# The Language of DNA

## Can Machines Understand DNA Like English?

Justin Jia

Brinkman Lab
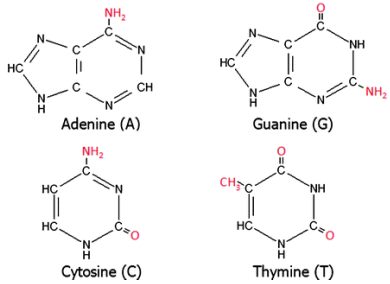
Sept 21, 2022

# Natural Language

- Any language that has evolved naturally in humans through use and repetition without conscious planning or premeditation

Alphabet

Vocabulary

Sentence

# The Language of DNA

### Alphabet



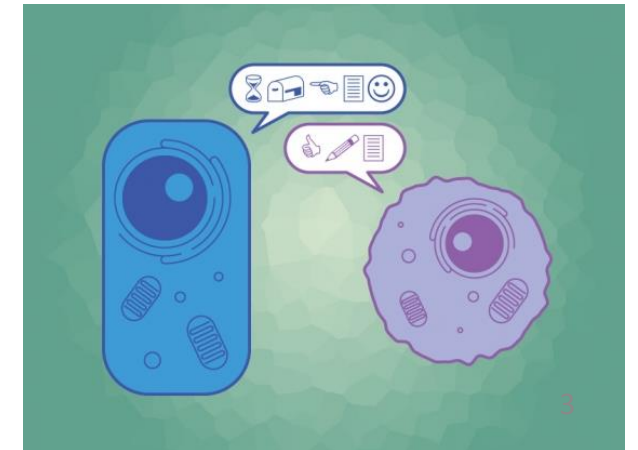Adenine (A)   Guanine (G)
Cytosine (C)   Thymine (T)

Purines/Pyrimidines

### Vocabulary

AGTTGA
AGTTGA
GTTGAG
GTTGAG
TTGAGT
TGAGTT
GAGTTG

Recognition sites
Promoters
Start codons
Etc…

### Sentence



Insulator  Enhancer      Genes        Insulator
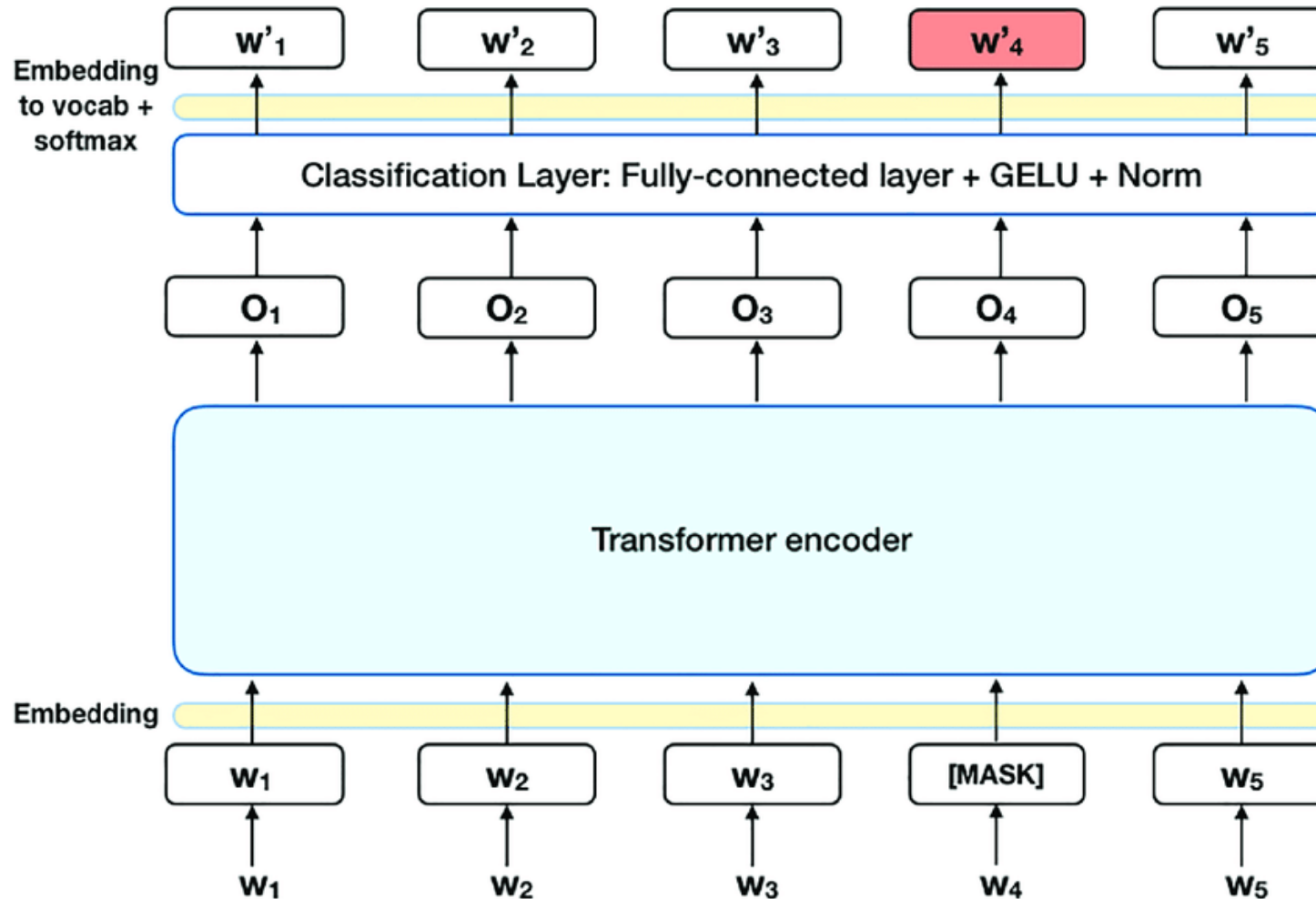
# Natural Language Processing (NLP)

- Programming computers to process and analyze large amounts of natural language data.

- The goal is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them.

# BERT: Bidirectional Encoder Representation from Transformers



- 2019 Google's research on natural language processing

- Self-learning of natural languages
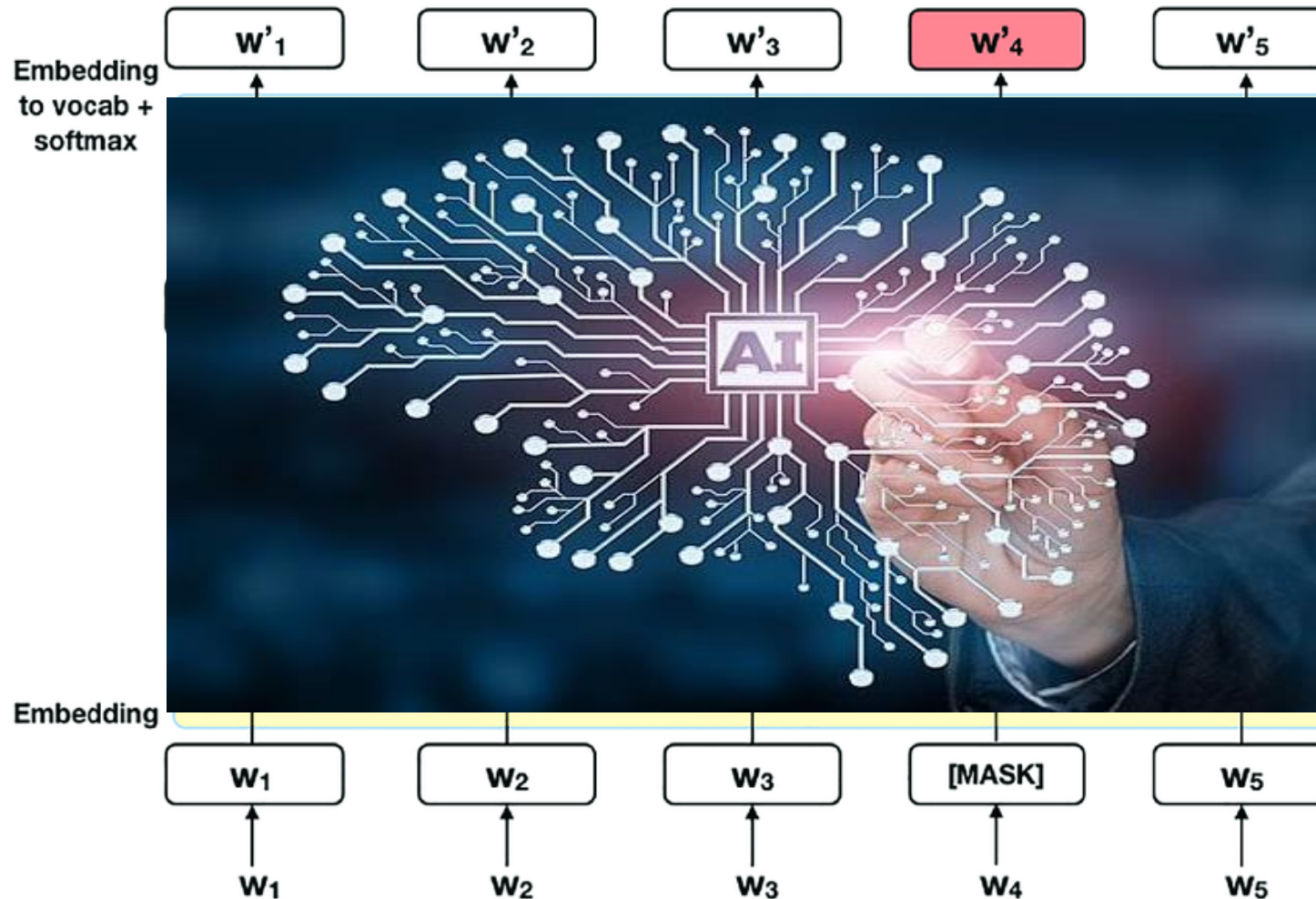- Learns contextual relations between words in a text.

# BERT: Bidirectional Encoder Representation from Transformers
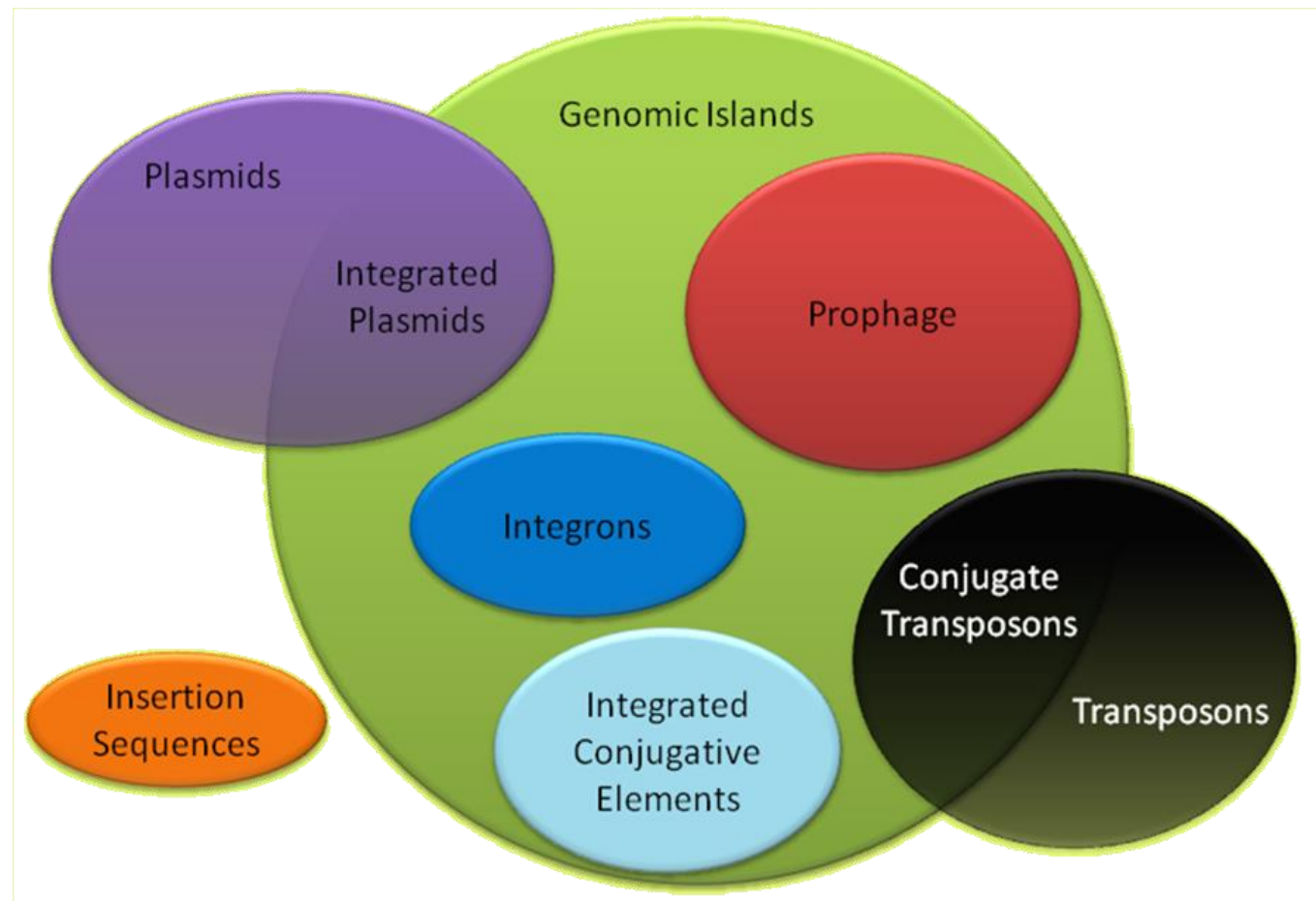


- 2019 Google's research on natural language processing

- Self-learning of natural languages
- Learns contextual relations between words in a text.

If machines can utilize NLP for human languages, can the same process be used for DNA?
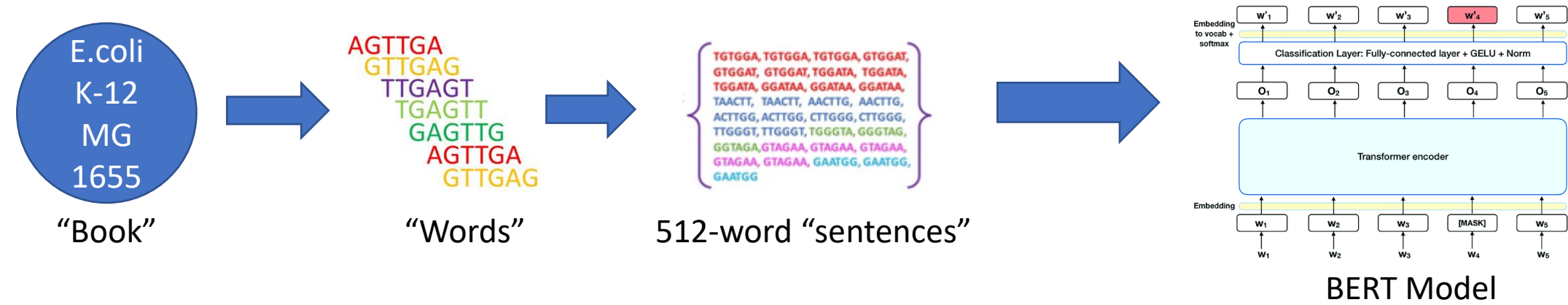
# Genomic Islands (GIs)

- Large segments (>8KB) of a genome that has evidence of horizontal origins

# Teaching the BERT Model the DNA Language

Step 1: Generalized self-learning of Bacterial DNA "words":



"Book"

"Words"

512-word "sentences"

BERT Model

# Predicting Genomic Islands

Bertelli et al. 2018. Bioinformatics.

Positive (GIs)

Literature (GIs)

Negative (non-GIs)

Step 2: fine-tuning to predict GI:



BERT Model from Step 1

3003-word "sentences"

Training:
80 GIs (all) from Literature Dataset +
80 non-GIs from Negative Dataset (Same Accession)

Evaluation:
9000 x "3003-word sentences"
Randomly chosen from Positive + Negative
dataset (non training accessions)
Never seen to model*

Visualize

*Example*

# Performance - Preliminary

## Yes. NLP works on DNA

| Method | MCC | F-Score | Accuracy | Precision | Recall |
|---|---|---|---|---|---|
| GI-BERT Unoptimized | 0.48 | 0.69 | 0.71 | 0.77 | 0.71 |
| **GI-BERT Optimized** | 0.67 | **0.82** | 0.82 | 0.82 | **0.82** |
| IslandViewer4 | **0.70** | 0.78 | **0.89** | 0.90 | 0.73 |
| SIGI-HMM | 0.35 | 0.37 | 0.73 | 0.92 | 0.26 |
| IslandPath-DIMOB v1 | 0.49 | 0.55 | 0.77 | 0.87 | 0.47 |
| MTGpick | 0.32 | 0.56 | 0.70 | 0.55 | 0.68 |
| ZislandExplorer | 0.2 | 0.23 | 0.69 | 0.85 | 0.18 |
| Islander | 0.19 | 0.20 | 0.7 | **0.97** | 0.14 |
| MSGIP | 0.15 | 0.20 | 0.68 | 0.87 | 0.16 |

# Why does it work?

- Mostly picking up patterns
  - E.g. repetitive sequences that flank GIs
  - E.g. presence of certain gene sequences in a GI.

- Need in depth evaluation.

# Acknowledgements

- Brinkman lab

- Omar Nassif