

An Evaluation of the Quality of Codon Deviation Coefficient as a Measurement of Codon Usage Bias



Amaral Lab

Ali Saeed¹, Sophia Liu²

¹Adlai E. Stevenson High School, Lincolnshire, IL

²Northwestern University, Dept. Chemical and Biological Engineering

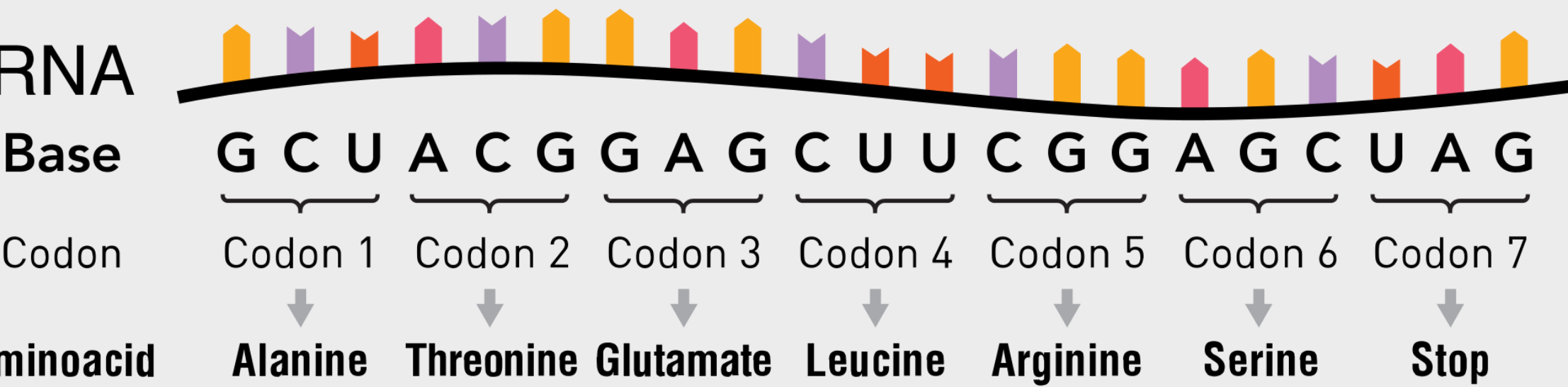
Stevenson
High School



Background

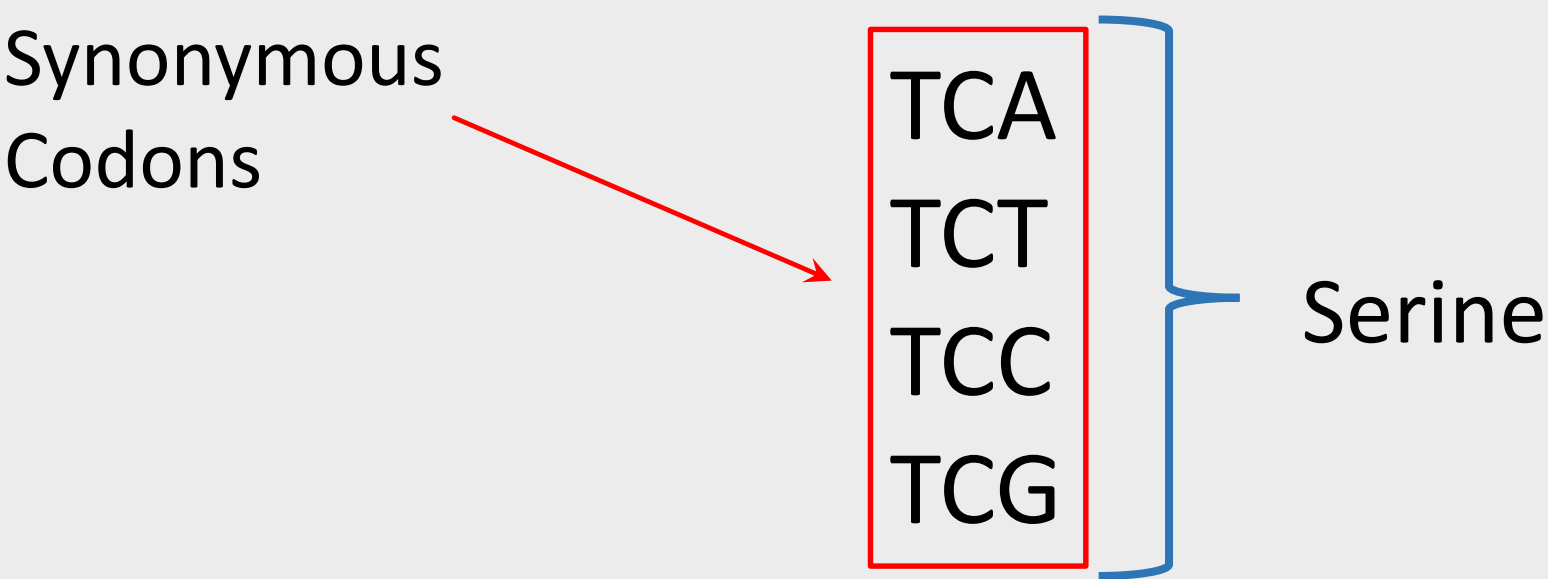
DNA is composed of four nitrogenous bases: adenine (A), thymine (T), guanine (G), and cytosine (C), which form **nucleotides**.

Sequences of three nucleotides, called **codons**, are translated into amino acids. There are 64 possible codon combinations; the 61 of which that code for amino acids are called **sense codons**.



However, there are only 20 translated amino acids. Different codons which translate into the same amino acid are called **synonymous codons**.

Some synonymous codons are used with greater frequency than others, a phenomenon known as **codon usage bias (CUB)**.



Why is Codon Usage Bias a Big Deal?

- Increased translational efficiency and accuracy
- Provides evolutionary information
- Optimize genetic engineering

Codon Deviation Coefficient

Codon Deviation Coefficient (CDC) is a metric used to calculate CUB using adenine, thymine, guanine, and cytosine content (A , T , G , and C), **GC content** (S), and **purine content** (R).

Calculating CDC

- Calculate **GC content** (S) and **purine content** (R) for each of the three codon positions (i)

$$S_i = \frac{G + C}{\text{number of sense codons}} \quad R_i = \frac{A + G}{\text{number of sense codons}}$$

- Calculate **expected nucleotide contents** (A_i , T_i , G_i , C_i) for each of the three codon positions

$$A_i = (1 - S_i)R_i, \quad T_i = (1 - S_i)(1 - R_i) \\ G_i = S_iR_i, \quad C_i = S_i(1 - R_i)$$

- Calculate the **expected codon usage** (π) for each of the 61 sense codons

For example, using the codon AGC:

$$\pi_{AGC} = \frac{A_1 G_2 C_3}{\sum x_1 y_2 z_3}, \text{ where } x, y, z \in \{A, T, G, C\}$$

- Calculate the **observed codon usage** ($\tilde{\pi}$) for each of the 61 sense codons

For example, using the codon AGC:

$$\tilde{\pi}_{AGC} = \frac{\text{number of times AGC appears in sequence}}{\text{number of codons in sequence}}$$

- Calculate the CDC from the observed and expected codon usages of each of the 61 sense codons

$$CDC = 1 - \frac{\sum_{xyz} \pi_{xyz} \times \tilde{\pi}_{xyz}}{\sqrt{\sum_{xyz} \pi_{xyz}^2 \times \tilde{\pi}_{xyz}^2}} \quad \pi_{AGC} \approx .0096, \quad \tilde{\pi}_{AGC} = 1/6$$

Pos	Pos	Pos	i	S _i	R _i	i	A _i	T _i	G _i	C _i
1	2	3	1	2/3	5/6	1	5/18	1/18	5/9	1/9
A	G	C	2	5/6	1/2	2	1/12	1/12	5/12	5/12
G	G	A	3	1/2	5/6	3	5/12	1/12	5/12	1/12
A	A	A								
C	C	A								
G	C	G								
G	C	G								

$$CDC \approx 0.490$$

CDC values range from 0 (no codon usage bias) to 1 (maximum bias)

Goal

- Design tests to determine the quality of CDC as a metric for measuring CUB

Hypotheses

How would CDC behave if it were a robust metric?

- CDC will converge to **zero** if the sequence reflects **no codon usage bias**
- CDC will converge towards **one** if the sequence reflects **maximum bias**
- CDC will be **independent of the length** of the sequence
- CDC will indicate **no correlation to the GC content** of the sequence

Method

- Use Python to generate random sequences from predetermined codon probabilities
- Measure the CDC as more codons are appended to the sequence

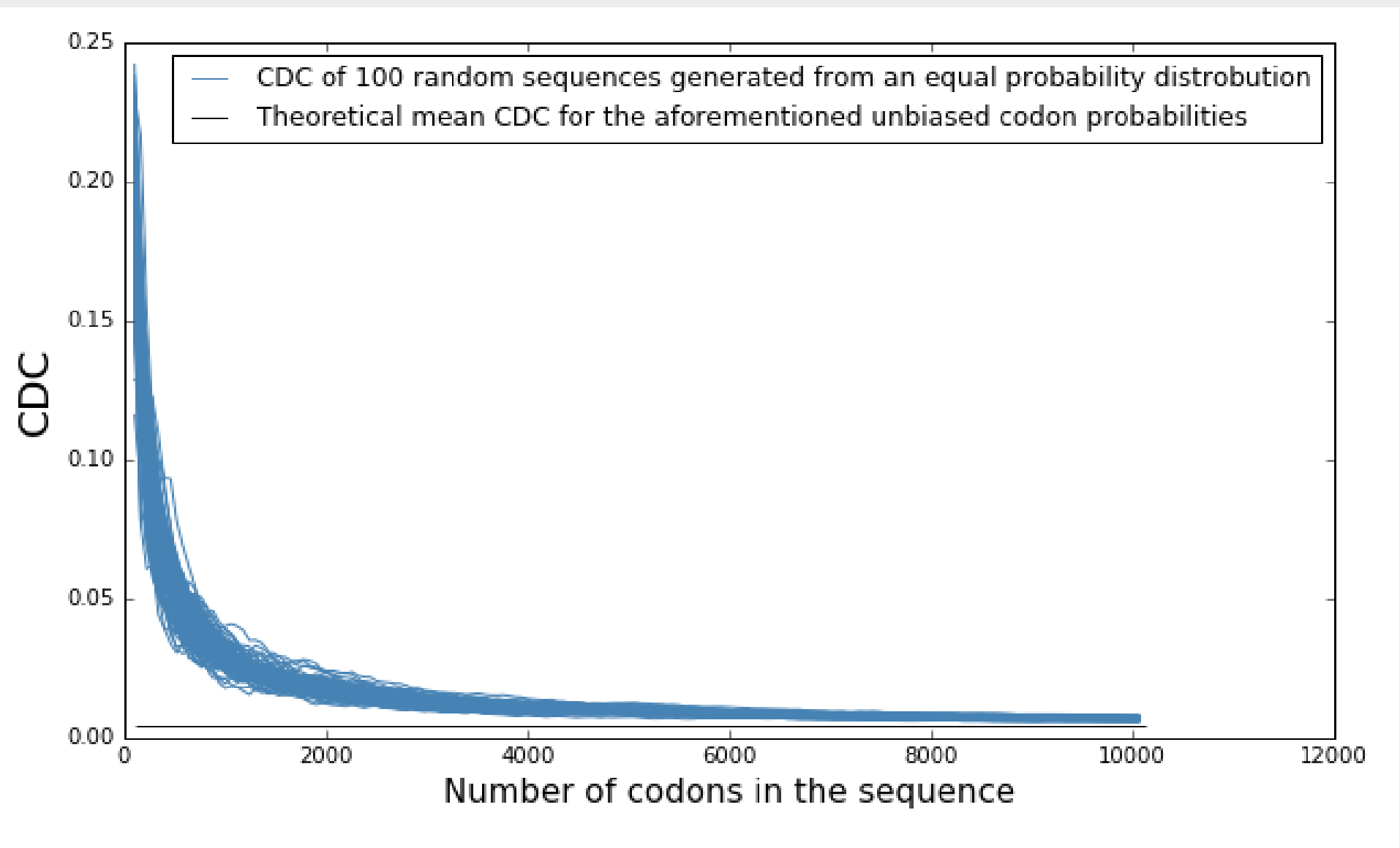
Hypothesis #1

METHOD

How can we test the hypothesis?

- CDC values will be calculated along a randomly generated sequence in which each codon has an equal probability of appearance

$$P(x) = \begin{cases} 1/61, & x \in ATT \\ 1/61, & x \in ATC \\ 1/61, & x \in ATA \\ \dots & \end{cases}$$



CDC values converged towards **zero**, reflecting no bias

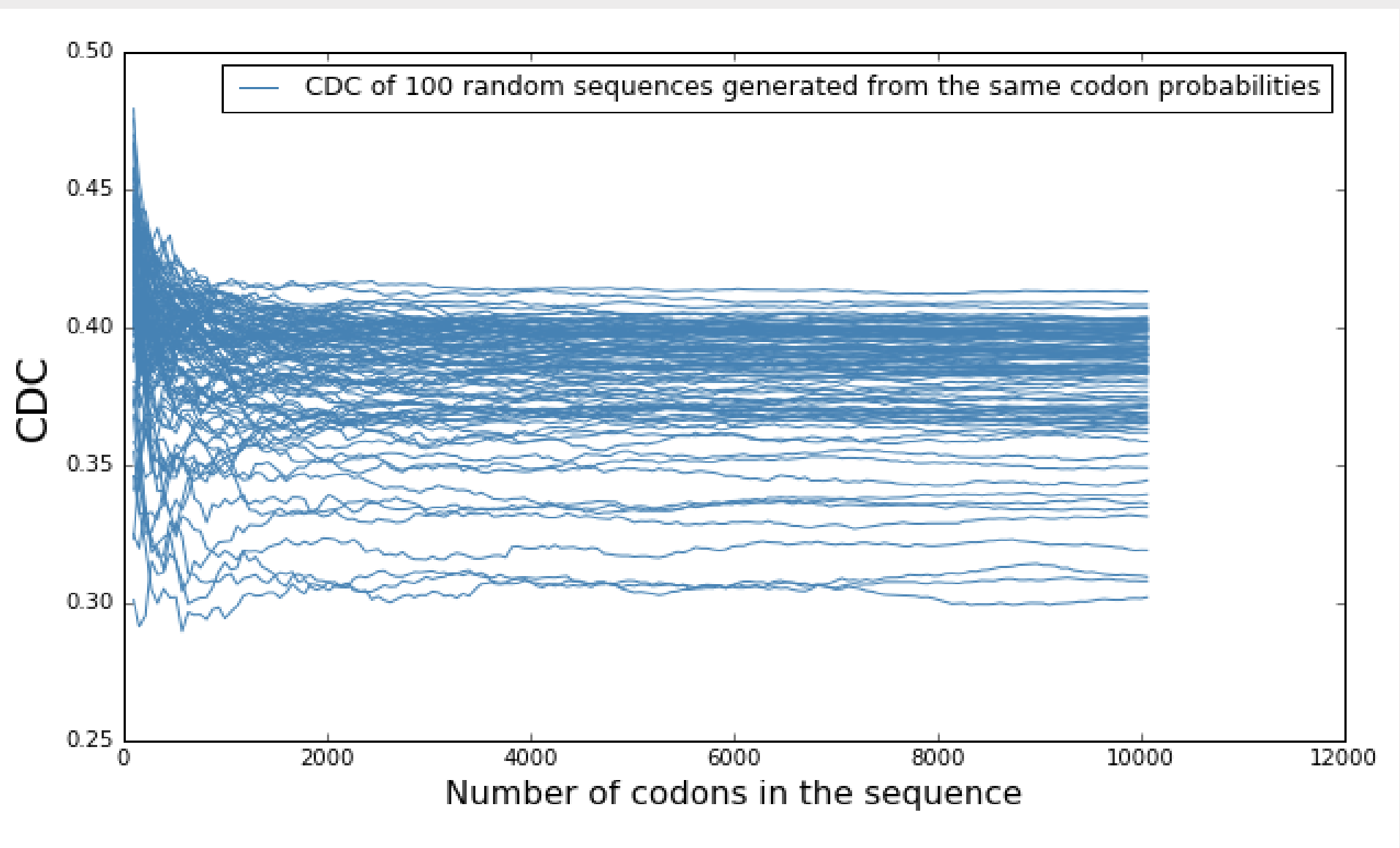
Hypothesis #2

METHOD

How can we test the hypothesis?

- CDC values will be calculated along a randomly generated sequence in which **only one random synonymous codon** per each set of 20 synonymous codon groups has an equal probability of occurrence

$$P(x) = \begin{cases} 0, & x \in ATT \\ 0.05, & x \in ATC \\ 0, & x \in ATA \\ \dots & \end{cases} \rightarrow \text{Isoleucine}$$



CDC values failed to converge towards one and reflect the highly biased sequence

CDC values failed to converge to any value

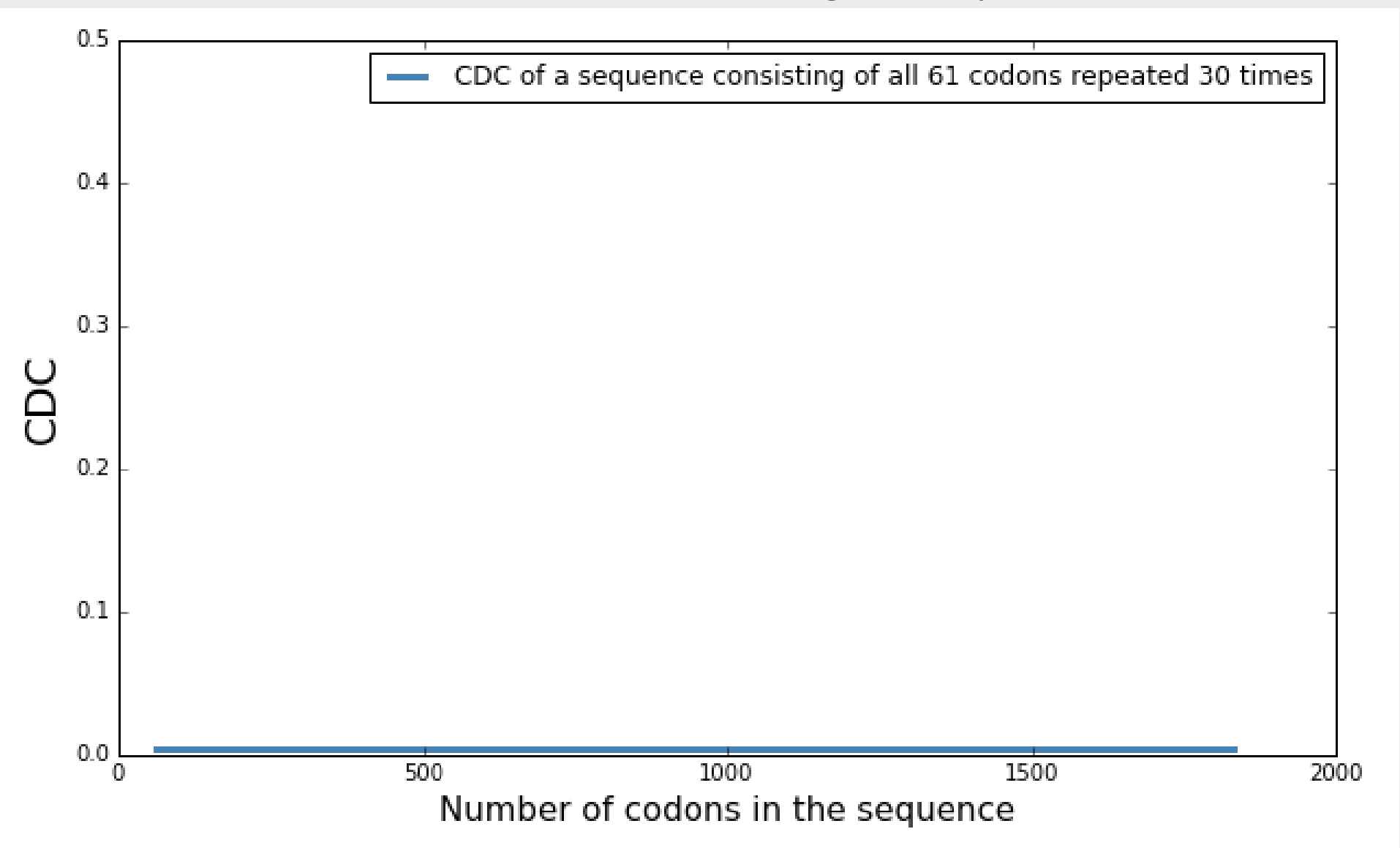
Hypothesis #3

METHOD

How can we test the hypothesis?

- CDC values will be calculated for a sequence that consists of **repeated sequences of all 61 codons**
 - CDC will be calculated at every repeat

AAA AAT AAC AAG...CCG CCC AAA AAT AAC AAG...
↑
Calculate CDC!



CDC values remained constant after every repeat, showing no correlation to length

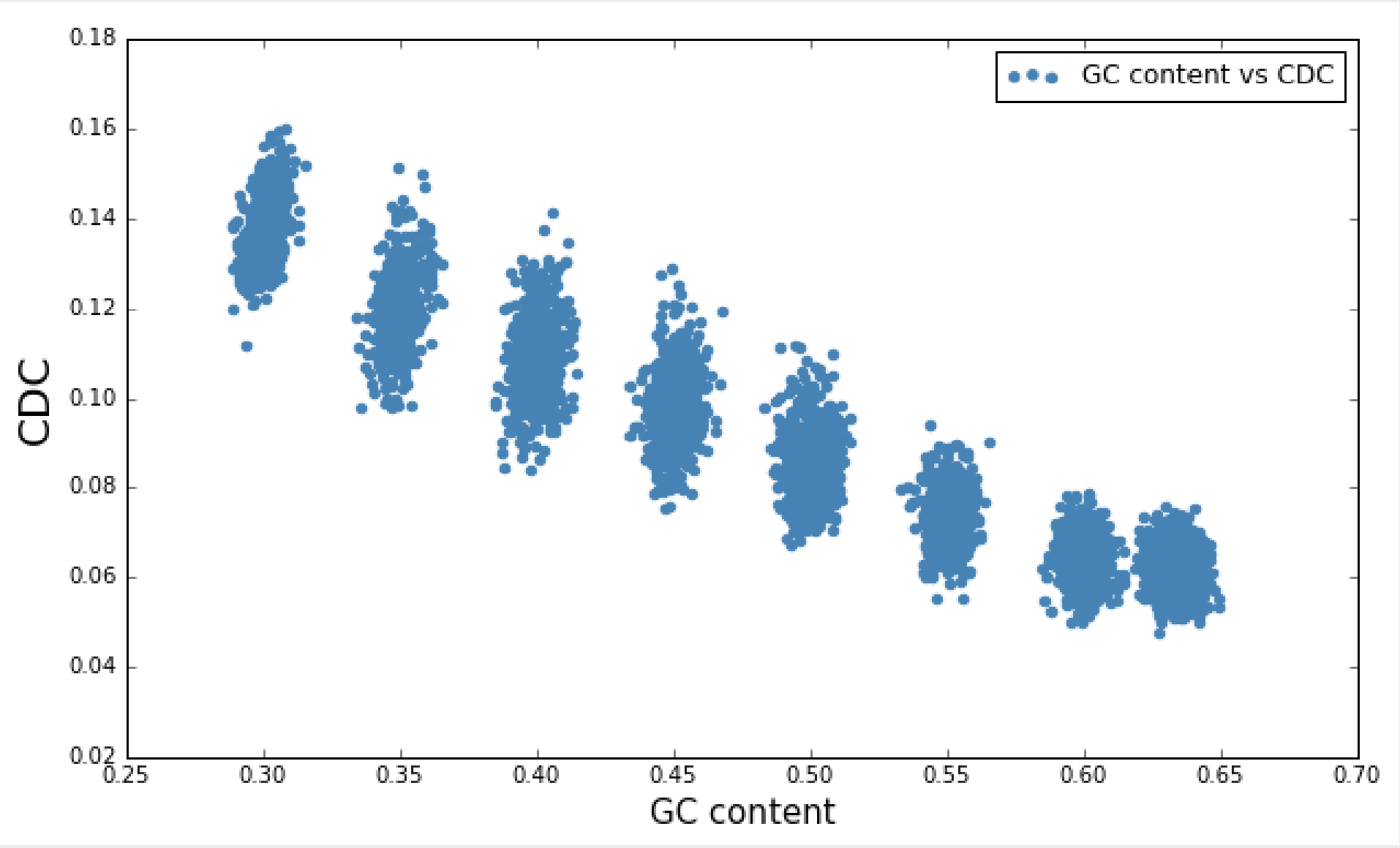
CDC values did not exactly equal zero

Hypothesis #4

METHOD

How can we test the hypothesis?

- CDC values will be calculated for completely unbiased sequences with equal amino acid usage but **variable GC contents**



CDC values demonstrated an undesirable **correlation to the GC content**

Conclusions

- CDC accurately measures unbiased sequences

When measuring sequences in which every codon is used with equal probability, the CDC value should and does converge to zero.

- CDC shows no correlation to the length of the sequence

Codon usage bias values should not demonstrate correlation to sequence length, which is an extraneous quality.

- CDC fails to accurately and precisely measure highly biased sequences

The CDC metric fails to converge to one while measuring the most highly biased situation in which only one synonymous codon is assigned a probability of occurrence.

- CDC demonstrates an undesirable correlation to the GC contents of sequences

Codon usage bias values should not demonstrate correlation towards GC content, which is an extraneous quality.

Acknowledgements

I would like to thank Dr. Luis Amaral and the members of his lab at Northwestern for allowing me to conduct research and being friendly and helpful along the way. In particular, I sincerely thank graduate student Sophia Liu for coordinating my project and assisting me through every step.

References

Zhang Z, Li J, Cui P, Ding F, Li A, Townsend JP, Yu J. 2012. Codon Deviation Coefficient: a novel measure for estimating codon usage bias and its statistical significance. BMC Bioinformatics.

Sharp PM, Bailes E, Grocock RJ, Peden JF, Sockett RE. 2005. Variation in the strength of selected codon usage bias among bacteria. Nucl. Acids Res. 33(4): 1141-1153

Akashi H. 1994. Synonymous codon usage in Drosophila melanogaster: natural selection and translational accuracy. GENETICS. 136(3): 927-935