The Codon Statistics Database: A Database of Codon Usage Bias

Krishnamurthy Subramanian, †, † Bryan Payne, † Felix Feyertag, and David Alvarez-Ponce (6)*

Biology Department, University of Nevada, Reno, Reno, NV 89557, USA

Associate editor: Naruya Saitou

Abstract

We present the Codon Statistics Database, an online database that contains codon usage statistics for all the species with reference or representative genomes in RefSeq (over 15,000). The user can search for any species and access two sets of tables. One set lists, for each codon, the frequency, the Relative Synonymous Codon Usage, and whether the codon is preferred. Another set of tables lists, for each gene, its GC content, Effective Number of Codons, Codon Adaptation Index, and frequency of optimal codons. Equivalent tables can be accessed for (1) all nuclear genes, (2) nuclear genes encoding ribosomal proteins, (3) mitochondrial genes, and (4) chloroplast genes (if available in the relevant assembly). The user can also search for any taxonomic group (e.g., "primates") and obtain a table comparing all the species in the group. The database is free to access without registration at http://codonstatsdb.unr.edu.

Key words: codon bias, codon usage, database, synonymous codons.

Introduction

Most amino acids are encoded by multiple synonymous codons. Despite encoding for the same amino acid, some synonymous codons are used significantly more often than others, a phenomenon known as codon usage bias. Species significantly differ in their codon preferences—for instance, glutamic acid is preferentially encoded by GAG in human, whereas the same amino acid is preferentially encoded by GAA in Escherichia coli (Ikemura 1982; Sharp et al. 2010). In addition, genes within any given genome differ in their patterns of codon usage. In particular, gene expression levels significantly correlate with gene-specific metrics of codon usage such as the Effective Number of Codons (ENC; Wright 1990), the Codon Adaptation Index (CAI; Sharp and Li 1987), or the frequency of optimal codons $(F_{op}; Ikemura 1985)$ (e.g., Gouy and Gautier 1982).

Codon preferences can be affected by a number of factors, including the genome's nucleotide composition (e.g., AT-rich genomes tend to use codons ending in A or T) and translational selection (codons that are translated by highly abundant tRNAs are translated faster and with fewer errors; e.g., Ikemura 1985; Hershberg and Petrov 2008).

Understanding codon preferences across the different species and genes is important not only to understanding genome evolution, but can also inform tasks such as heterologous expression, gene prediction, or phylogenetic inference (e.g., Gustafsson et al. 2004; Christianson 2005; Rota-Stabelli et al. 2013). In addition, the patterns of codon usage of viruses tend to be similar to those of their host species (e.g., Shackelton et al. 2006). Despite the relevance of maintaining species- and gene-specific codon usage information, existing databases have not been updated in a long time, focus on specific taxa, and/or do not provide gene-specific metrics (Nakamura et al. 2000; Hilterbrand et al. 2012; Athey et al. 2017).

Implementation

For each of the species with reference or representative genomes in the RefSeq database (release 207), we chose one full assembly (in order of preference, the one labeled as "reference," the one with the highest assembly level, or the most recent one) and retrieved the corresponding coding sequences (CDSs) file. Using that file as input, a number of tables were pre-computed using an R pipeline. To avoid codon redundancy, only one CDS per gene was used (if multiple were available, the longest one was chosen). The web interface was created using PERL CGI.

For each species, we computed the total frequency of each codon, and used the information to compute the Relative Synonymous Codon Usage (RSCU) of each codon. For each gene, we computed the GC content for the entire CDS (GC), the GC content at third codon positions (GC3), the ENC, and the RSCU for each codon.

For species with over 1,000 genes, we also compared genes inferred to be highly expressed (bottom 10% ENC values) with genes inferred to be lowly expressed (top 10% ENC values). Codons with significantly higher RSCU values in the highly expressed gene set (according to a

© The Author(s) 2022. 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 journals.permissions@oup.com

Open Access

^{*}Corresponding author: E-mail: dap@unr.edu.

[†]These authors contributed equally to this work.

[†]Present address: Department of Genetics, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA

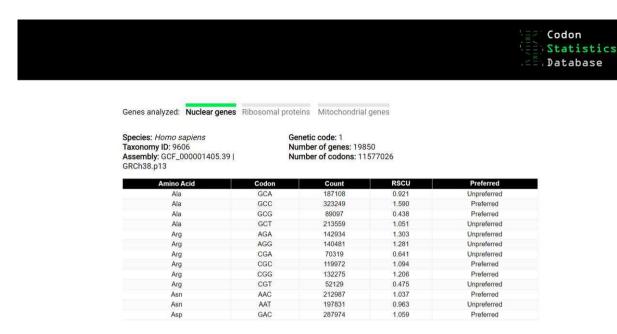


Fig. 1. Species summary. Codon statistics corresponding to all human nuclear genes are shown.

Mann–Whitney U test) were considered preferred/optimal. We then computed the $F_{\rm op}$ for each gene. The highly expressed gene set was also used as reference to compute the CAI of each gene.

The Codon Statistics Database

We present the Codon Statistics Database, an online database that contains codon usage information for all species with reference or representative genomes in RefSeq (over 15,000). The user can search for any species or taxonomic group by taxonomic ID (e.g., "9606"), scientific name (e.g., "Homo sapiens"), or common name (e.g., "human"), and select an option from a drop-down menu.

If a species is selected, the user is directed to a table that lists, for each codon, the encoded amino acid, the total count in the genome, the RSCU, and whether the codon is preferred or unpreferred (fig. 1). The user can visualize and download equivalent tables for (1) all nuclear genes (default option), (2) nuclear genes encoding ribosomal proteins (this subset is included since such proteins are often highly expressed and thus subjected to strong codon bias), (3) mitochondrial genes, and (4) chloroplast genes (if such gene sets are available in the relevant genome assembly). For viruses, only one option including all genes is available. Additionally, for each gene set, the user can download a tab-delimited file (.tsv) listing the following statistics for each gene: GC, GC3, ENC, CAI, and $F_{\rm op}$.

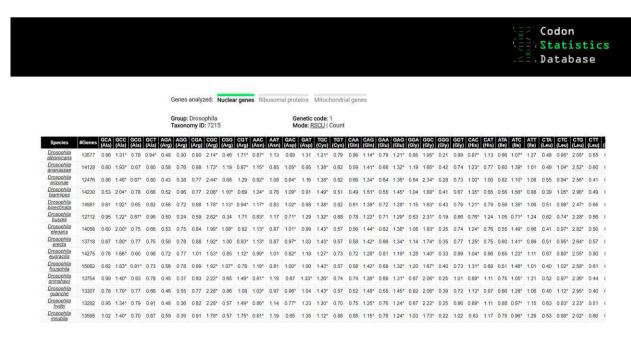


Fig. 2. Taxonomic group summary. Codon preferences for species in the genus Drosophila are shown.

If a taxonomic group with multiple species is selected (e.g., "7215," "Drosophila," or "fruit flies"), the user is presented with a table comparing all the species in the group (fig. 2). The user has the option to visualize either codon counts or RSCU values. Preferred codons in each species are marked with asterisks.

Acknowledgments

We are grateful to Alejandra Nores for assistance with web design, and to the Office of Information Technology of the University of Nevada, Reno for computational resources. This work was supported by Nevada INBRE (funded by grant P20GM103440, National Institute of General Medical Sciences, National Institutes of Health) and by the National Science Foundation (grant MCB 1818288).

Data Availability

All data used in this work are derived from the RefSeq database.

References

Athey J, Alexaki A, Osipova E, Rostovtsev A, Santana-Quintero LV, Katneni U, Simonyan V, Kimchi-Sarfaty C. 2017. A new and

- updated resource for codon usage tables. BMC Bioinform. 18: 1–10
- Christianson ML. 2005. Codon usage patterns distort phylogenies from or of DNA sequences. Am J Bot. 92:1221–1233.
- Gouy M, Gautier C. 1982. Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Res.* **10**:7055–7074.
- Gustafsson C, Govindarajan S, Minshull J. 2004. Codon bias and heterologous protein expression. *Trends Biotech.* **22**(7):346–353.
- Hershberg R, Petrov DA. 2008. Selection on codon bias. *Annu Rev Genet.* **42**:287–299.
- Hilterbrand A, Saelens J, Putonti C. 2012. CBDB: the codon bias database. *BMC Bioinform.* **13**:1–7.
- Ikemura T. 1985. Codon usage and tRNA content in unicellular and multicellular organisms. *Mol Biol Evol.* **2**:13–34.
- Nakamura Y, Gojobori T, Ikemura T. 2000. Codon usage tabulated from international DNA sequence databases: status for the year 2000. *Nucleic Acids Res.* **28**:292.
- Rota-Stabelli O, Lartillot N, Philippe H, Pisani D. 2013. Serine codon-usage bias in deep phylogenomics: pancrustacean relationships as a case study. Syst Biol. **62**(1):121–133.
- Shackelton LA, Parrish CR, Holmes EC. 2006. Evolutionary basis of codon usage and nucleotide composition bias in vertebrate DNA viruses. J Mol Evol. 62:551–563.
- Sharp PM, Emery LR, Zeng K. 2010. Forces that influence the evolution of codon bias. *Philos Trans R Soc Lond B Biol Sci.* **365**: 1203–1212
- Sharp PM, Li WH. 1987. The codon adaptation index: a measure of directional synonymous codon usage, and its potential applications. *Nucleic Acids Res.* **15**:1281–1295.
- Wright F. 1990. The "effective number of codons" used in a gene. *Gene* 87:23–29.