# Exploring WORD2VEC models for Capturing the Similarity of Codon Embeddings

Anđa Denić, Jelena Pejić, Aleksandar Trokicić

University of Niš
Faculty of Sciences and Mathematics

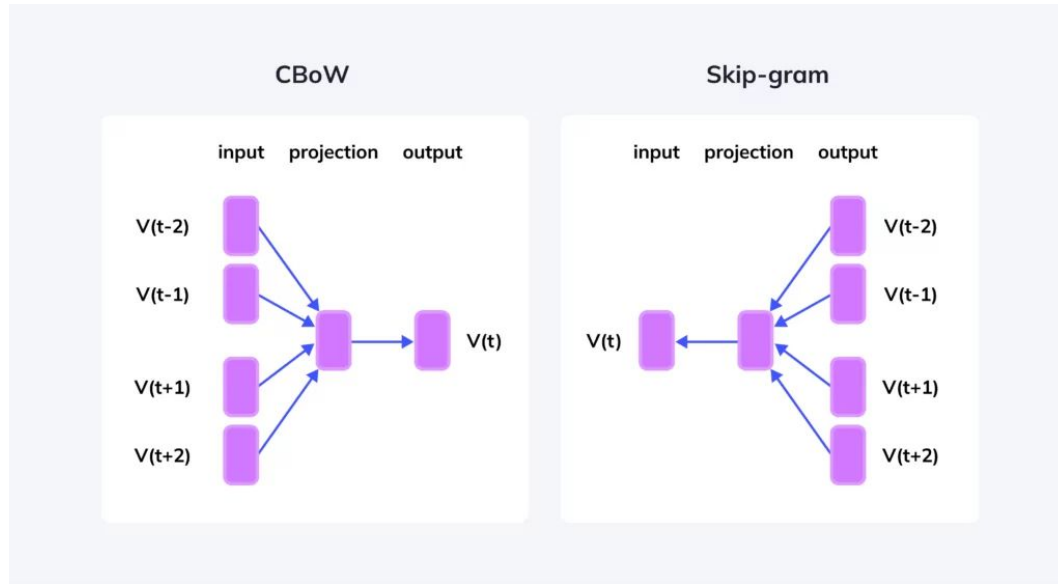# Words as high dimensional vectors

vec('cup') = [1.22, 44.444, ..., 5.78]

vec('coffee') = [3.23, 50.4, ..., 5.78]

# Word2vec

model that learns a vector representation of words



CBoW

input    projection    output

V(t-2)

V(t-1)

V(t+1)    V(t)

V(t+2)

Skip-gram

input    projection    output

V(t-2)

V(t-1)

V(t)

V(t+1)

V(t+2)

# Word2vec

**Word** meaning is based on the **surrounding words** (**context**).

Today is a **beautiful** **and** **sunny** **day** **in** **Kragujevac**, **the** **fourth** largest city in Serbia.

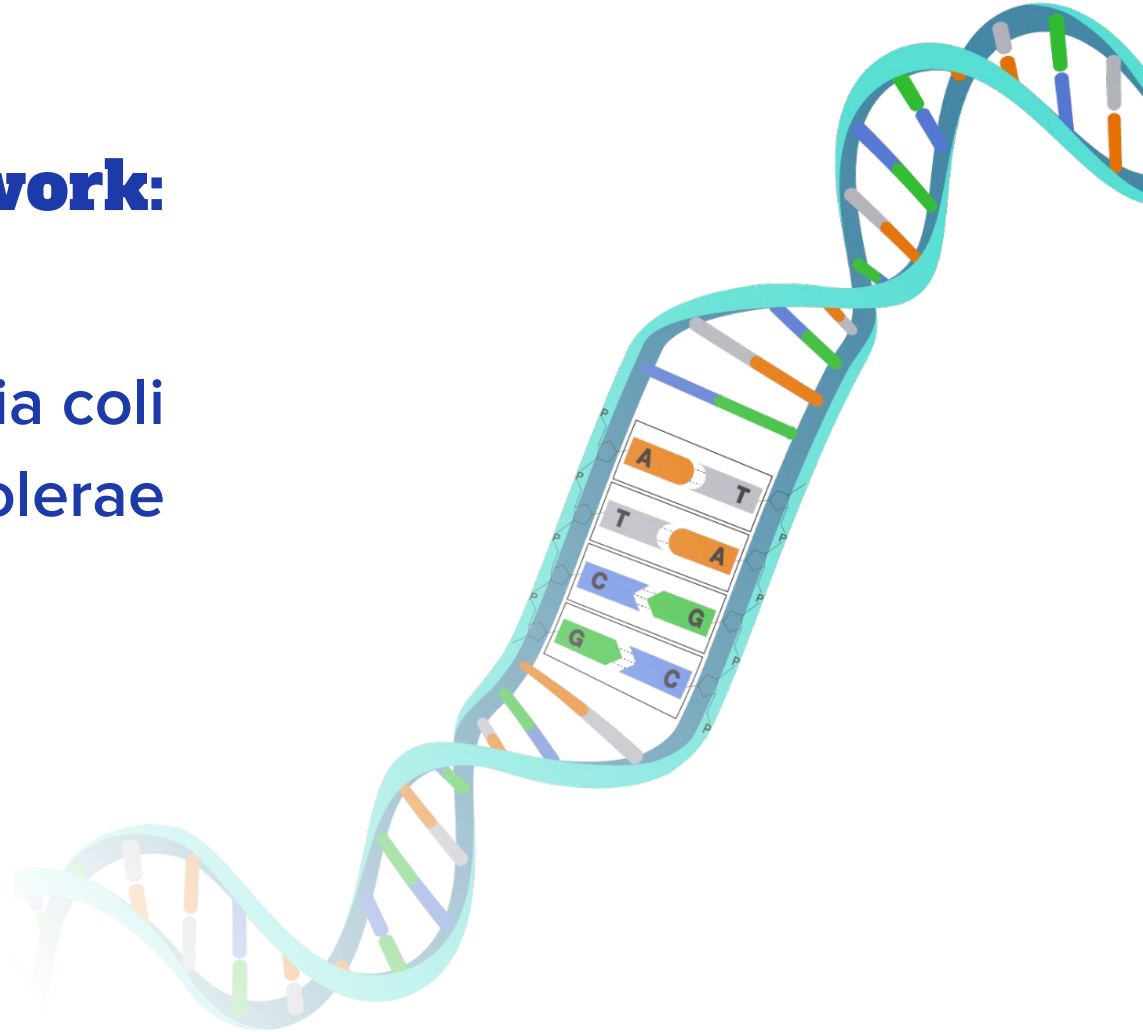# Word2vec in bioinformatics

scientific reports

OPEN **A deep learning framework combined with word embedding to identify DNA replication origins**

**Assessment of Mutation Susceptibility in DNA Sequences w Word Vectors**

**dna2vec: Consistent vector representations of variable-length k-mers**

**Word2vec based deep learning network for DNA N4-methylcytosine sites identification**

# Question:

AI can `understand` language, can it `understand` DNA? 🤔

Genome

# Genome

**DNA strands:**

Coding strand 5' A T G A T C T C G T A A 3'

Template strand 3' T A C T A G A G C A T T 5'

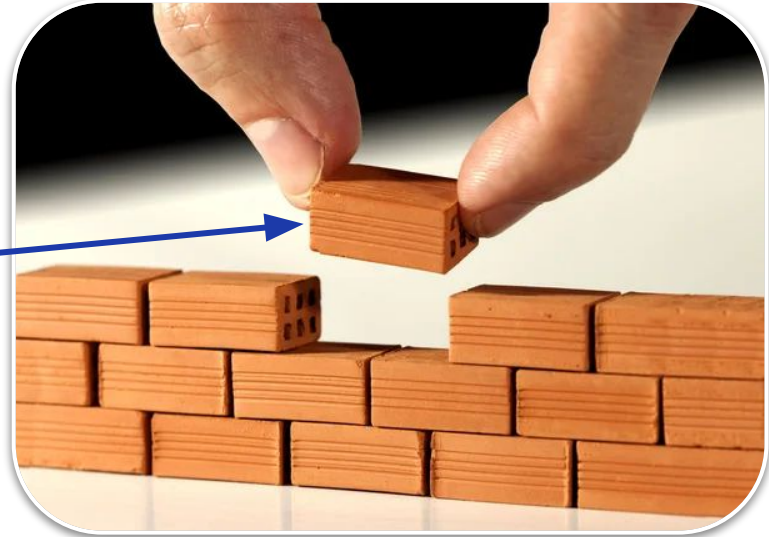**DNA sequence representation:**

Coding strand 5' A T G A T C T C G T A A 3'

# Genome
## (few millions characters)

ACAATGAGGTCACTATGTTCGAGCTCTTCAAACCGGCTGCGCATACGCAGCGGCTGCCATCCGATAAGGTGGA · CGTCTATTCACGC

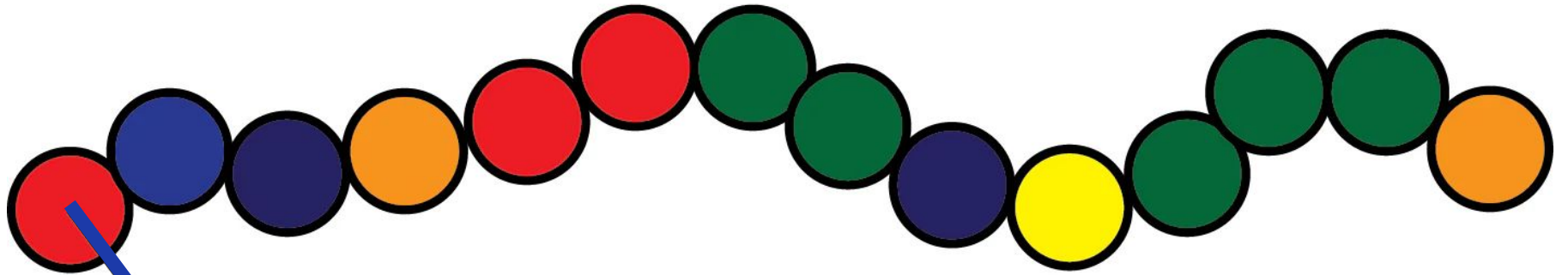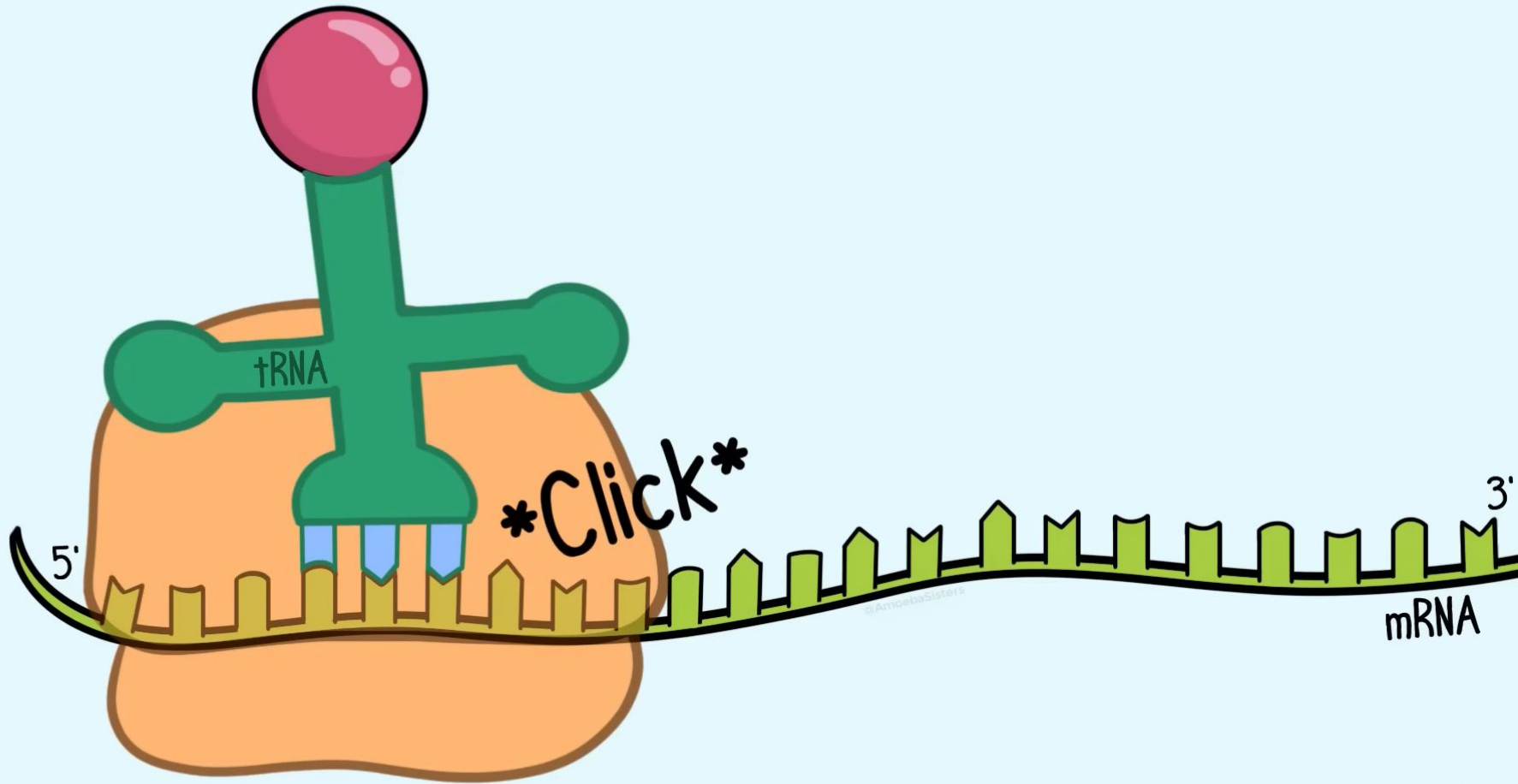First **gene**
(few thousands characters)
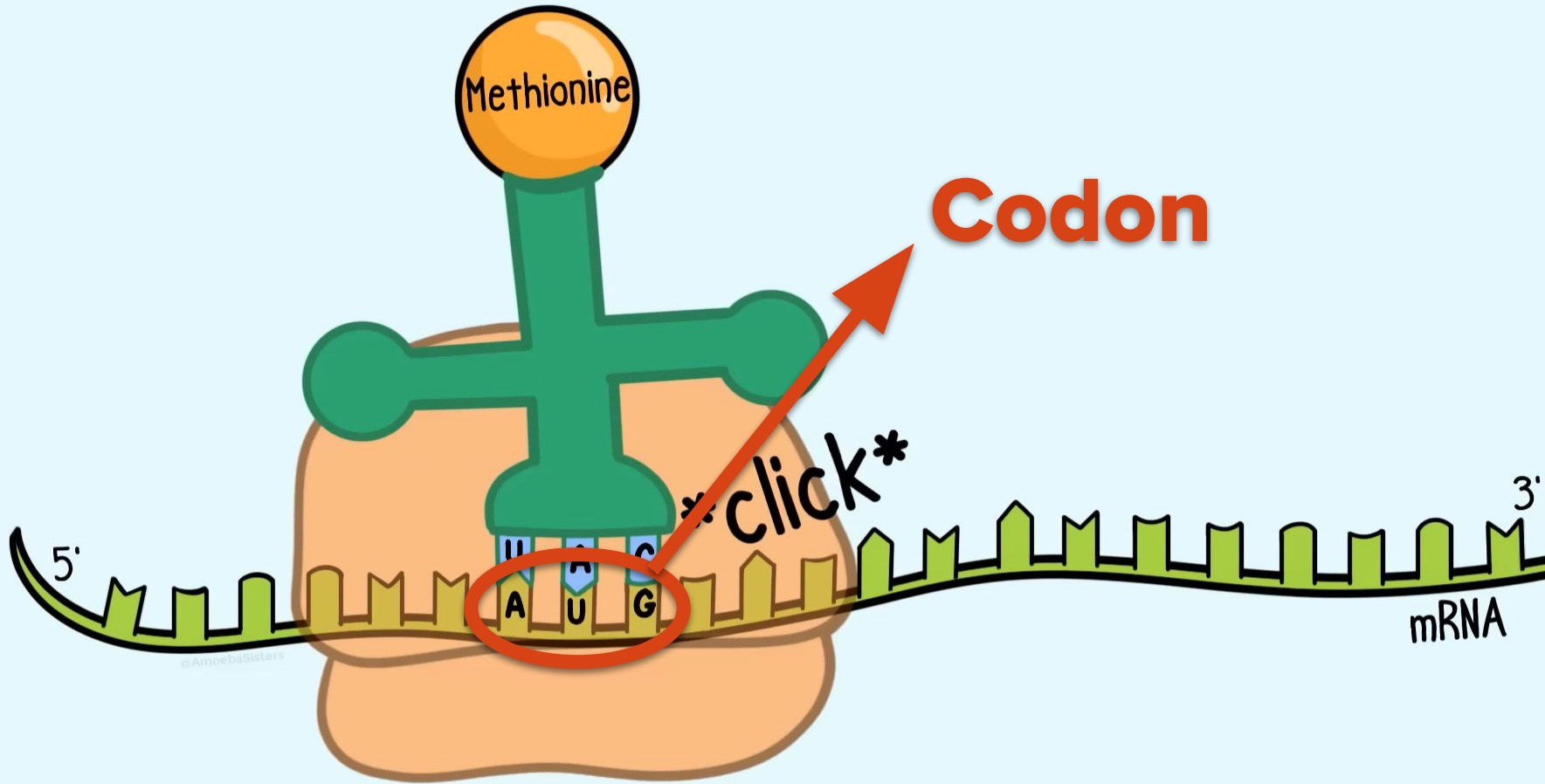
Second **gene** ...

ACAATGAGGTCACTATGTTCGAGCTCTTCAAACCGGCTGCGCATACGCAGCGGCTGCCATCCGATAAGGTGGA ·      CGTCTATTCACGC

**Genes**
code for
**proteins**

**Protein**

**Amino acid**

RNA 5'  A U G A U C  → 3'

Template strand 3'  T A C T A G A G C A T T  5'

Number of all possible codons:

____  ____  ____

# Codon chart:

# Context in NLP

**Context** of a **word 'day'**:

Today is a **beautiful** **and** **sunny** **day** **in** **Kragujevac**, **the** **fourth** largest city in Serbia.

# Context in bioinformatics?

**Context** of a **codon 'AUU'**?

UUCAACCACG AUU GCGCCGCUUU

# Context in bioinformatics?

**Context** of a **codon 'AUU'**?

k=3    UUCAACCACG AUU GCGCCGCUUU

k=4    UUCAACCACG AUU GCGCCGCUUU

k=5    UUCAACCACG AUU GCGCCGCUUU

# Context in bioinformatics?

**Context** of a **codon 'AUU'**?

$m=1$  UUCAACC ACG AUU GCG CCGCUUU

$m=3$  UUCAACCACG AUU GCGCCGCUUU

# Context in bioinformatics?

**Context** of a **codon 'AUU'**?

overlap

UUCAACCACG AUU GCGCCGCUUU

non-overlap

UUCAACCACG AUU GCGCCGCUUU

# What is the best context?

# Two word2vec architectures



CBoW

| input | projection | output |
|-------|-----------|--------|

V(t-2)
V(t-1)
V(t+1)
V(t+2)

V(t)

Skip-gram

| input | projection | output |
|-------|-----------|--------|

V(t)

V(t-2)
V(t-1)
V(t+1)
V(t+2)

# Skip gram architecture

optimize the objective of **output word being in the context of the input codon**.



One-hot
Single word '1' all other zeros
Input vector

3 nodes in hidden layer

$W^1_{NxV}$

$W^2_{NxV}$

$W^2_{NxV}$

Output layer
Predicted target context word

Actual target Word

# Training

UUCA ACC ACG AUU GCG CCG CUUU



One-hot
Single word '1' all other zeros
Input vector

$W^1_{NxV}$

3 nodes in hidden layer

$W^2_{NxV}$

$W^2_{NxV}$

Output layer
Predicted target context word

Actual target Word

# Codon Embedding



One-hot
Single word '1' all other zeros
Input vector

$W^1_{NxV}$

3 nodes in hidden layer

$W^2_{NxV}$

$W^2_{NxV}$

Output layer
Predicted target context word

Actual target Word

# Cosine similarity

is a measure between vectorised codons.

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$
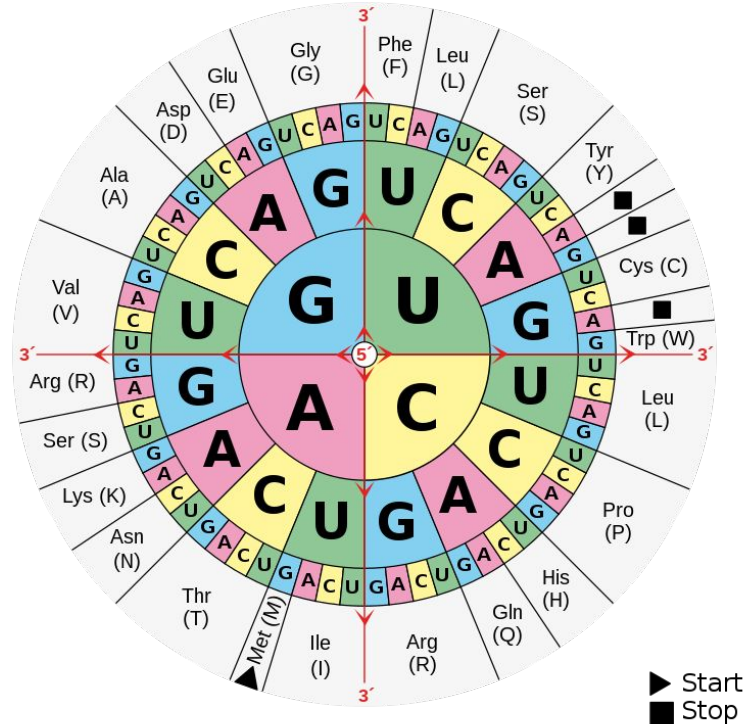
# Cosine similarity

is a measure between vectorised codons.

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$



▶ Start
■ Stop

Score of the trained model?

$$\text{Recall} = \frac{\text{\# codons from the same group in top 6 neighbors}}{\text{total number of codons in the group}}$$

AUA 0.515

CUC 0.584

CUU 0.624

CUA

CUG 0.595

GUA 0.589

UUA 0.584

$$\text{Recall ('CUA')} = \frac{4}{5} = 80\%$$

1. Find recall for each codon
2. Calculate the average recall for each model
3. Compare the results

# Results

## the best combinations of hyperparameters

| | Overlapping | | | Non-overlapping | | |
|---|---|---|---|---|---|---|
| | k=3 | k=4 | k=5 | k=3 | k=4 | k=5 |
| Vibrio cholerae | m=3<br>84.38% | m=3<br>42.66% | m=3<br>40.78% | m=3<br>45.52% | m=3<br>38.54% | m=3<br>35.56% |
| Escherichia coli | m=3<br>84.38% | m=5<br>33.18% | m=5<br>32.29% | m=10<br>42.97% | m=3<br>35.66% | m=3<br>32.23% |

# Results

## the best combinations of hyperparameters

| | Overlapping | | | Non-overlapping | | |
|---|---|---|---|---|---|---|
| | k=3 | k=4 | k=5 | k=3 | k=4 | k=5 |
| Vibrio cholerae | m=3 84.38% | m=3 42.66% | m=3 40.78% | m=3 45.52% | m=3 38.54% | m=3 35.56% |
| Escherichia coli | m=3 84.38% | m=5 33.18% | m=5 32.29% | m=10 42.97% | m=3 35.66% | m=3 32.23% |

# Results

## 3-mers as neighbors

| | Overlapping | | | | | Non-overlapping | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | m=3 | m=5 | m=10 | m=20 | m=40 | m=3 | m=5 | m=10 | m=20 | m=40 |
| Vibrio cholerae | 84.38% | 84.38% | 84.38% | 83.28% | 80.94% | 45.52% | 37.66% | 40.52% | 34.48% | 35.37% |
| Escherichia coli | 84.38% | 84.38% | 82.81% | 81.66% | 79.27% | 41.46% | 42.66% | 42.97% | 39.27% | 35.78% |

# Results

## 3-mers as neighbors, overlapping window

| | Overlapping | | | | | Non-overlapping | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | m=3 | m=5 | m=10 | m=20 | m=40 | m=3 | m=5 | m=10 | m=20 | m=40 |
| Vibrio cholerae | 84.38% | 84.38% | 84.38% | 83.28% | 80.94% | 45.52% | 37.66% | 40.52% | 34.48% | 35.37% |
| Escherichia coli | 84.38% | 84.38% | 82.81% | 81.66% | 79.27% | 41.46% | 42.66% | 42.97% | 39.27% | 35.78% |

# Results

## 3-mers as neighbors, overlapping window and smaller contect size

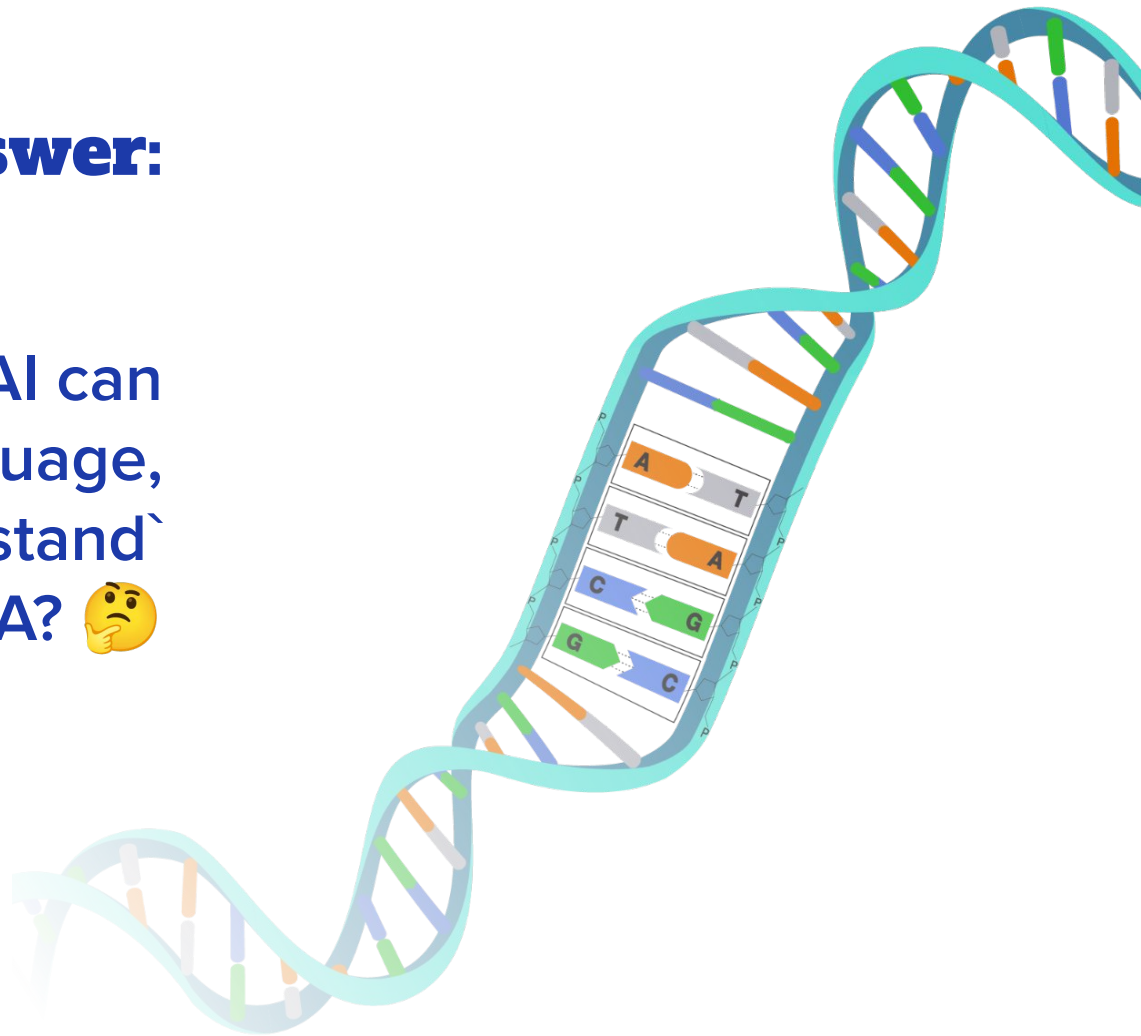| | Overlapping | | | | | Non-overlapping | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | m=3 | m=5 | m=10 | m=20 | m=40 | m=3 | m=5 | m=10 | m=20 | m=40 |
| Vibrio cholerae | 84.38% | 84.38% | 84.38% | 83.28% | 80.94% | 45.52% | 37.66% | 40.52% | 34.48% | 35.37% |
| Escherichia coli | 84.38% | 84.38% | 82.81% | 81.66% | 79.27% | 41.46% | 42.66% | 42.97% | 39.27% | 35.78% |

# Conclusion

Using:
- 3-mers as neighbors
- overlapping windows
- about 5 neighbors

gave the highest similarity among vectorized codons from the same group.

| | Overlapping | | | | | Non-overlapping | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | m=3 | m=5 | m=10 | m=20 | m=40 | m=3 | m=5 | m=10 | m=20 | m=40 |
| Vibrio cholerae | 84.38% | 84.38% | 84.38% | 83.28% | 80.94% | 45.52% | 37.66% | 40.52% | 34.48% | 35.37% |
| Escherichia coli | 84.38% | 84.38% | 82.81% | 81.66% | 79.27% | 41.46% | 42.66% | 42.97% | 39.27% | 35.78% |

# Conclusion

AI (word2vec) can learn the code written in DNA! 🎉

# Future work

🧬 Train models on different organism's genome
🧬 The best context may vary between organisms

# Potential applications

🧬 Gene function prediction
🧬 Detection of start and end of a gene

# Thank you

Question, comment?

✉ andjadenic@gmail.com