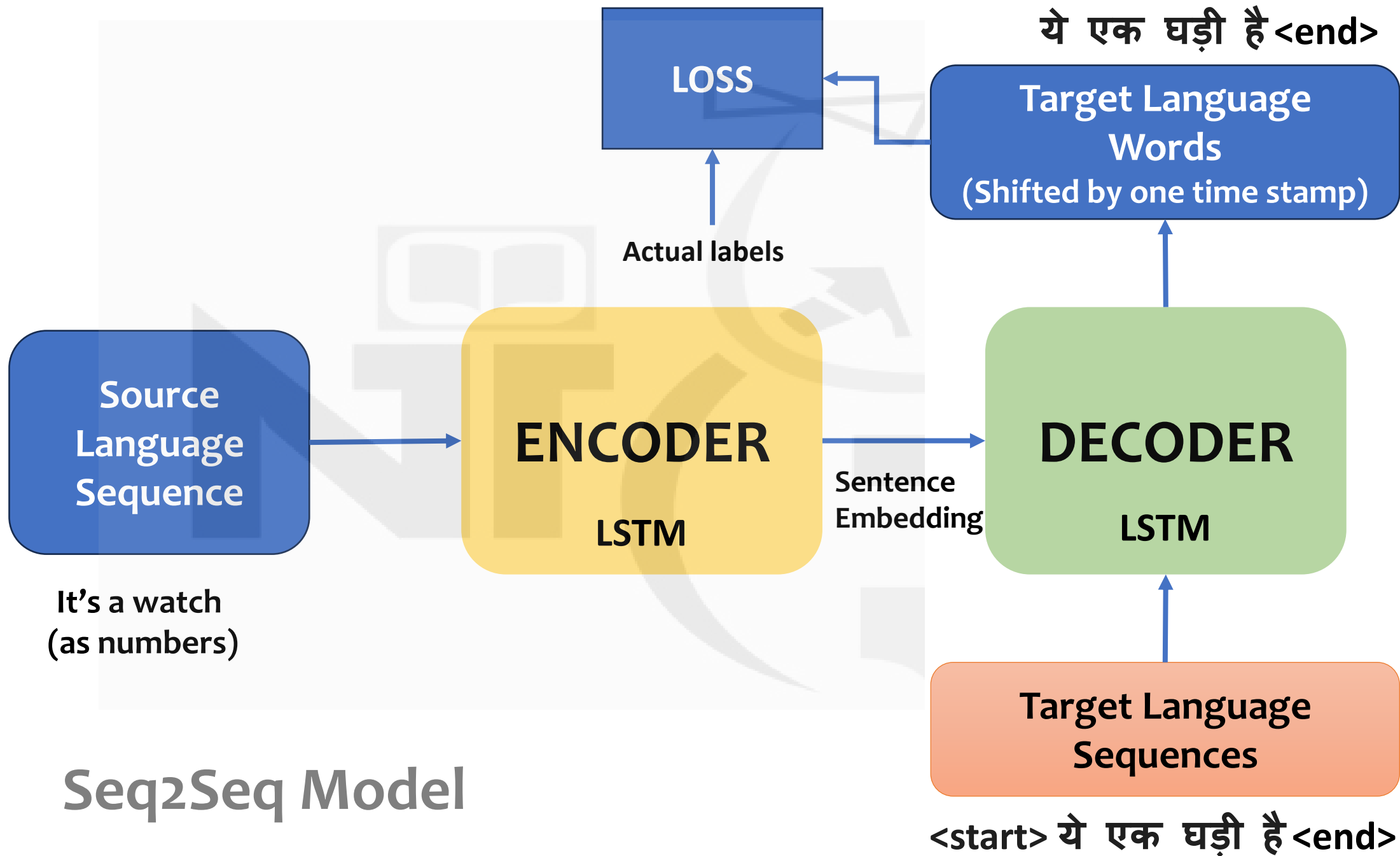


The background features a large, faint watermark of the NITCE logo. The logo consists of the letters 'NITCE' in a stylized font, with a circular emblem containing a book and a torch above the 'E'.

# Transformers

**MUKESH KUMAR**



# Key issue with LSTM

- What is the major issue with LSTM
- It cannot take advantage of parallel processing available in GPU
- And also, large data take more time because it process data word by word

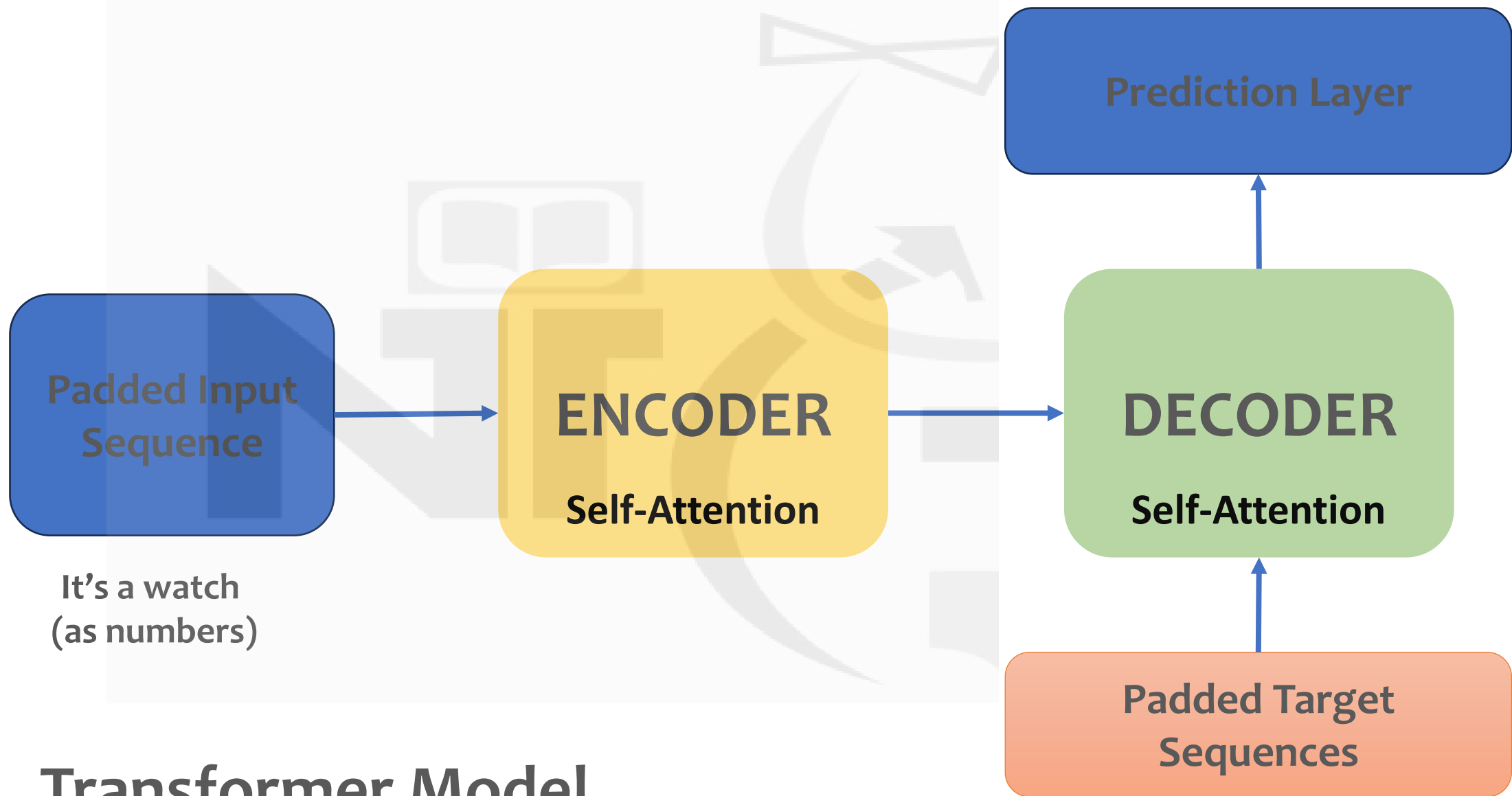
# Attention is all you need

<https://arxiv.org/abs/1706.03762>

**Moving away from LSTM**

# Attention

- Came around 2017
- Build without LSTM



# Transformer Model

# Self Attention

- Self attention layer is a dense layer

How does Encoder understand sequence in transformer model?

It tries to understand sequence using **Self Attention**



# Self-Attention Vs LSTM

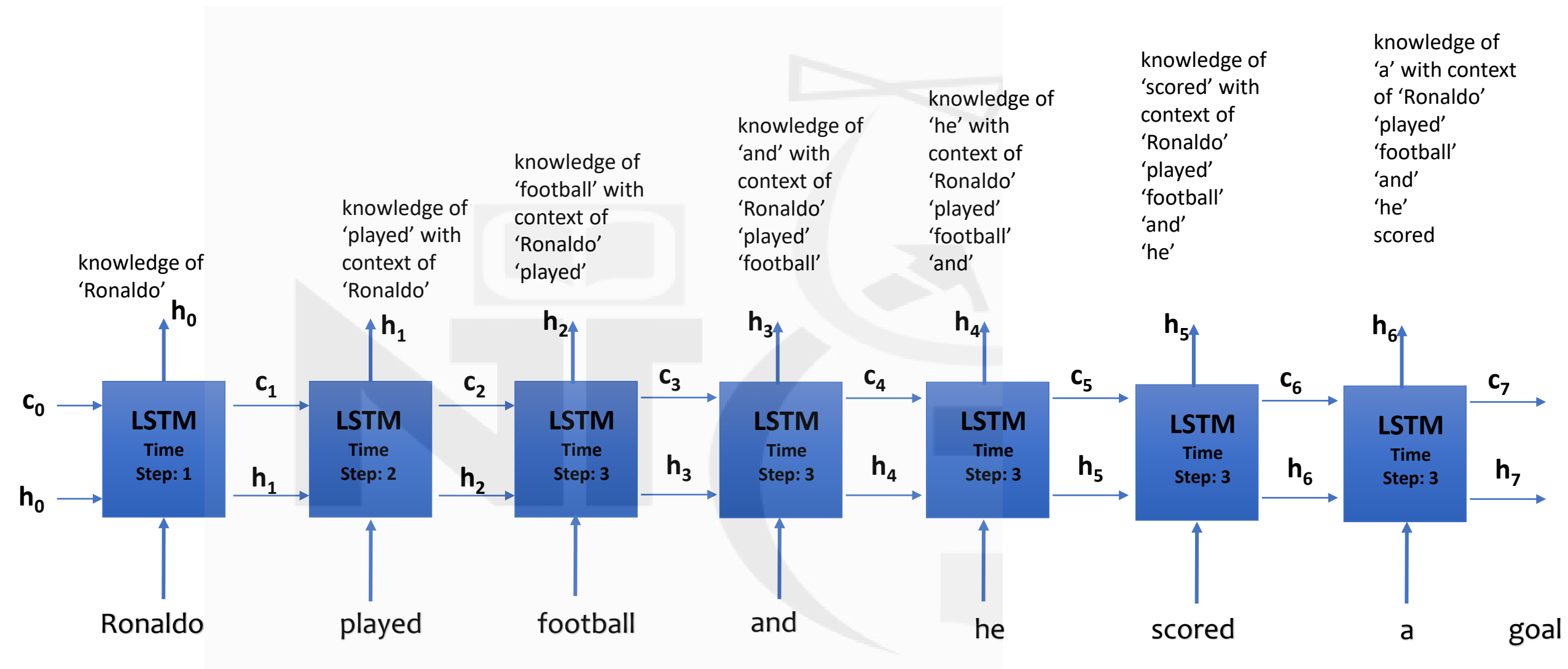
- Let's consider an example:

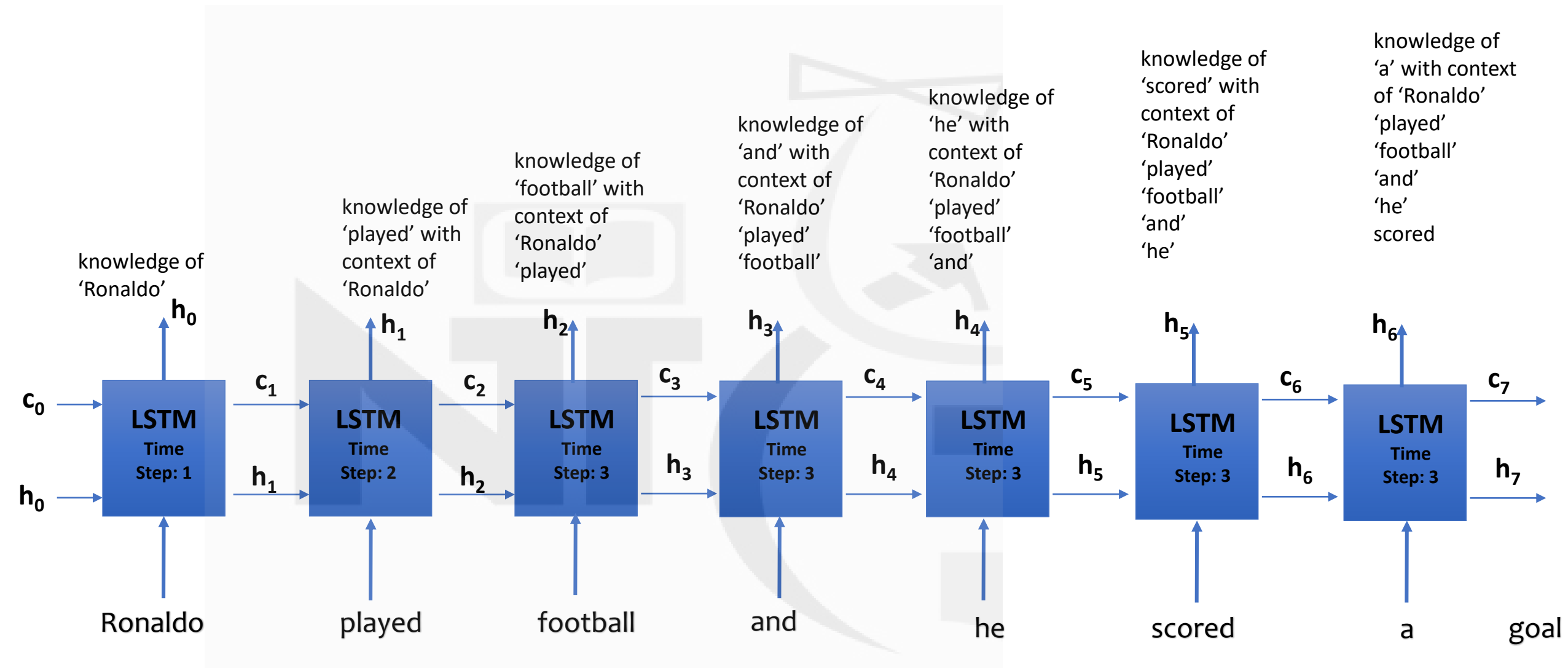
Ronaldo played football and he scored a goal

**How do we understand  
this sentence?**

LSTM







LSTM builds a memory (hidden and cell state) as it's new data elements in the sequence. This allows LSTM to understand individual data elements e.g words based on the context.

LSTM looks at Sequence one step at a time!!!

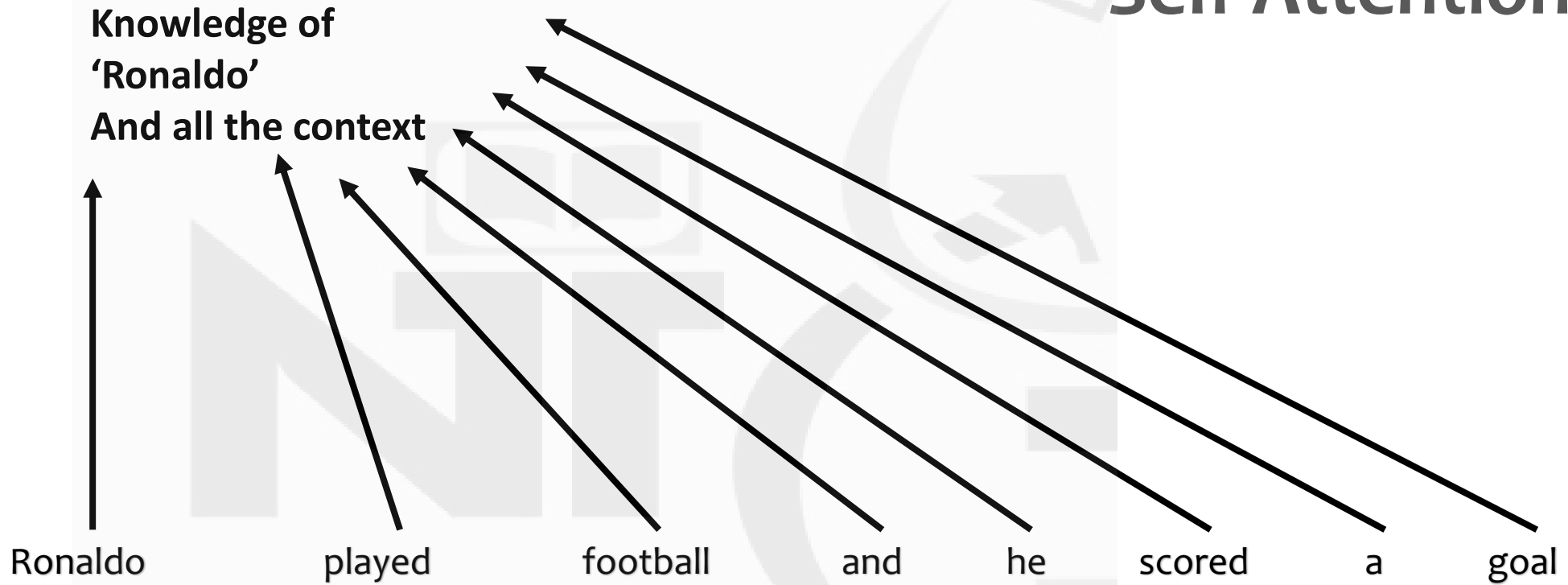
How do we get 'hidden State' type on information in  
Transformer Encoder

**We will need to accomplish the  
same without 'memory'**



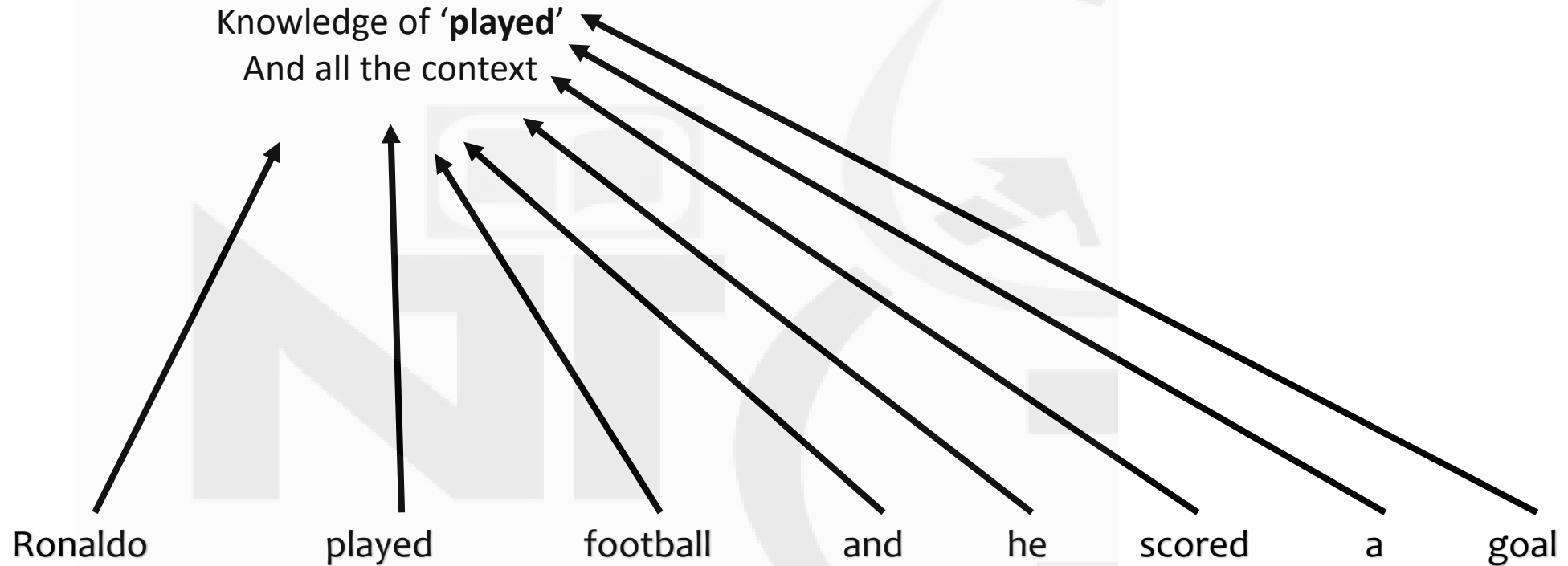
# Self-Attention

## Self Attention



To create encoding for word 'Ronaldo'  
Self-Attention looks at all the words in the  
sequence

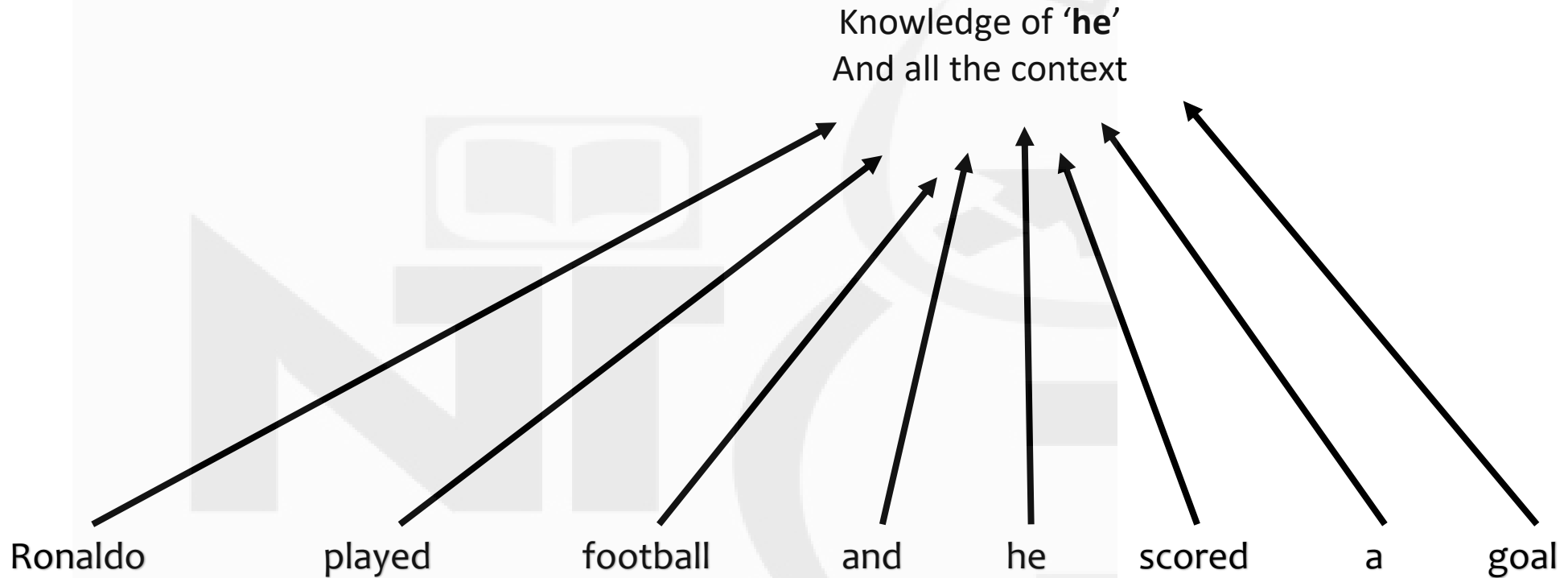
# Self Attention



To create encoding for word 'played'  
Self-Attention looks at all the words in the  
sequence



# Self Attention



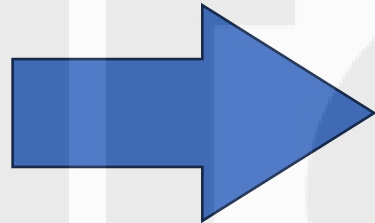
**Great news encoding for each word in sequence can be done in parallel**

# Understanding Encoder

Transformer Model

**ENCODER**

With Dense layers

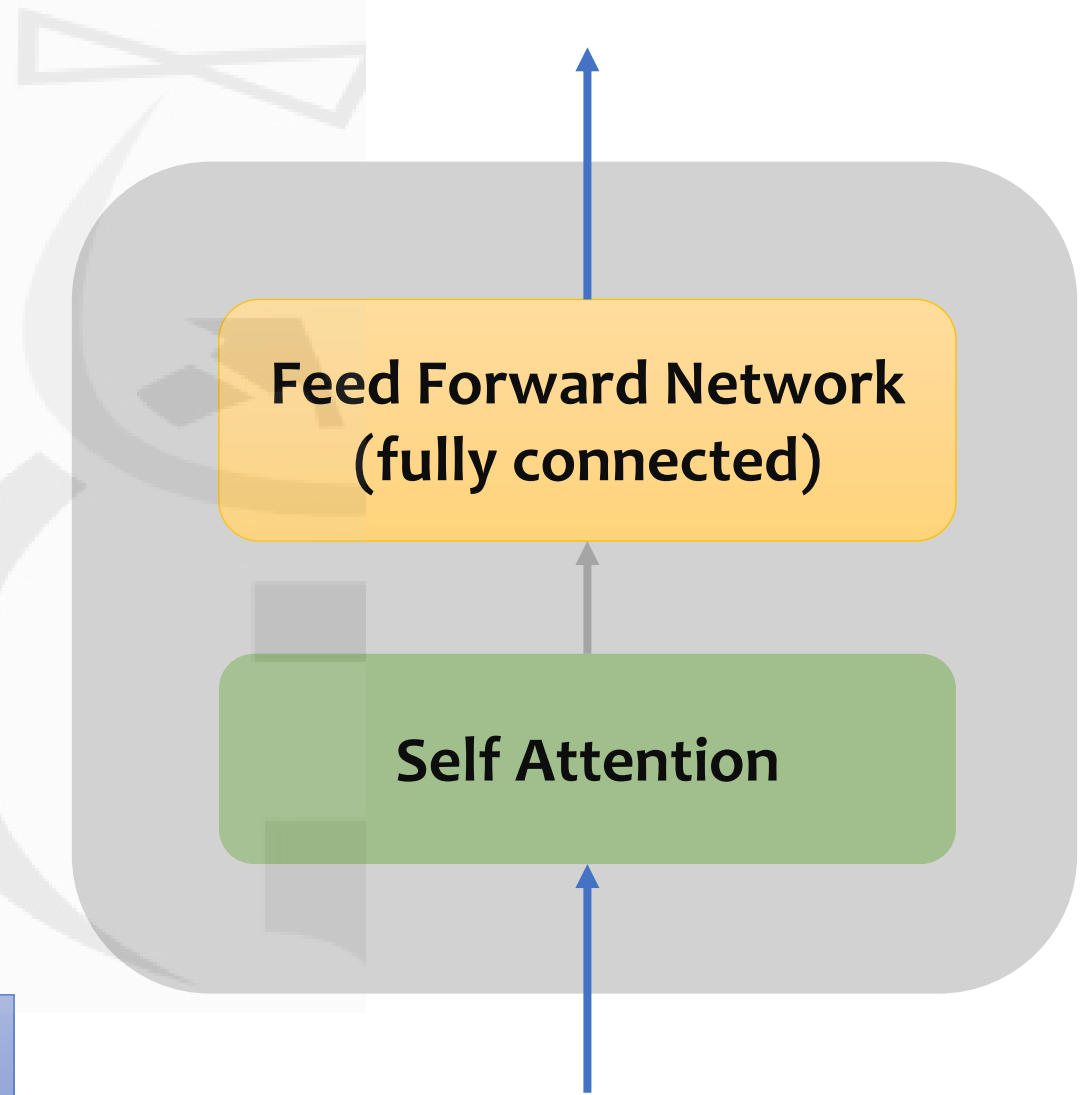


**Feed Forward Network  
(fully connected)**

**Self Attention**

**Embeddings**

In Transformer model, Encoder is made up of two types of layers as shown on right





# Implementing Self Attention

## Implementing Self Attention



Input Sentence

Who

is

Ronaldo

## Implementing Self Attention

Embedding

$x_1$



$x_2$



$x_3$



Embedding is created  
for each word

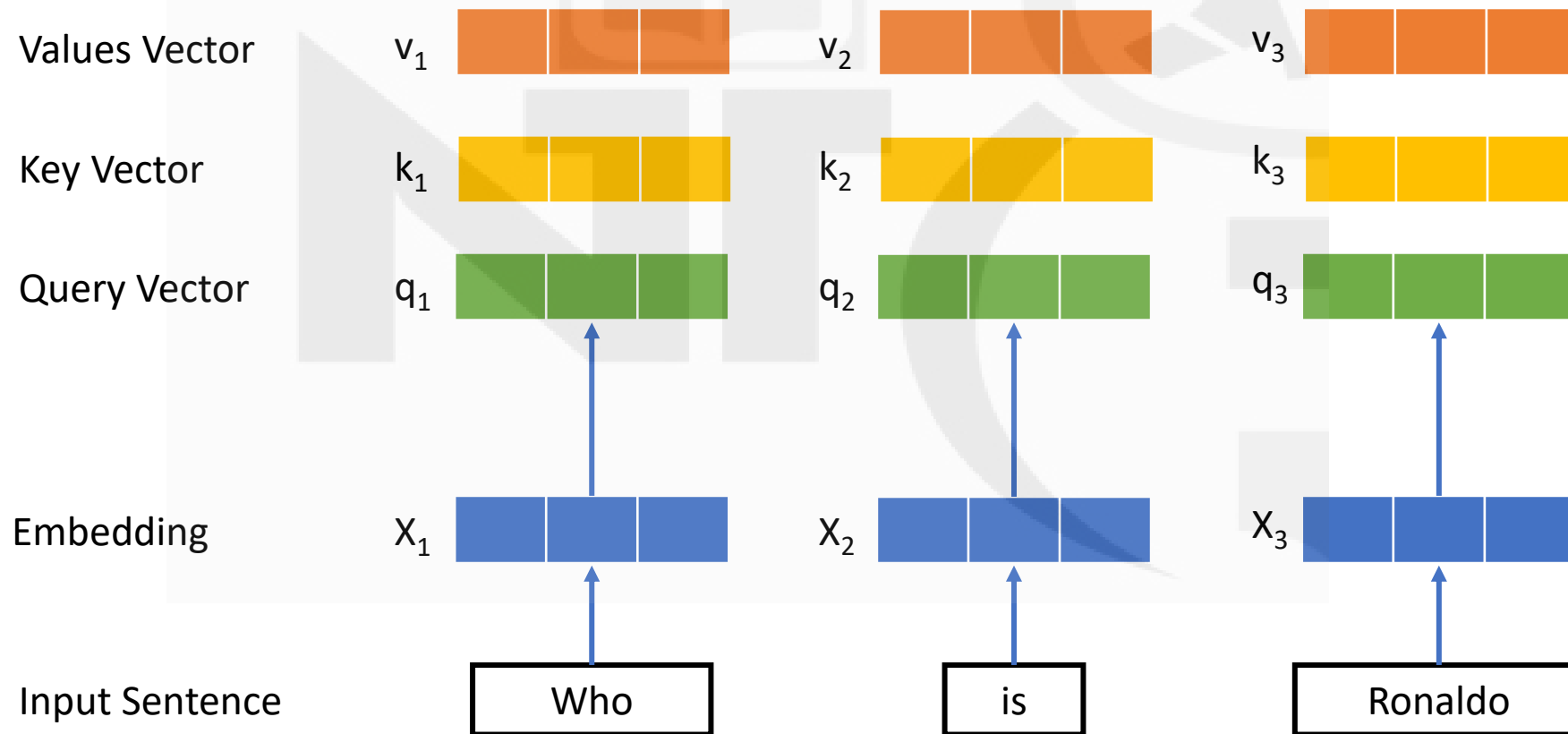
Input Sentence

Who

is

Ronaldo

# Implementing Self Attention

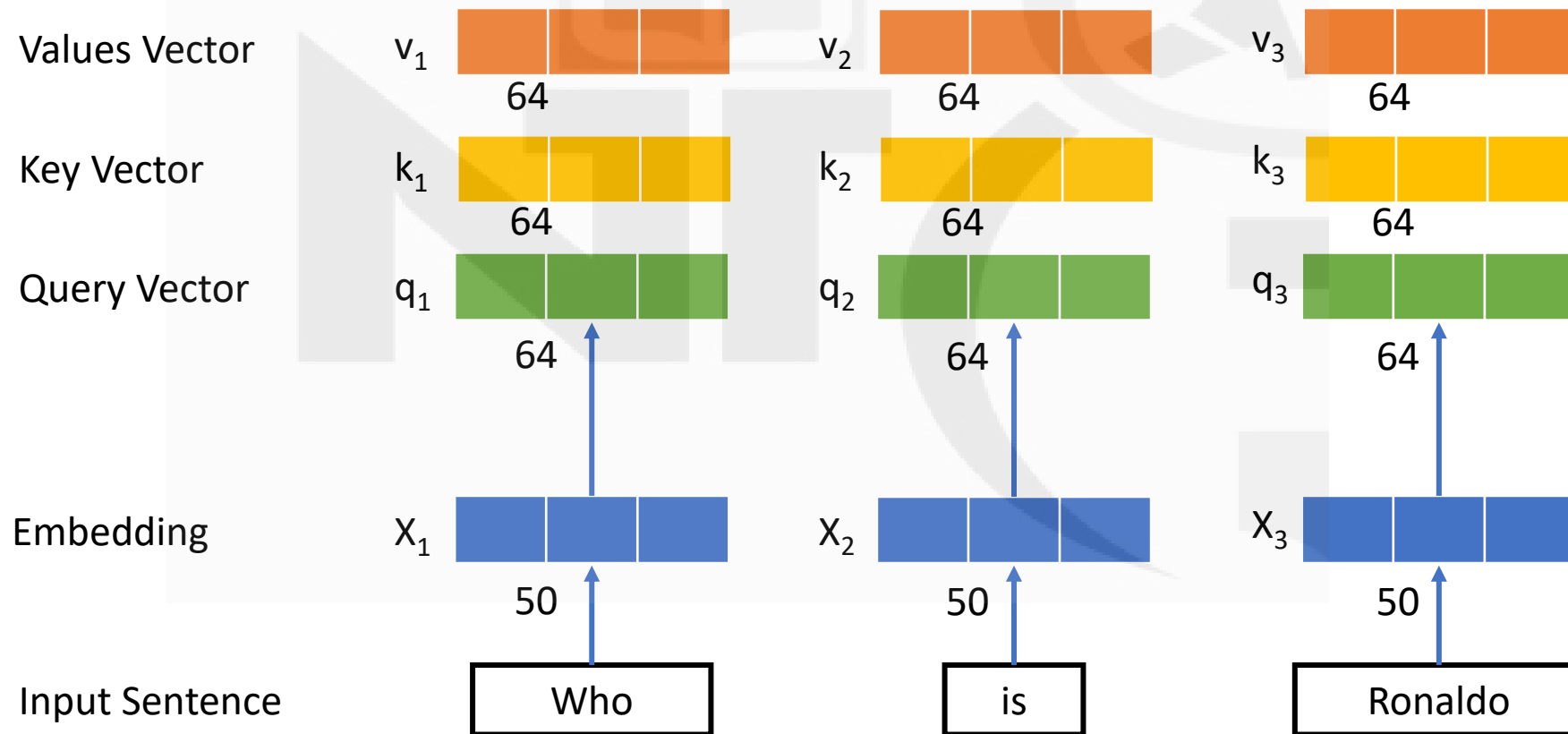


For each word, 3 vectors are created using the embedding as input. These vectors are created using 3 different Dense layers

Embedding is created for each word

# Implementing Self Attention

Say embedding =50, and  
the size of vectors =64

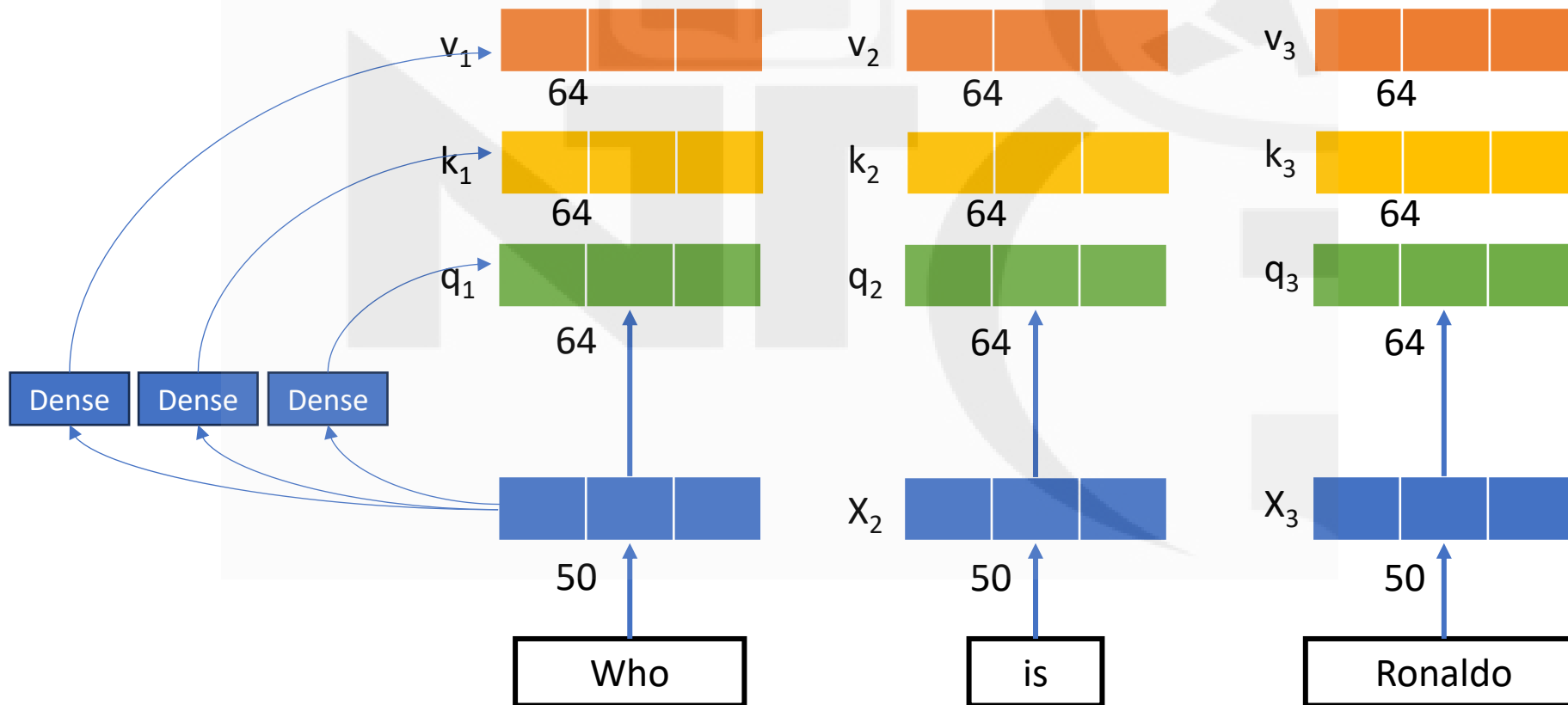


For each word,  
3 vectors are created  
using the embedding  
as input. These  
vectors are created  
using 3 different  
Dense layers

Embedding is created  
for each word

# Implementing Self Attention

How do we convert 50 num embedding to 64 number vectors



For each word, 3 vectors are created using the embedding as input. These vectors are created using 3 different Dense layers

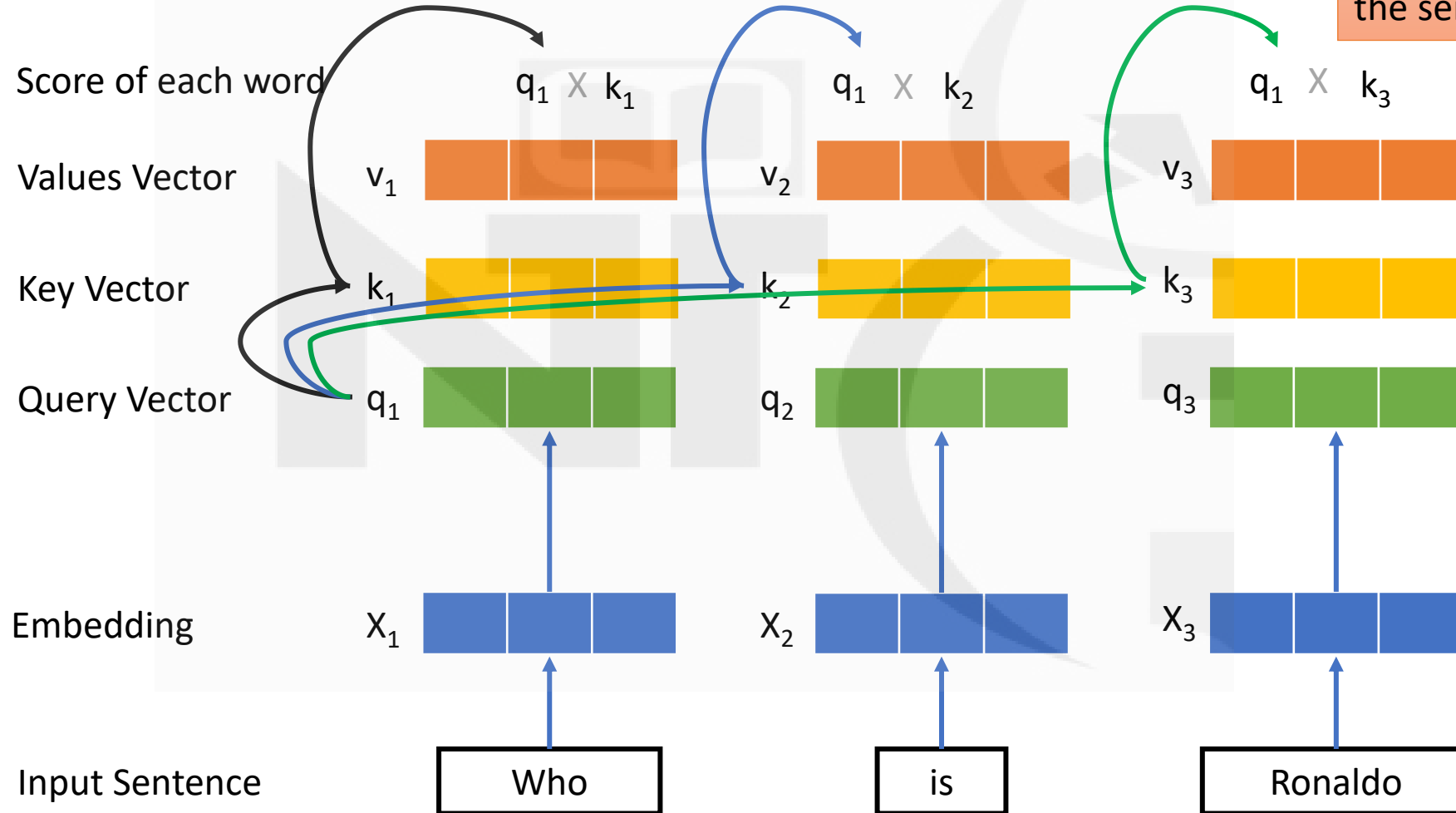
Embedding is created for each word



## Implementing Self Attention

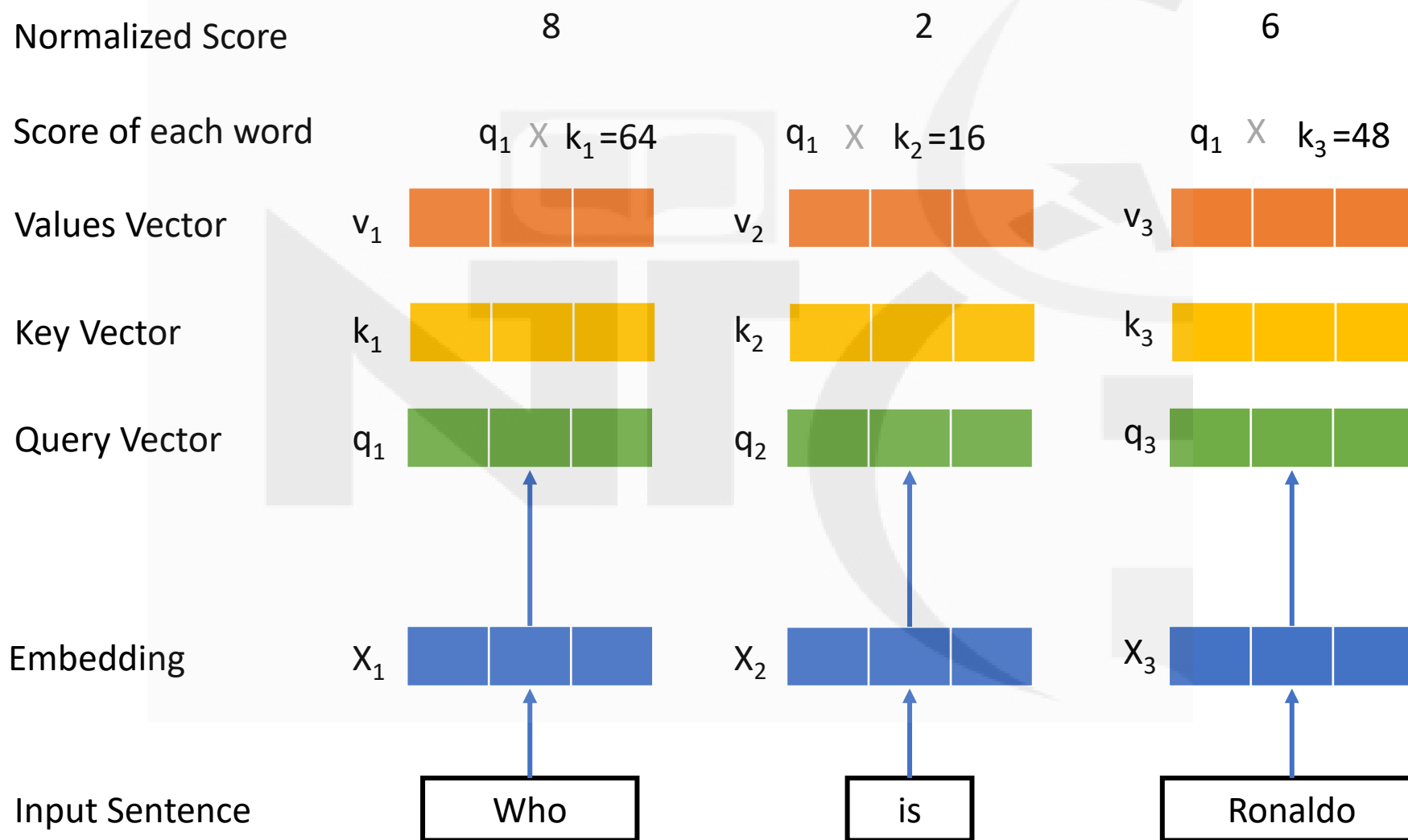
We multiply Query and key vector to check if there is a match

Dot product of query vector of 1st word and key vector of each word in the sentence.



**Calculating Self  
Attention of First  
word**

# Implementing Self Attention



Scores are divided by square root of length of key vector e.g 8 (square root of 64)

**Calculating Self Attention of First word**

# Implementing Self Attention

Softmax

0.01

0.0

0.99

Normalized Score

8

2

6

Score of each word

$q_1 \times k_1 = 64$

$q_1 \times k_2 = 16$

$q_1 \times k_3 = 48$

Values Vector

$v_1$

$v_2$

$v_3$

Key Vector

$k_1$

$k_2$

$k_3$

Query Vector

$q_1$

$q_2$

$q_3$

Embedding

$x_1$

$x_2$

$x_3$

Input Sentence

Who

is

Ronaldo

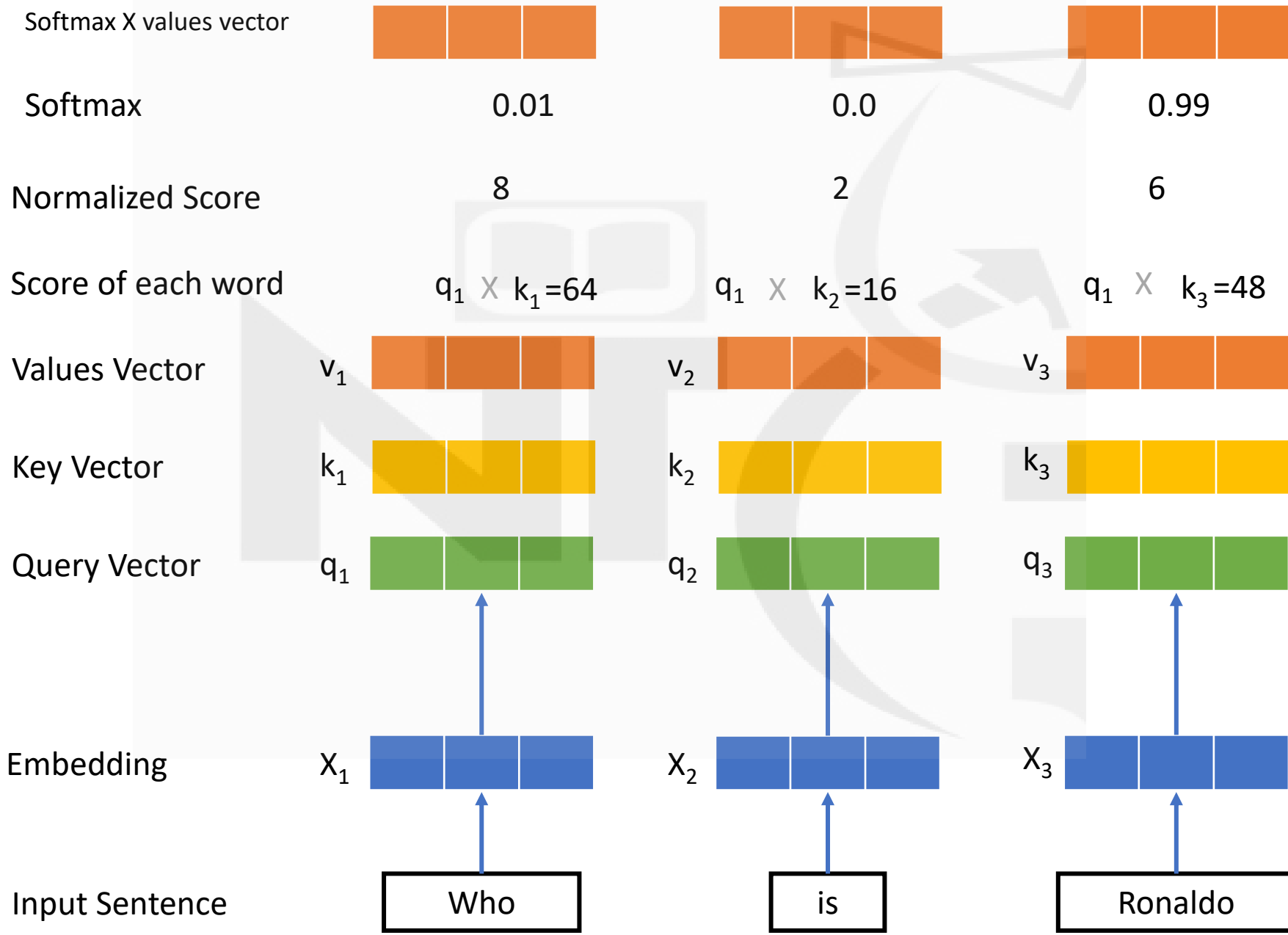
Values for each word are multiplied with respective softmax score

Softmax provides relative importance of each word for 1st word

Scores are divided by square root of length of key vector e.g 8 (square root of 64)

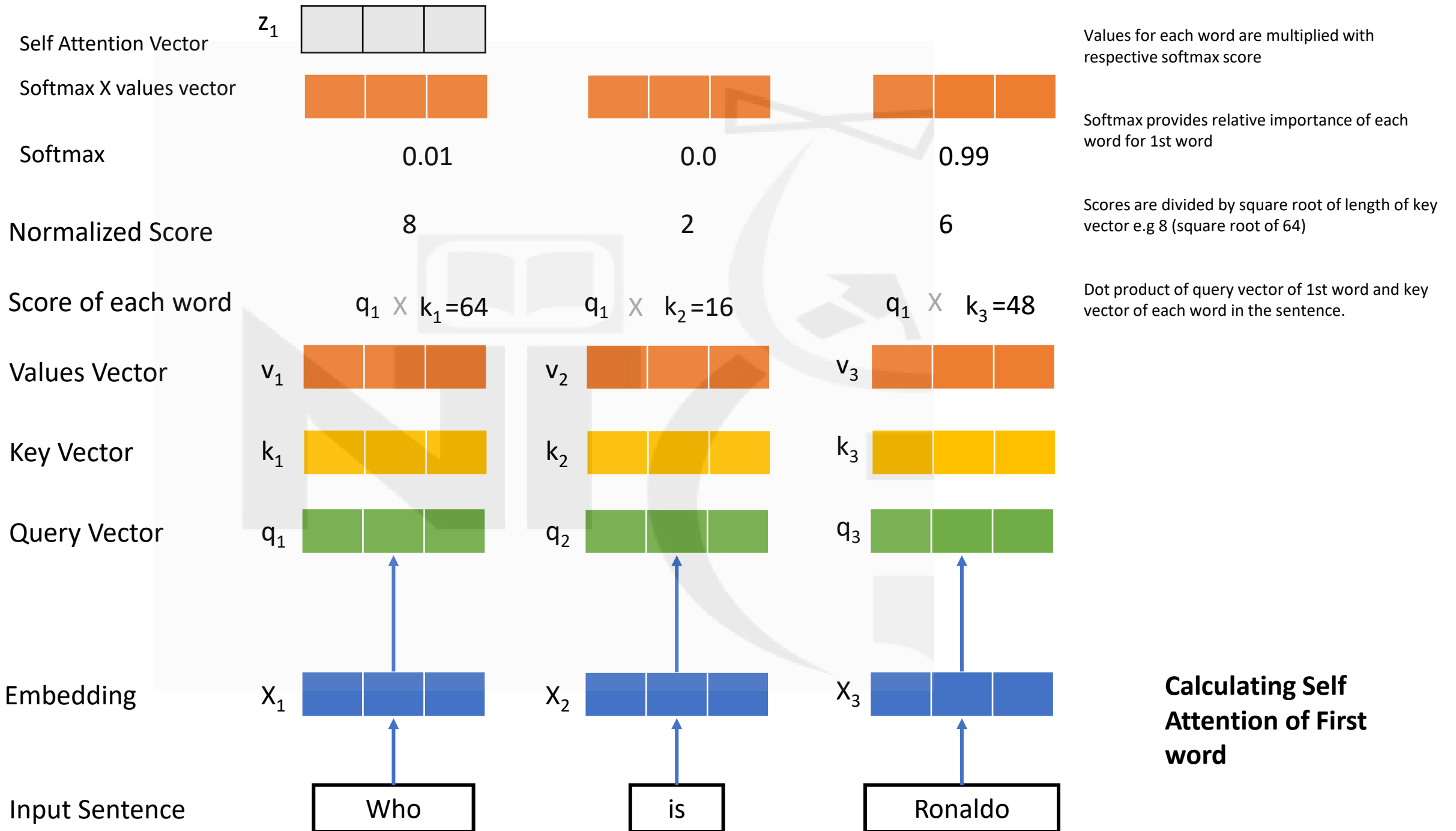
Dot product of query vector of 1st word and key vector of each word in the sentence.

**Calculating Self Attention of First word**



Scores are divided by square root of length of key vector e.g 8 (square root of 64)

**Calculating Self Attention of First word**



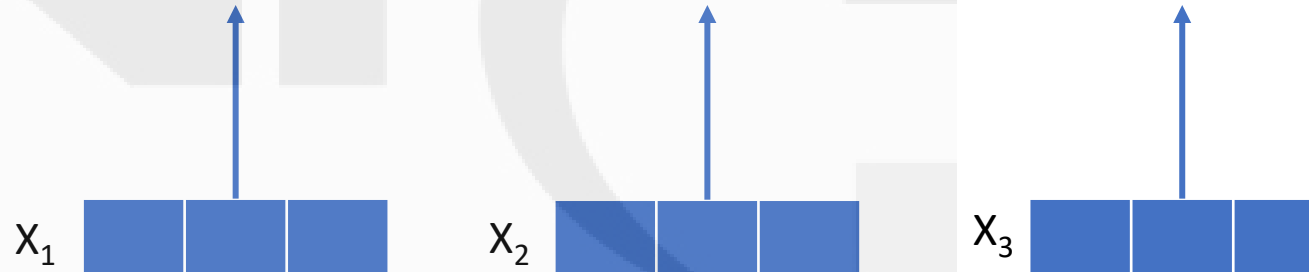
Repeat the process to calculate the self Attention vectors of all the words

Self Attention Vector



**Self Attention Layer**

Embedding

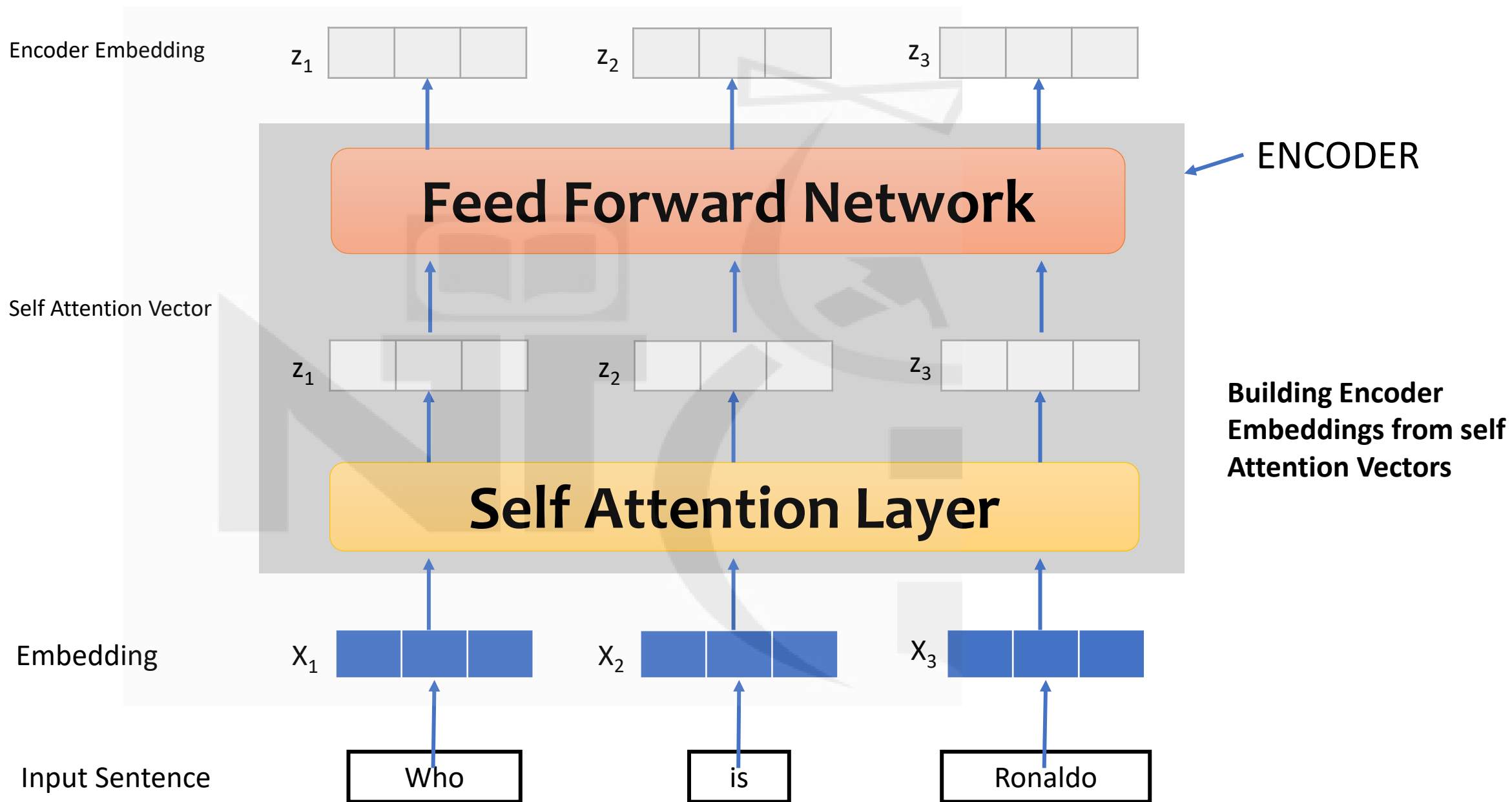


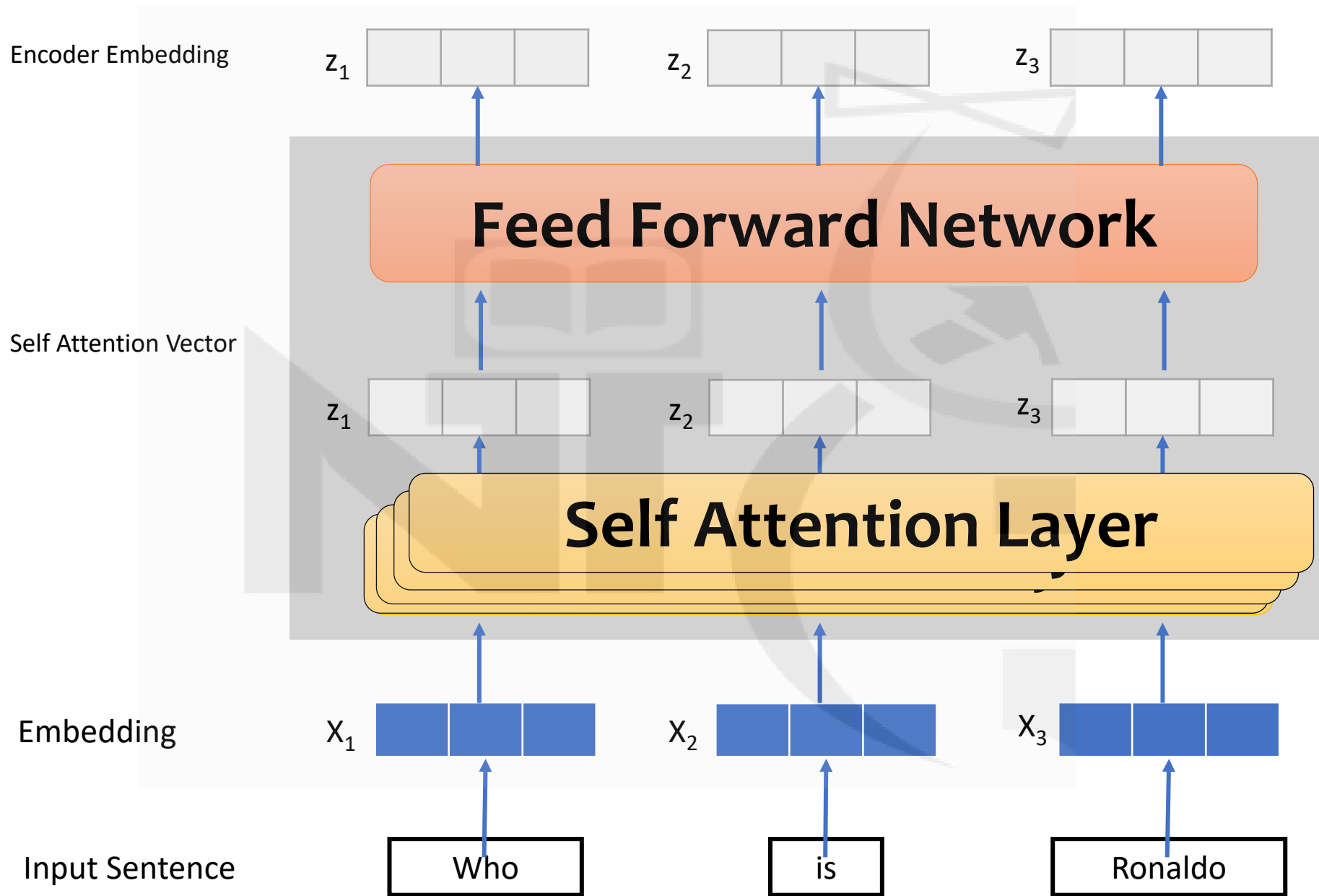
Input Sentence

Who

is

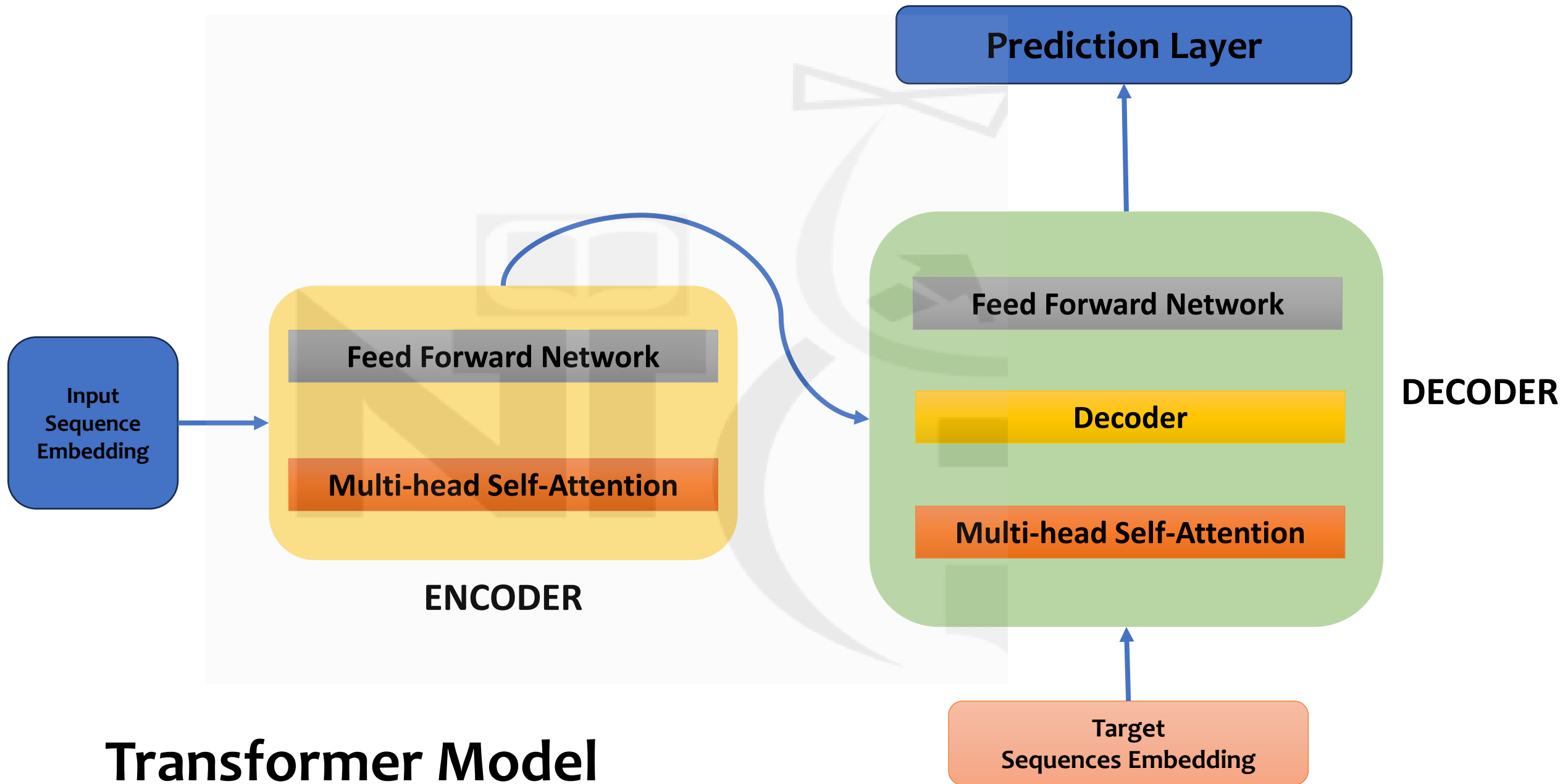
Ronaldo





**Encoder uses multiple Self-Attention layers each with its own Key, Query and Values vector**

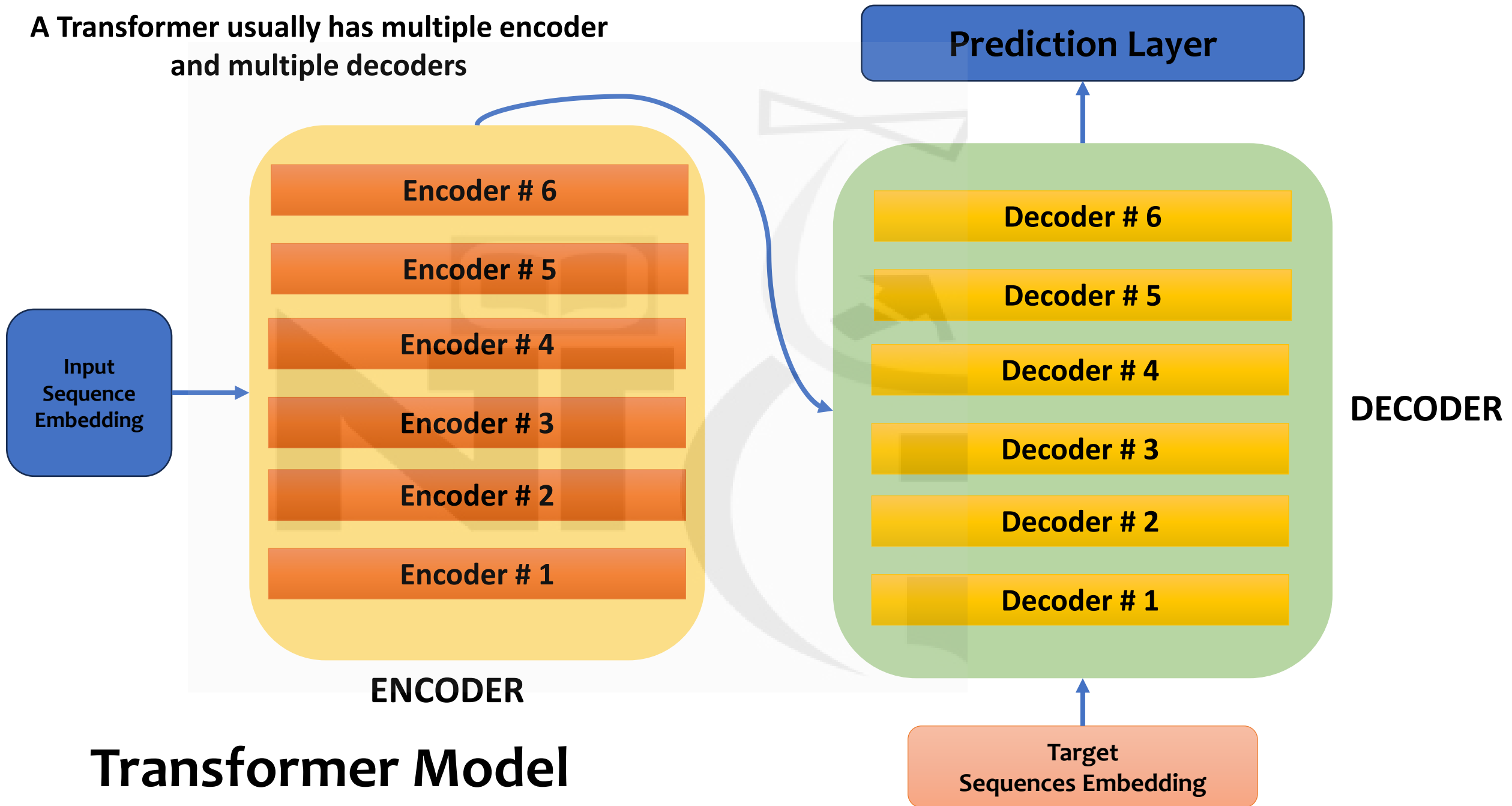


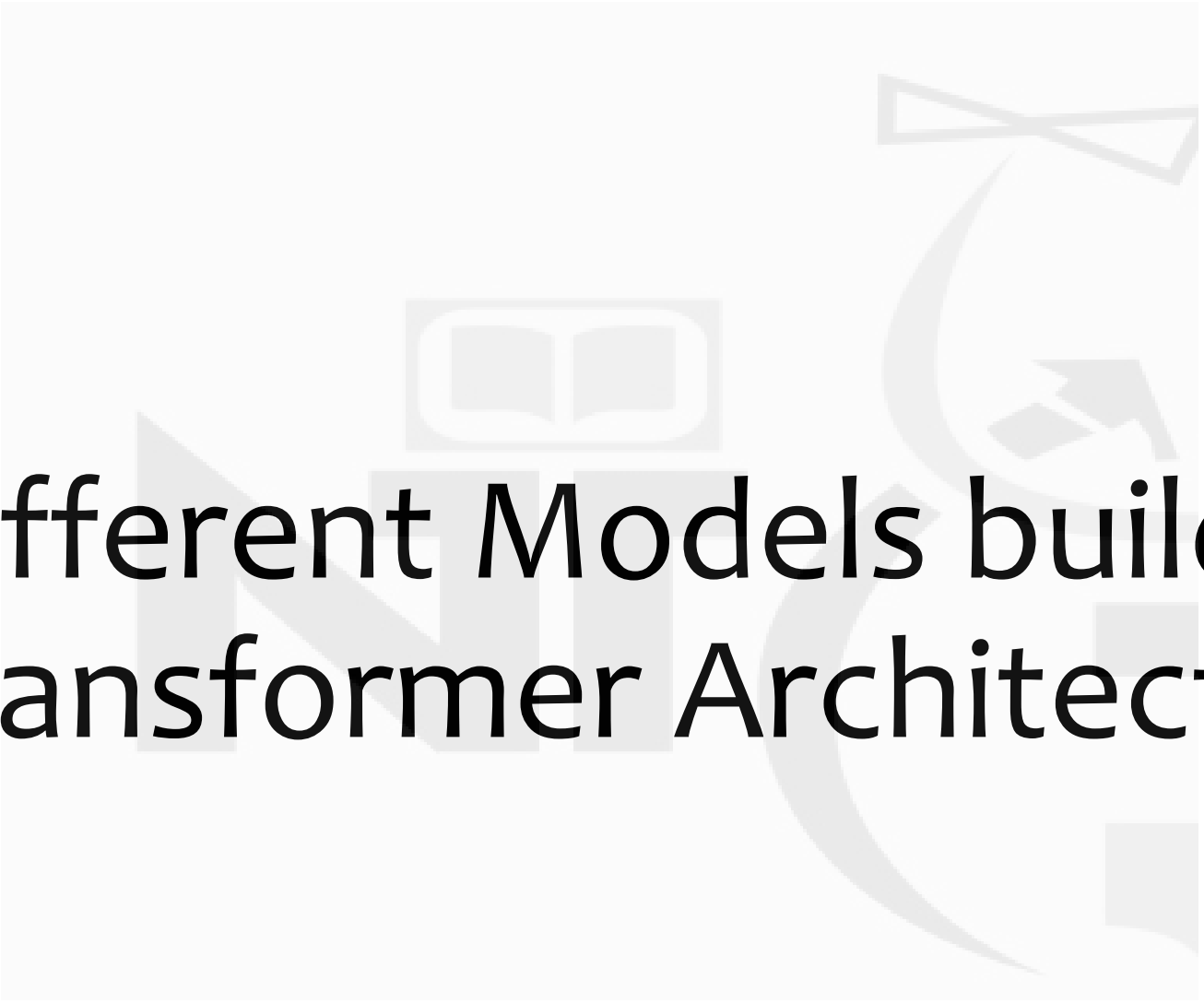




The authors of Transformers dint stop here

A Transformer usually has multiple encoder and multiple decoders





# Different Models build using Transformer Architecture

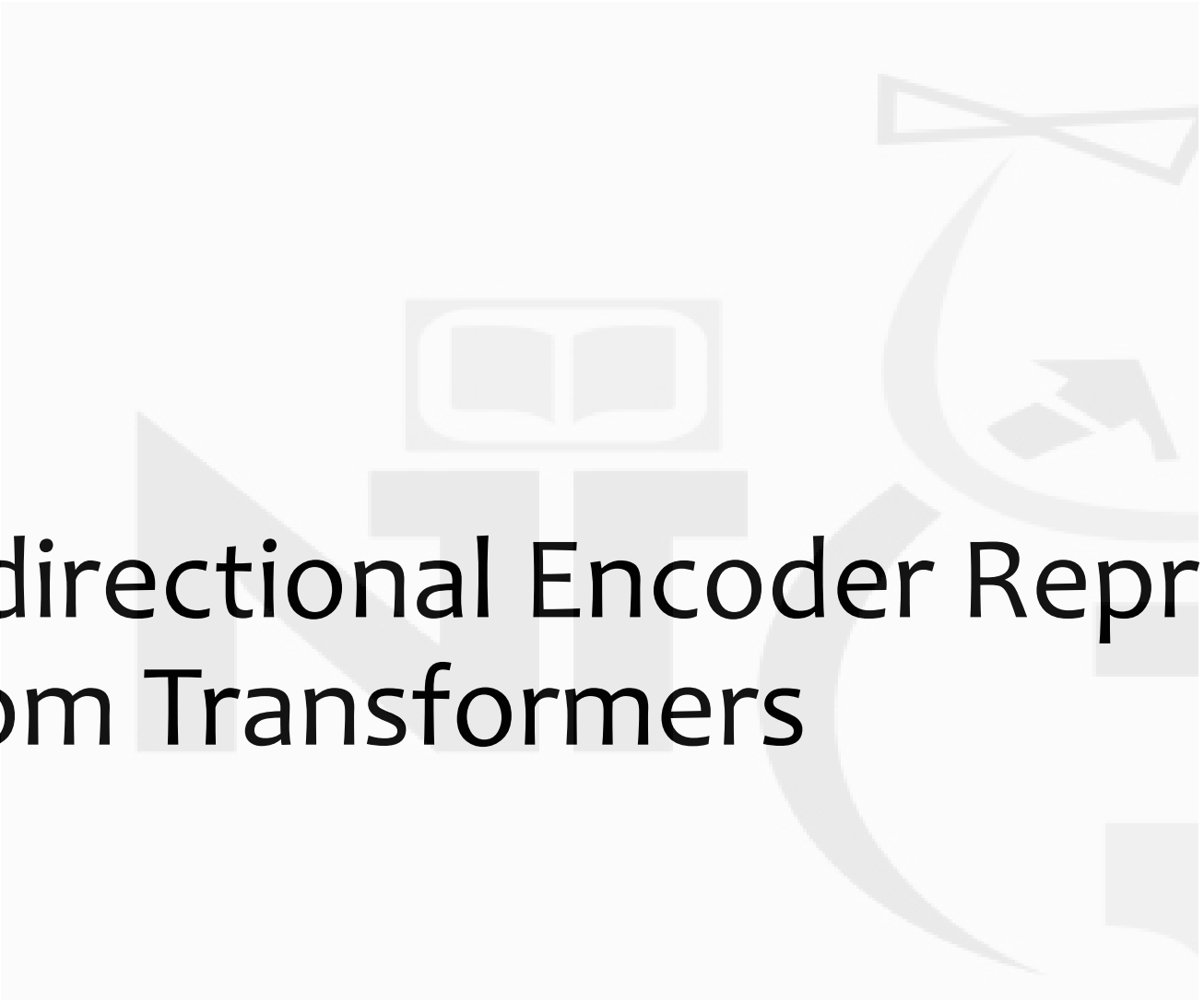
Transformer  
XL

BERT

GPT

Rober Ta

Albert



# Bidirectional Encoder Representations from Transformers

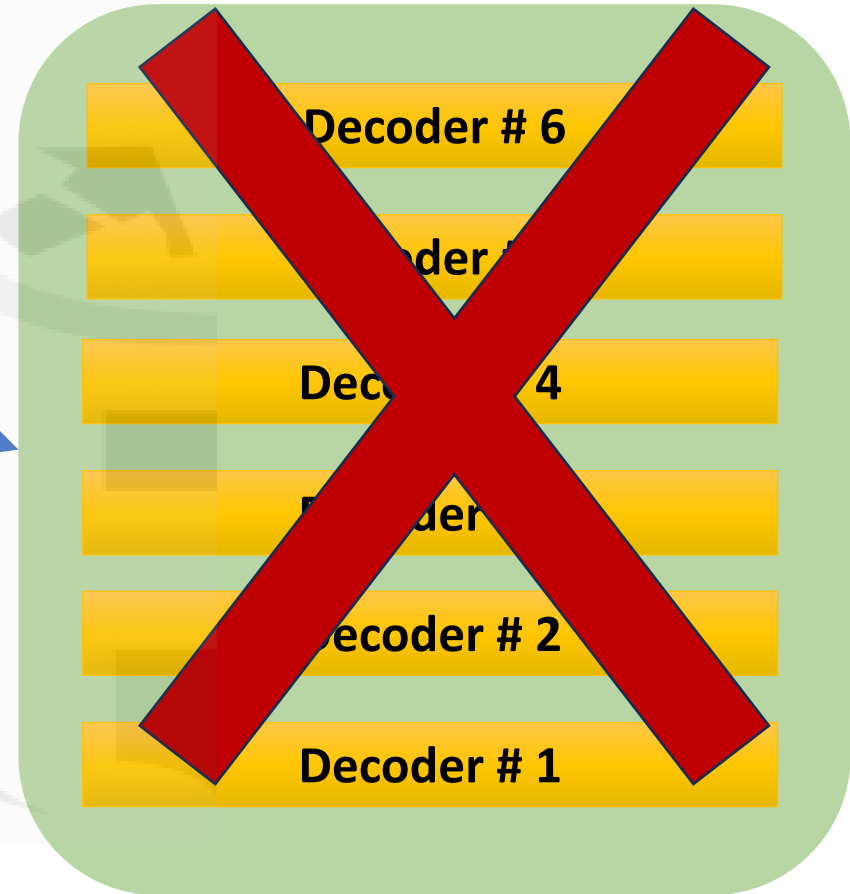
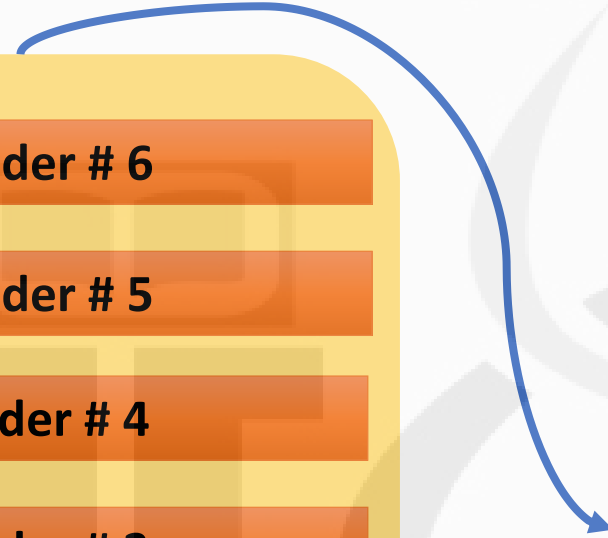
# BERT

Only involves Encoders

ENCODER



Input Sequence Embedding



DECODER

# How do we train BERT?

- With tons of data but in self-supervised way



# How BERT is trained in self supervised way

- Masked language Modelling (MLM)
- Next Sentence Prediction

# **BERT**

**Masked Language  
Modelling**

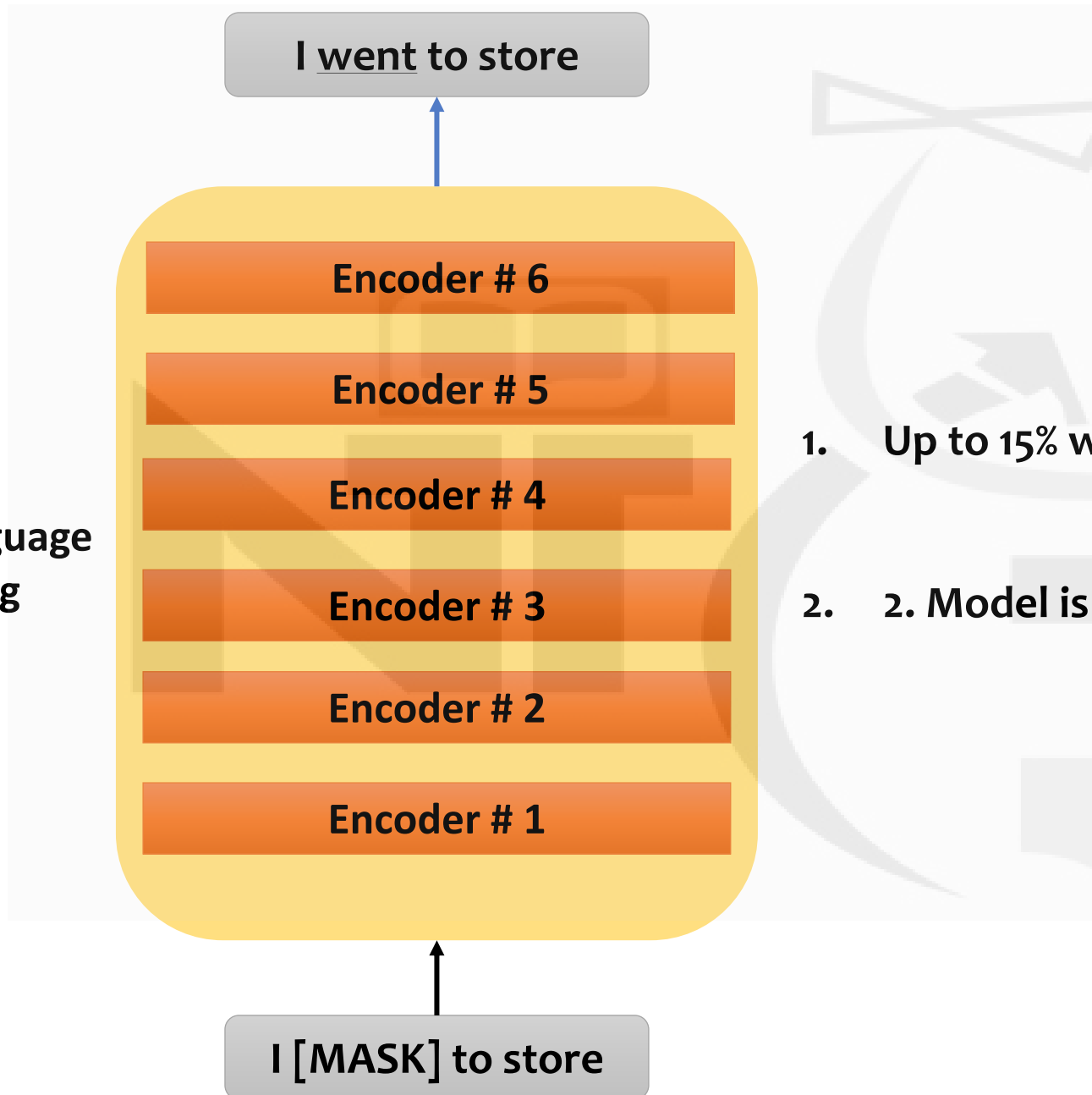


Input Sequence Embedding

1. Up to 15% words in the input sequence are masked.
2. Model is expected to predict masked word(s).

# BERT

Masked Language  
Modelling



1. Up to 15% words in the input sequence are masked.
2. Model is expected to predict masked word(s).

# BERT

Next Sentence  
Prediction(NSP)



1. Feed two sentences to the model
2. Model will predict if 2nd sentence should follow 1st sentence.