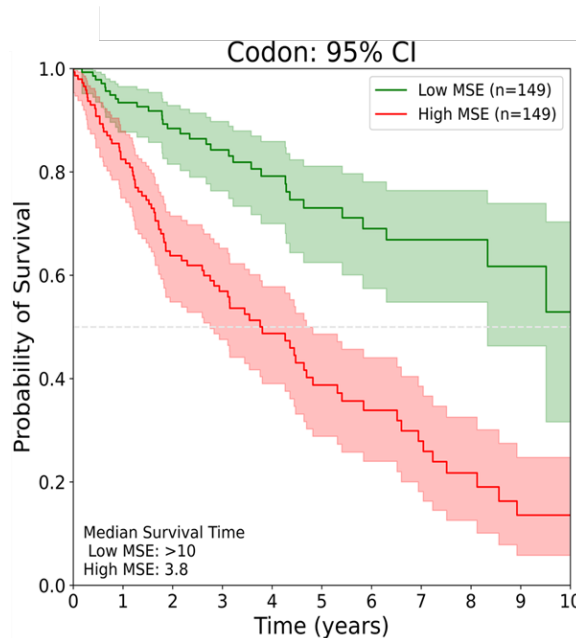
# Example: Change in Codon Usage Associated with Patient Survival Across All Cancers

MSE (Mean Squared Error):
Degree of Change



Meyer D., Kames J. *et. al.*, *Genome Medicine*, 2021

# Conclusions

Codon Usage Tables provide the frequency of occurrence of codons and is essential in many biological studies and applications, e.g.
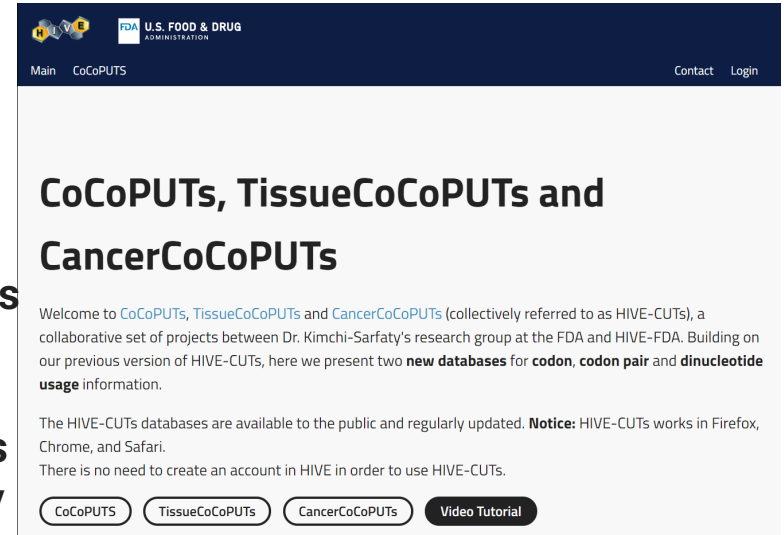
- Drug Development
- Gene Therapy
- Vaccine Development
- Implications of synonymous variants
- Recombinant Therapeutic protein
- Personalized Cancer Medicine
- Virus de-optimization design
- Evolutionary and Translational studies

# Availability

https://dnahive.fda.gov/dna.cgi?cmd=cuts_main

1. Clickable heatmap displaying codon pair information
2. Text-based codon usage table
3. Graphs plotting the codon frequencies per 1000 codons
4. Graphs showing total dinucleotide and junction dinucleotide frequencies
5. Graphs showing the GC% content
6. Effective Number of Codons (ENC) and Codon Pairs (ENCP), metrics measuring codon and codon pair usage bias
7. Taxonomy tree that displays each query and traces the route back to their last shared classification (only available for species related queries)

## CoCoPUTs, TissueCoCoPUTs and CancerCoCoPUTs

Welcome to CoCoPUTs, TissueCoCoPUTs and CancerCoCoPUTs (collectively referred to as HIVE-CUTs), a collaborative set of projects between Dr. Kimchi-Sarfaty's research group at the FDA and HIVE-FDA. Building on our previous version of HIVE-CUTs, here we present two **new databases** for **codon**, **codon pair** and **dinucleotide usage** information.

The HIVE-CUTs databases are available to the public and regularly updated. **Notice:** HIVE-CUTs works in Firefox, Chrome, and Safari.
There is no need to create an account in HIVE in order to use HIVE-CUTs.

CoCoPUTS     TissueCoCoPUTs     CancerCoCoPUTs     Video Tutorial

# Acknowledgements

# Thank you!