# TRANSFORMER INTUITION : A prerequisite for LLMs and Biomedical Foundational Models

ThinkBio.Ai

Vasudev R
Empid – 10073
Software Developer Trainee In Machine Learning
Machine Learning Department
Feathersoft Info IT Solutions Kochi

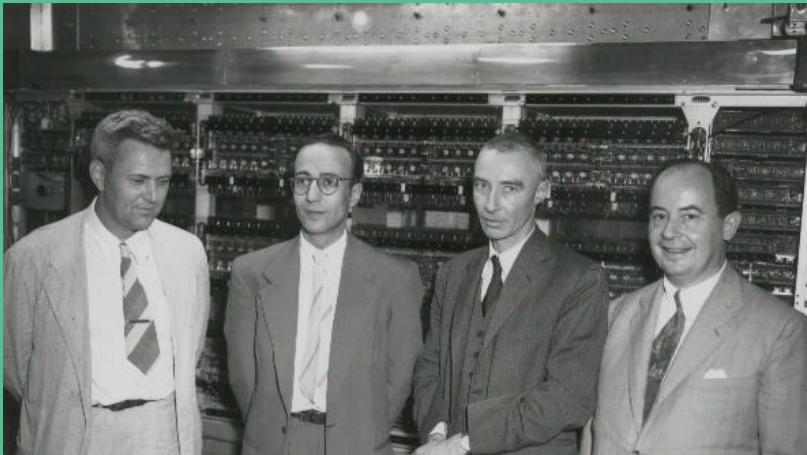# Agenda For Discussion

ThinkBio.Ai

# History of Seq2Seq Models

**What is Seq2seq ?**
- Input : a sequence (like sentence)
- Output : a sequence

**Applications include**
- Machine Translation
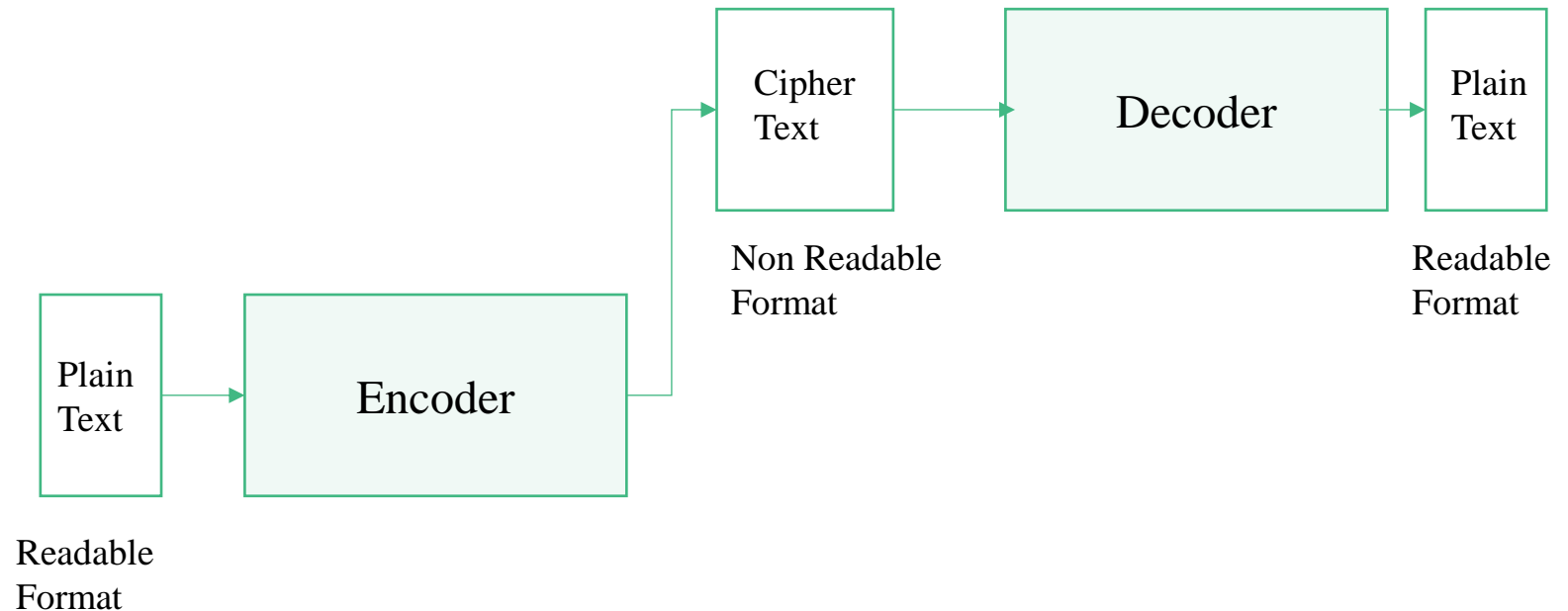- Text Summarization
- Conversational Models



Warren Weaver, Letter to Norbert Wiener, March 4, 1947

" *One naturally wonders if the **problem of translation** could conceivably be treated as a **problem in cryptography.** * "

# Cryptography Technique

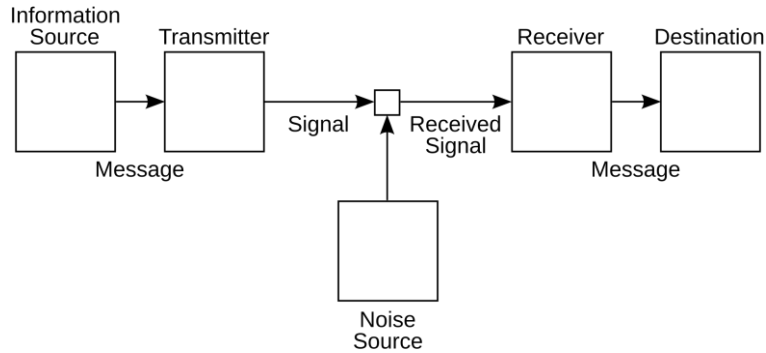High level Architecture – Encoder Decoder Model

Plain Text → Encoder → Cipher Text → Decoder → Plain Text

Readable Format

Non Readable Format

Readable Format

Two process is involved – Encryption and Decryption

**Key Features**

- Confidentiality
- Integrity
- Authentication

# Sequential to Sequential Models
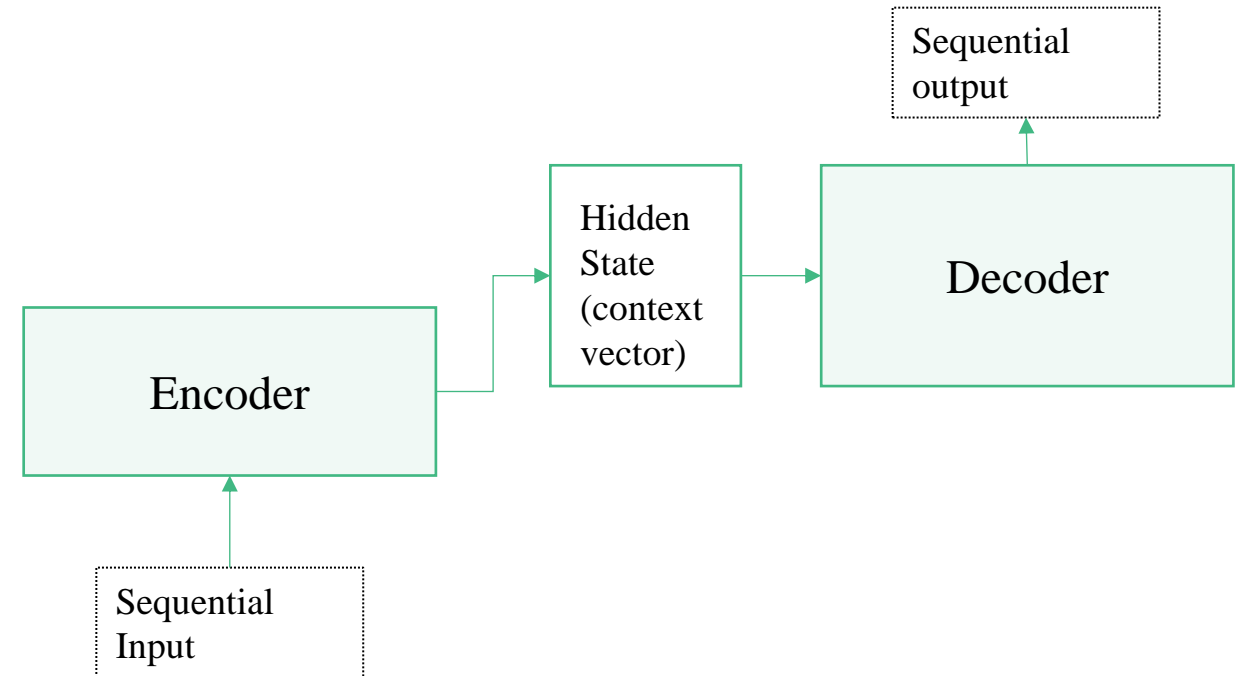
Signal Transduction Process



**Sequential to Sequential Model in Natural Lang Processing**

- Input : Sequential data (sequence of words or sentences)
  Output: Sequential data

- Before Seq2seq
  1. Statistical Methods
  2. Phrase Based Methods

  Unable to handle long term dependencies

- Seq2Seq Model
  Use RNN based Networks for input processing and as well as Output generation



- Encoder
  - Process the input sequence
  - convert into a fixed size hidden representation

- Decoder
  - use hidden state representation
  - produce target sequence

- Context vector – semantic information and other important information
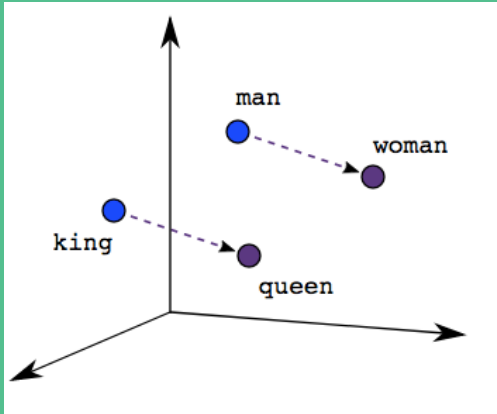- Advanced version - **Transformer**

# Agenda For Discussion

ThinkBio.Ai

# Pre-requisite : Embeddings

- Bridge between humans and computers
  Text to numbers





- Better the embedding better the model
  prediction will be

# Why Attention Mechanism ?

**Limitations with the Word Embeddings**

Where will you keep the word Cell ?



Prison

Jailer

Warden

Prisoner

College

Students

Faculty

Department

Tissues

Organs

Blood vessels

# What Attention does ?

Uses the context of the whole sentence to know what we are talking about by adjusting the embeddings values

- "The cell houses essential components like the nucleus, mitochondria, and cytoplasm..."

- Each department formed a cell to drive innovation and quick solutions

- The prisoner sat alone in his small cell.

# What about other words ?

The prisoner sat alone in his small cell.



**Single-head attention**

$$\text{Attention}(Q, K, V) = \text{softmax}_{\text{row}} \left( \frac{QK^\top}{\sqrt{d}} \right) V$$

# Is Single Head Attention Enough ?

- No its not, You need multiple attention aka multihead Attention

- More specifically multiple embeddings are created

- Among the multiple embeddings we select the embeddings with the best separation of cluster



- Select the best clusters from a set of clusters for clustering
- One embedding is created – Apply linear transformations on it
  - Shear
  - Stretch
  - Rotate
  - Combination of all

# Agenda For Discussion

ThinkBio.Ai

# Mathematical Intuition Of Attention Mechanism

**Two Process is involved**

1. Similarity capturing between every words
2. Normalization and Exponential

**Similarity measures in Euclidian space**
- Dot product
- Cosine Similarity
- Scaled Dot Product

$$\overrightarrow{A^T} = \begin{bmatrix} A_1 & A_2 & A_3 \end{bmatrix} \qquad \overrightarrow{B} = \begin{bmatrix} B_1 \\ B_2 \\ B_3 \end{bmatrix}$$

$$\begin{bmatrix} A_1 & A_2 & A_3 \end{bmatrix} \begin{bmatrix} B_1 \\ B_2 \\ B_3 \end{bmatrix} = A_1B_1 + A_2B_2 + A_3B_3 = \overrightarrow{A}.\overrightarrow{B}$$

Eg. The prisoner sat in his cell.

| | The | Prisoner | Sat | In | His | Cell | Tissue | College |
|---|---|---|---|---|---|---|---|---|
| The | 1 | 0.2 | 0 | 0 | 0 | 0 | 0 | 0 |
| Prisoner | 0.2 | 1 | 0 | 0 | 0 | 0.75 | 0 | 0 |
| Sat | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| In | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| The | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Cell | 0 | 0.75 | 0 | 0 | 0 | 1 | 0 | 0 |
| Tissue | 0 | 0 | 0 | 0 | 0 | 0.8 | 1 | 0 |
| College | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 1 |

The =

Prisoner = 1 * Prisoner + 0.75 * Cell

Sat =

in =

his =

Cell = 0.75 * Prisoner + 1 * Cell

# Normalization and exponentiation

Eg. The prisoner sat in his  cell.

| | The | Prisoner | Sat | In | His | Cell | Tissue | College |
|---|---|---|---|---|---|---|---|---|
| The | 1 | 0.2 | 0 | 0 | 0 | 0 | 0 | 0 |
| Prisoner | 0.2 | 1 | 0 | 0 | 0 | 0.75 | 0 | 0 |
| Sat | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| In | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| The | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Cell | 0 | 0.75 | 0 | 0 | 0 | 1 | 0 | 0 |
| Tissue | 0 | 0 | 0 | 0 | 0 | 0.8 | 1 | 0 |
| College | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 1 |

The =

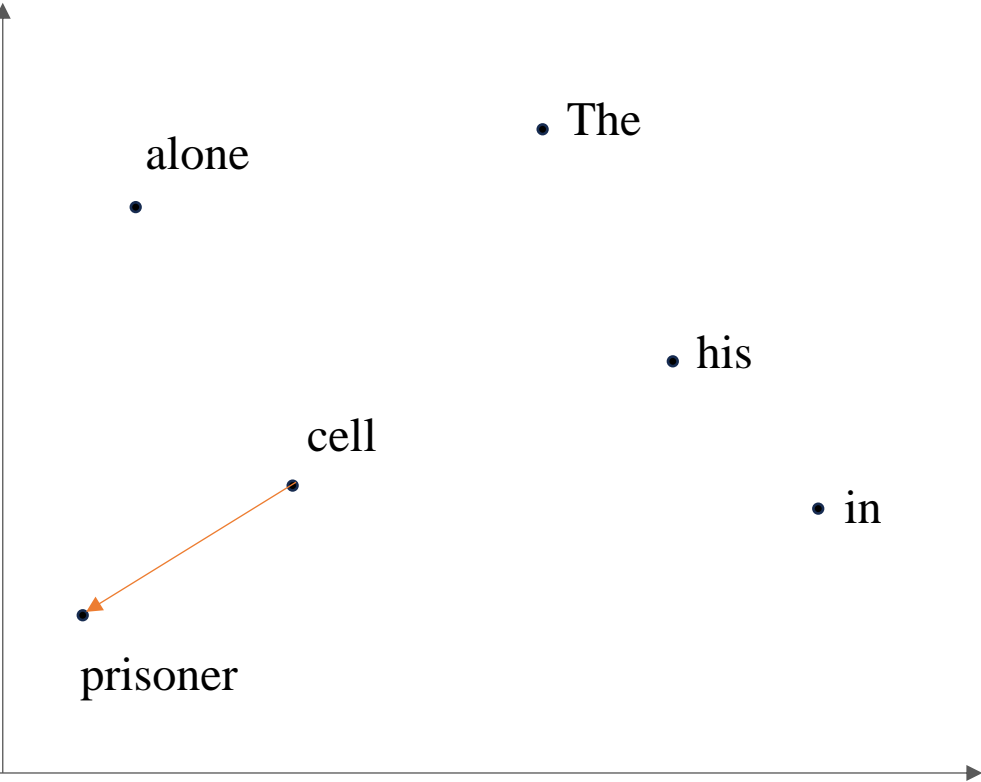Prisoner = 1 * Prisoner + 0.75 * Cell

Sat =

in =

his =

Cell = 0.75 * Prisoner + 1 * Cell

Prisoner = $\dfrac{1 * \text{Prisoner} + 0.75 * \text{Cell}}{1 + 0.75}$ = 0.58 Prisoner + 0.42 Cell
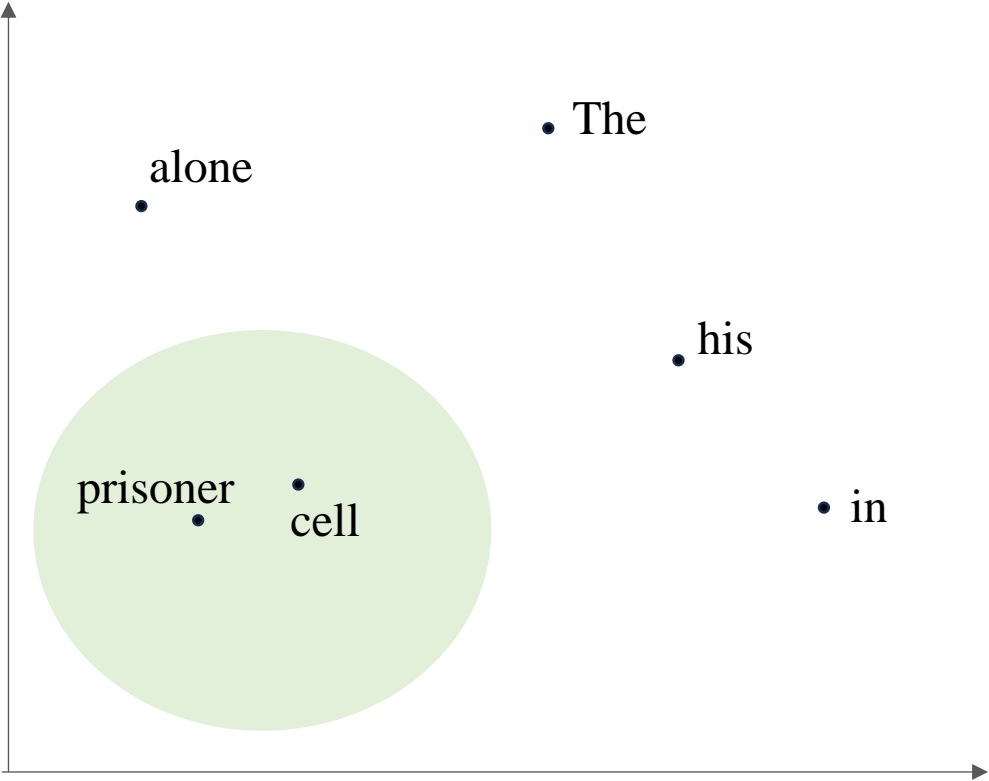
Prisoner = $\dfrac{(e \wedge 1) * \text{Prisoner} + (e \wedge 0.75) * \text{Cell}}{e \wedge (1 + 0.75)}$ = 0.58 Prisoner + 0.42 Cell

Eg. The prisoner sat in his cell.



Cell = 0.58 Prisoner + 0.42 Cell

Prisoner = 0.42 Prisoner + 0.58 Cell

# Agenda For Discussion

ThinkBio.Ai

# High level Architecture



# Whats inside Encoder and Decoder ?

# The Transformer Timeline

**1986**

### Word Embeddings

Hinton proposed the idea of "learning distributed representation of words"
- Representing semantics of a word by mapping it into a higher dimension space.
- Such that words that are together have similar meaning.

**2013**

### Word2vec

- This was a breakthrough in NLP
- Embeddings generated were called Neural Embeddings
- These embedding were of lower dimensions also.

**2017**

### Transformer ( Attention )

- Update the embedding values
- Updated values will be able to capture wrt to context of the sentence.

**Thank You** For your Valuable Time.