

Procesamiento de Lenguaje Natural

Tópicos Avanzados en Analítica
Maestría en Analítica para la Inteligencia de Negocios

Sergio Alberto Mora Pardo - H2 2023

Encoder-Decoder Attention

Encoder-Decoder

SeqToSeq

Sutskever et al., 2014

Sequence to Sequence Learning with Neural Networks

Ilya Sutskever Oriol Vinyals Quoc V. Le
Google Google Google
ilyasu@google.com vinyals@google.com qvl@google.com

Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT-14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous state of the art. The LSTM also learned sensible phrase and sentence representations that are sensitive to word order and are relatively invariant to the active and the passive voice. Finally, we found that reversing the order of the words in all source sentences (but not target sentences) improved the LSTM's performance markedly, because doing so introduced many short term dependencies between the source and the target sentence which made the optimization problem easier.

1 Introduction

Deep Neural Networks (DNNs) are extremely powerful machine learning models that achieve excellent performance on difficult problems such as speech recognition [13, 7] and visual object recognition [19, 6, 21, 20]. DNNs are powerful because they can perform arbitrary parallel computation for a modest number of steps. A surprising example of the power of DNNs is their ability to sort N N -bit numbers using only 2 hidden layers of quadratic size [27]. So, while neural networks are related to conventional statistical models, they learn an intricate computation. Furthermore, large DNNs can be trained with supervised backpropagation whenever the labeled training set has enough information to specify the network's parameters. Thus, if there exists a parameter setting of a large DNN that achieves good results (for example, because humans can solve the task very rapidly), supervised backpropagation will find these parameters and solve the problem.

Despite their flexibility and power, DNNs can only be applied to problems whose inputs and targets can be sensibly encoded with vectors of fixed dimensionality. It is a significant limitation, since many important problems are best expressed with sequences whose lengths are not known a-priori. For example, speech recognition and machine translation are sequential problems. Likewise, question answering can also be seen as mapping a sequence of words representing the question to a

1

Cho et al., 2014

Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation

Kyunghyun Cho
Bart van Merriënboer Caglar Gulcehre Dzmitry Bahdanau
Université de Montréal Université de Montréal, CIFAR Senior Fellow
firstname.lastname@umontreal.ca d.bahdanau@jacobs-university.de

Fethi Bougares Holger Schwenk Yoshua Bengio
Université du Maine, France Université de Montréal, CIFAR Senior Fellow
firstname.lastname@lum.univ-lemans.fr find.me@on-the.web

Abstract

Along this line of research on using neural networks for SMT, this paper focuses on a novel neural network architecture that can be used as a part of the conventional phrase-based SMT system. The proposed neural network architecture, which we will refer to as an *RNN Encoder–Decoder*, consists of two recurrent neural networks (RNN) that act as an encoder and a decoder pair. The encoder maps a variable-length source sequence to a fixed-length vector, and the decoder maps the vector representation back to a variable-length target sequence. The two networks are trained jointly to maximize the conditional probability of a target sequence given a source sequence. The performance of a statistical machine translation system is empirically found to improve by using the conditional probabilities of phrase pairs computed by the RNN Encoder–Decoder as an additional feature in the existing log-linear model. Qualitatively, we show that the proposed model learns a semantically and syntactically meaningful representation of linguistic phrases.

1 Introduction

Deep neural networks have shown great success in various applications such as objection recognition (see, e.g., (Krizhevsky et al., 2012)) and speech recognition (see, e.g., (Dahl et al., 2012)).

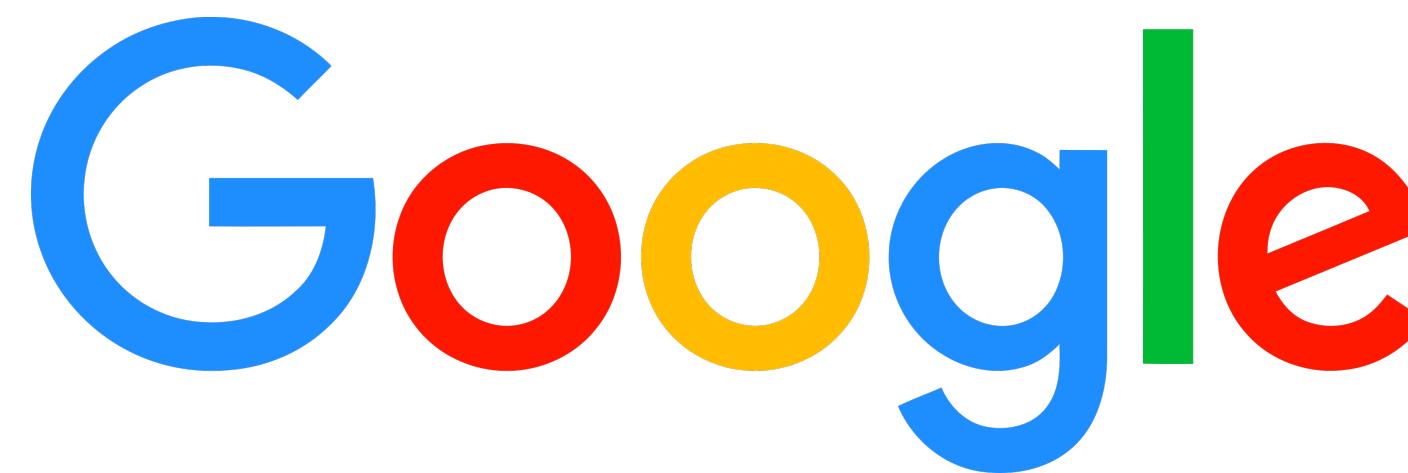
Furthermore, many recent works showed that neural networks can be successfully used in a number of tasks in natural language processing (NLP). These include, but are not limited to, language modeling (Bengio et al., 2003), paraphrase detection (Socher et al., 2011) and word embedding extraction (Mikolov et al., 2013). In the field of statistical machine translation (SMT), deep neural networks have begun to show promising results. (Schwenk, 2012) summarizes a successful usage of feedforward neural networks in the framework of phrase-based SMT system.

We qualitatively analyze the trained RNN Encoder–Decoder by comparing its phrase scores with those given by the existing translation model.

The qualitative analysis shows that the RNN Encoder–Decoder is better at capturing the linguistic regularities in the phrase table, indirectly explaining the quantitative improvements in the overall translation performance. The further analysis of the model reveals that the RNN Encoder–Decoder learns a continuous space representation of a phrase that preserves both the semantic and syntactic structure of the phrase.

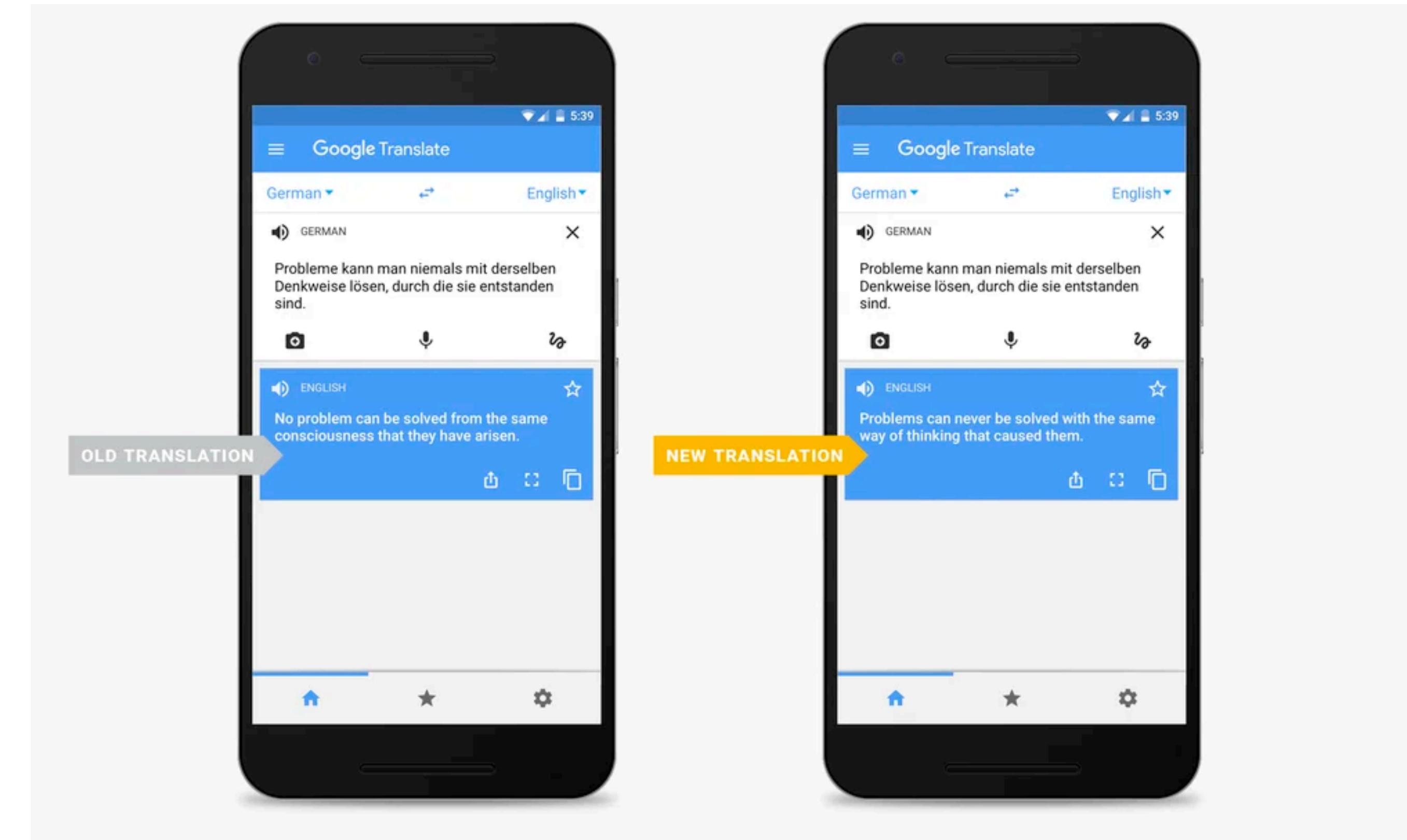
Encoder-Decoder Neural Machine Translation

*Found in translation: More accurate,
fluent sentences in Google Translate*



Cho et al., 2014

Sutskever et al., 2014



Encoder-Decoder

SeqToSeq

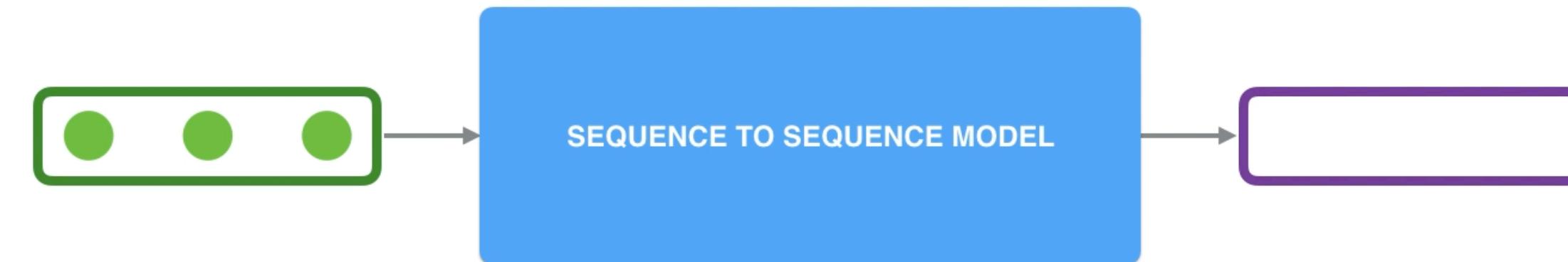
SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen



Encoder-Decoder

SeqToSeq

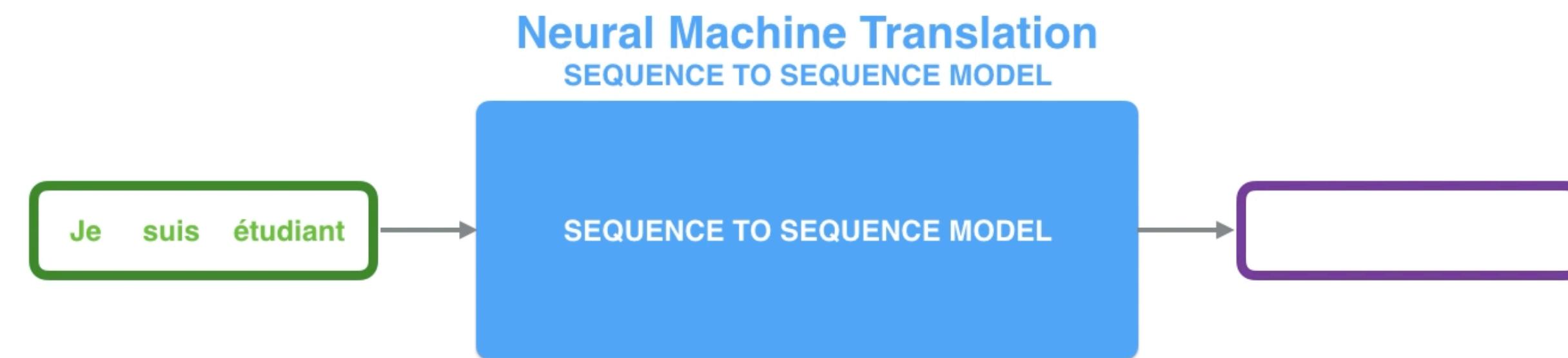
SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen



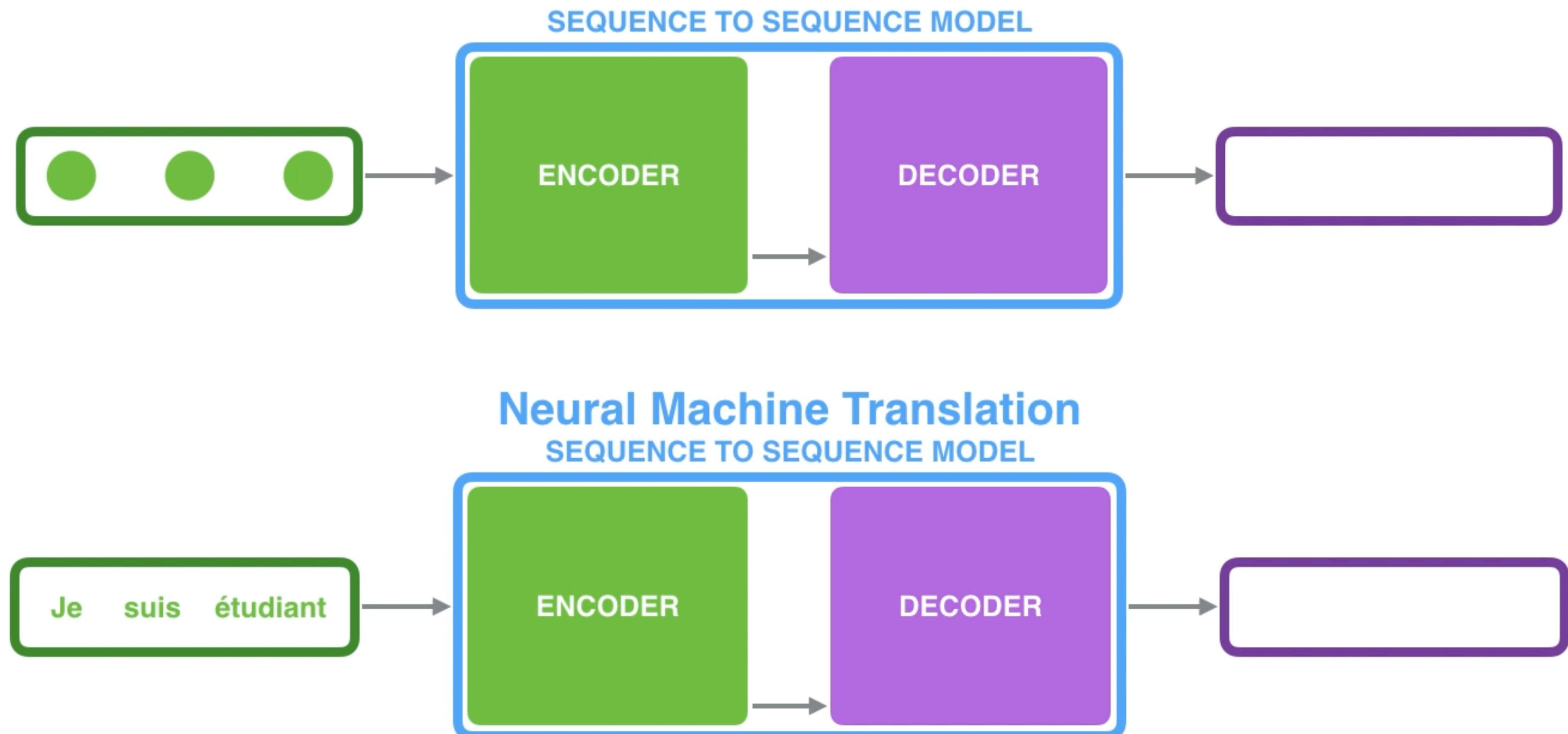
Encoder-Decoder

SeqToSeq

SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras
- letras
- Características de una imagen



Encoder-Decoder

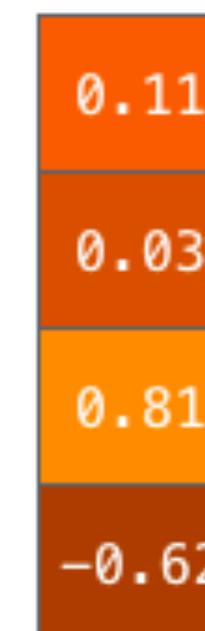
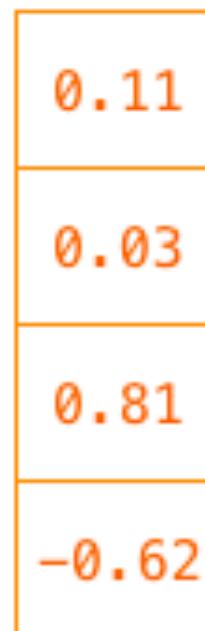
SeqToSeq

SeqtoSeq
(Neural Machine
Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras
- letras
- Características de una imagen

CONTEXT



Contexto es
un vector



Dimensión
Latente?

Tamaño regular:
256, 512 o 1024

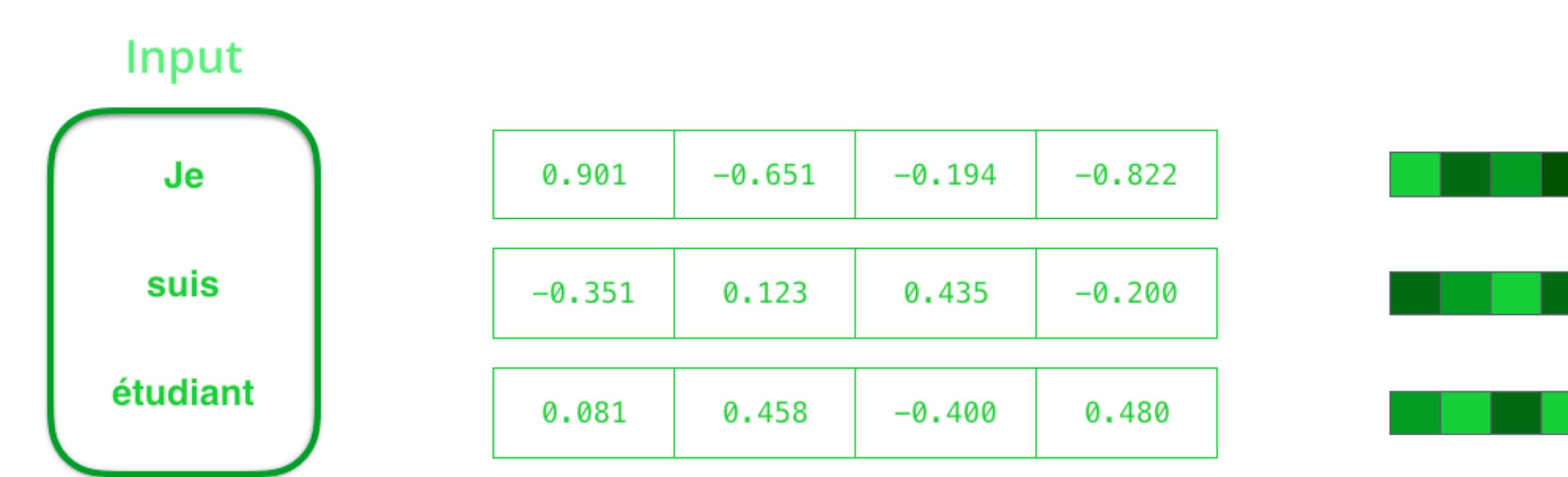
Encoder-Decoder

SeqToSeq

SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras
- letras
- Características de una imagen



Embedding



Tamaño regular: 200
o 300 dimensiones

Encoder-Decoder

SeqToSeq

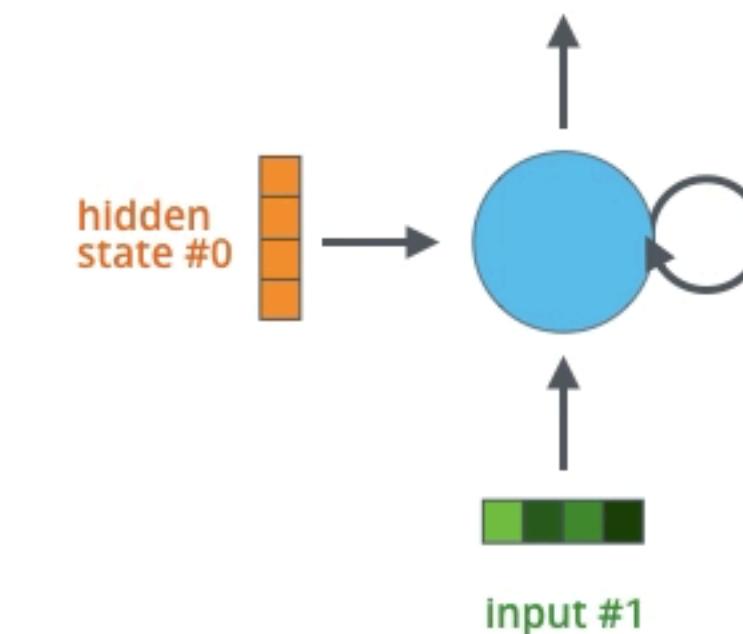
SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras
- letras
- Características de una imagen

Recurrent Neural Network

Time step #1:
An RNN takes two input vectors:



Encoder-Decoder

SeqToSeq

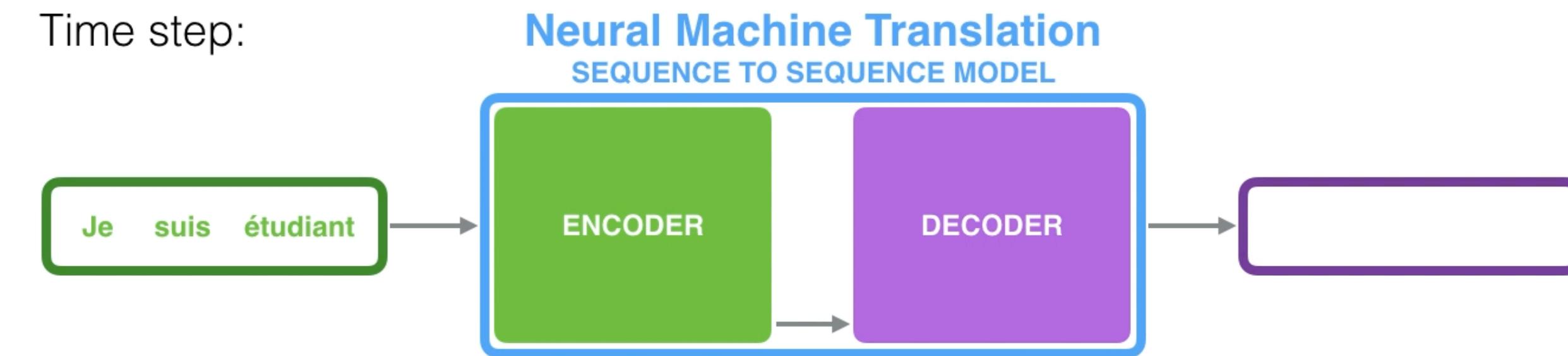
SeqtoSeq (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen



Encoder-Decoder

SeqToSeq

SeqtoSeq (Neural Machine Translation)

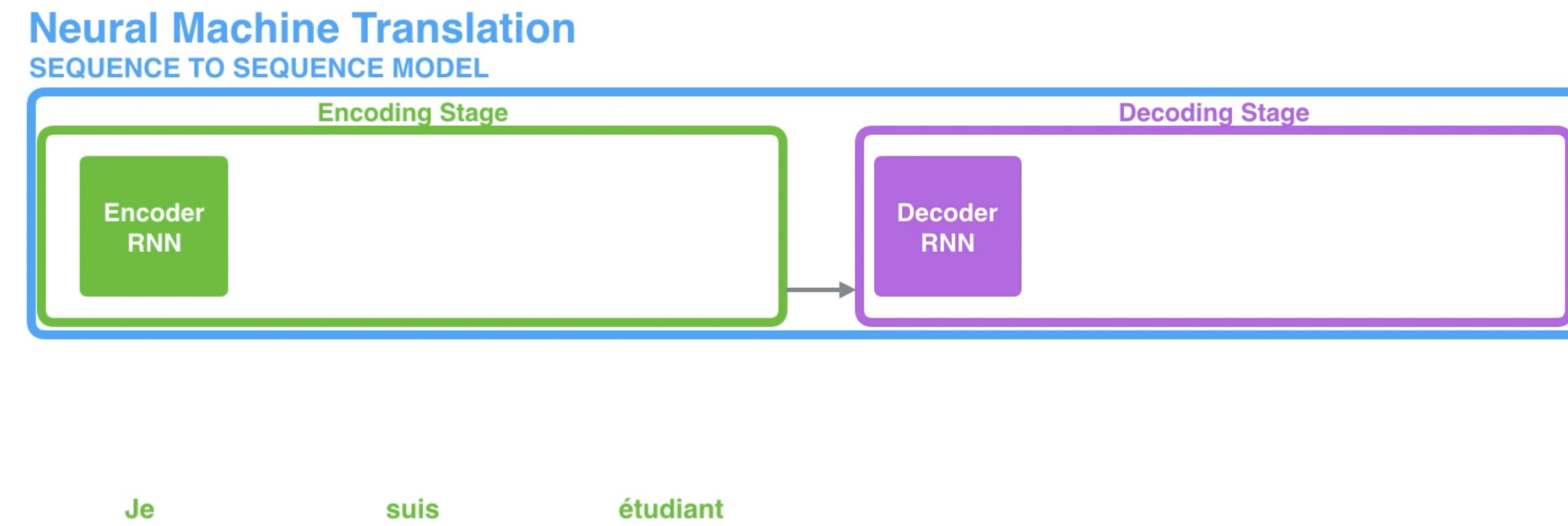
Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen

Desenrollando la RNN



Contexto se vuelve un “cuello de botella”

Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen



Bahdanau et al., 2014

Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang Luong Hieu Pham Christopher D. Manning
Computer Science Department, Stanford University, Stanford, CA 94305
{lmthang,hyhieu,manning}@stanford.edu

Abstract

An attentional mechanism has lately been used to improve neural machine translation (NMT) by selectively focusing on parts of the source sentence during translation. However, there has been little work exploring useful architectures for attention-based NMT. This paper examines two simple and effective classes of attentional mechanism: a *global* approach which always attends to all source words and a *local* one that only looks at a subset of source words at a time. We demonstrate the effectiveness of both approaches on the WMT15 translation tasks between English and German in both directions. With local attention, we achieve a significant gain of 5.0 BLEU points over non-attentional systems that already incorporate known techniques such as beam search. Our ensemble model using different attention architectures yields a new state-of-the-art result in the WMT'15 English to German translation task with 25.9 BLEU points, an improvement of 1.0 BLEU points over the existing best system backed by NMT and an *n*-gram reranker.¹

1 Introduction

Neural Machine Translation (NMT) achieved state-of-the-art performances in large-scale translation tasks such as from English to French (Luong et al., 2015) and English to German (Jean et al., 2015). NMT is appealing since it requires minimal domain knowledge and is conceptually simple. The model by Luong et al. (2015) reads through all the source words until the end-of-sentence symbol <eos> is reached. It then starts emitting one target word at a time, as illustrated in Figure 1. NMT is often a large neural network that is trained in an end-to-end fashion and has the ability to generalize well to very long word sequences. This means the model does not have to explicitly store gigantic phrase tables and language models as in the case of standard MT; hence, NMT has a small memory footprint. Lastly, implementing NMT decoders is easy unlike the highly intricate decoders in standard MT (Koehn et al., 2003).

In parallel, the concept of “attention” has gained popularity in training neural networks, allowing models to learn alignments between different modalities, e.g., between image objects and agent actions in the dynamic control problem (Mnih et al., 2014), between speech frames and text in the speech recognition task (?), or between visual features of a picture and its text description in the image caption generation task (Xu et al., 2015). In the context of NMT, Bahdanau et al. (2015) has successfully applied such attentional mechanism to jointly translate and align words. To the best of our knowledge, there has not been any other work exploring the use of attention-based architectures for NMT.

¹All our code and models are publicly available at <http://nlp.stanford.edu/projects/nmt/>.

In this work, we design, with simplicity and ef-



Luong et al., 2015
Bengio*

NEURAL MACHINE TRANSLATION
BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau
Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio*
Université de Montréal

ABSTRACT

Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this encoder-decoder architecture. We propose to alleviate this bottleneck by allowing the model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

1 INTRODUCTION

Neural machine translation is a newly emerging approach to machine translation, recently proposed by Kalchbrenner and Blunsom (2013), Sutskever et al. (2014) and Cho et al. (2014b). Unlike the traditional phrase-based translation system (see, e.g., Koehn et al., 2003) which consists of many small sub-components that are tuned separately, neural machine translation attempts to build and train a single large neural network that reads a sentence and outputs a correct translation.

Most of the proposed neural machine translation models belong to a family of *encoder-decoders* (Sutskever et al., 2014; Cho et al., 2014a), with an encoder and a decoder for each language, or involve a language-specific encoder applied to each sentence whose outputs are then compared (Hermann and Blunsom, 2014). An encoder neural network reads and encodes a source sentence into a fixed-length vector. A decoder then outputs a translation from the encoded vector. The whole encoder-decoder system, which consists of the encoder and the decoder for a language pair, is jointly trained to maximize the probability of a correct translation given a source sentence.

A potential issue with this encoder-decoder approach is that a neural network needs to be able to compress all the information of a source sentence into a fixed-length vector. This may make it difficult for the neural network to cope with long sentences, especially those that are longer than the sentences in the training corpus. Cho et al. (2014b) showed that indeed the performance of a basic encoder-decoder deteriorates rapidly as the length of an input sentence increases.

In order to address this issue, we propose an extension to the encoder-decoder model which learns to align and translate jointly. Each time the proposed model generates a word in a translation, it (soft-)searches for a set of positions in a source sentence where the most relevant information is concentrated. The model then predicts a target word based on the context vectors associated with these source positions and all the previous generated target words.

*CIFAR Senior Fellow

Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

- palabras

- letras

- Características
de una imagen

Time step: 7

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



'Atención' permite al
decodificador
centrarse en '**étudiant**'

Amplificar la señal de la parte relevante de
la secuencia de entrada hace que los
modelos con atención produzcan mejores
resultados que los modelos sin atención.

Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

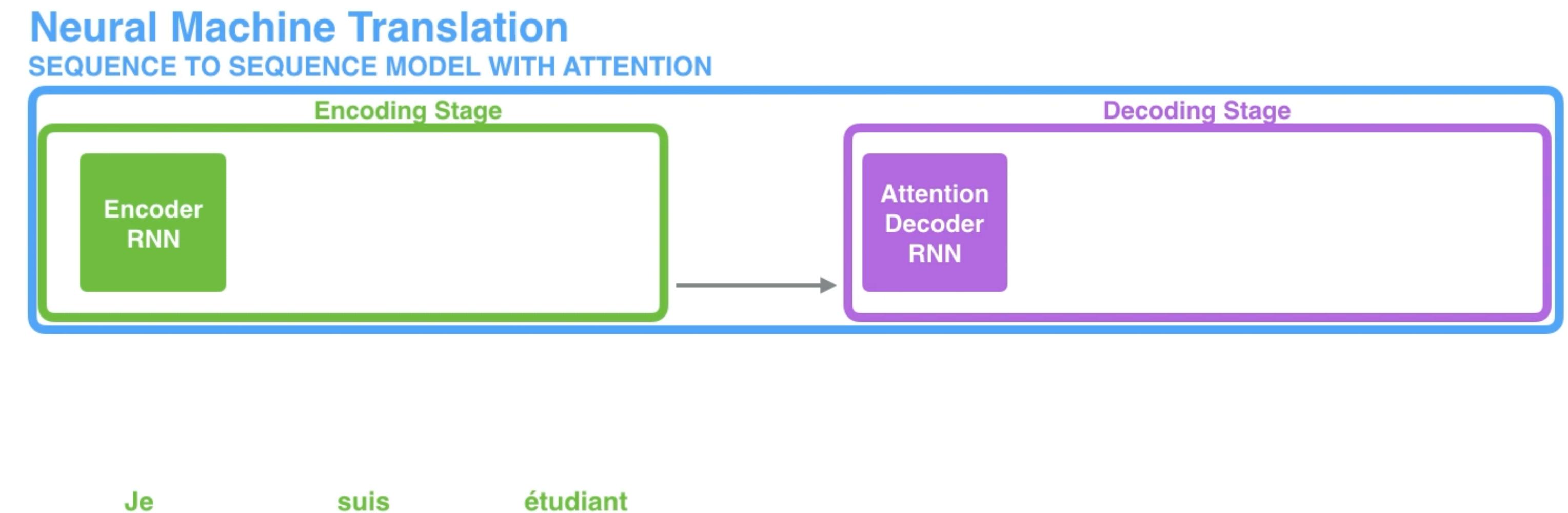
Toma una secuencia de elementos y genera otra secuencia de elementos

Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention



Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention

1. Pasamos todos los estados ocultos al decodificador

2. Centrarse en las partes de la entrada que son relevantes.

Revisión de estado oculto

Score para cada estado oculto

Multiplica ese score por una softmax

Cada estado oculto asociado a una palabra

Amplificando estados ocultos con puntuaciones altas

Ahogando estados ocultos con puntuaciones bajas

Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

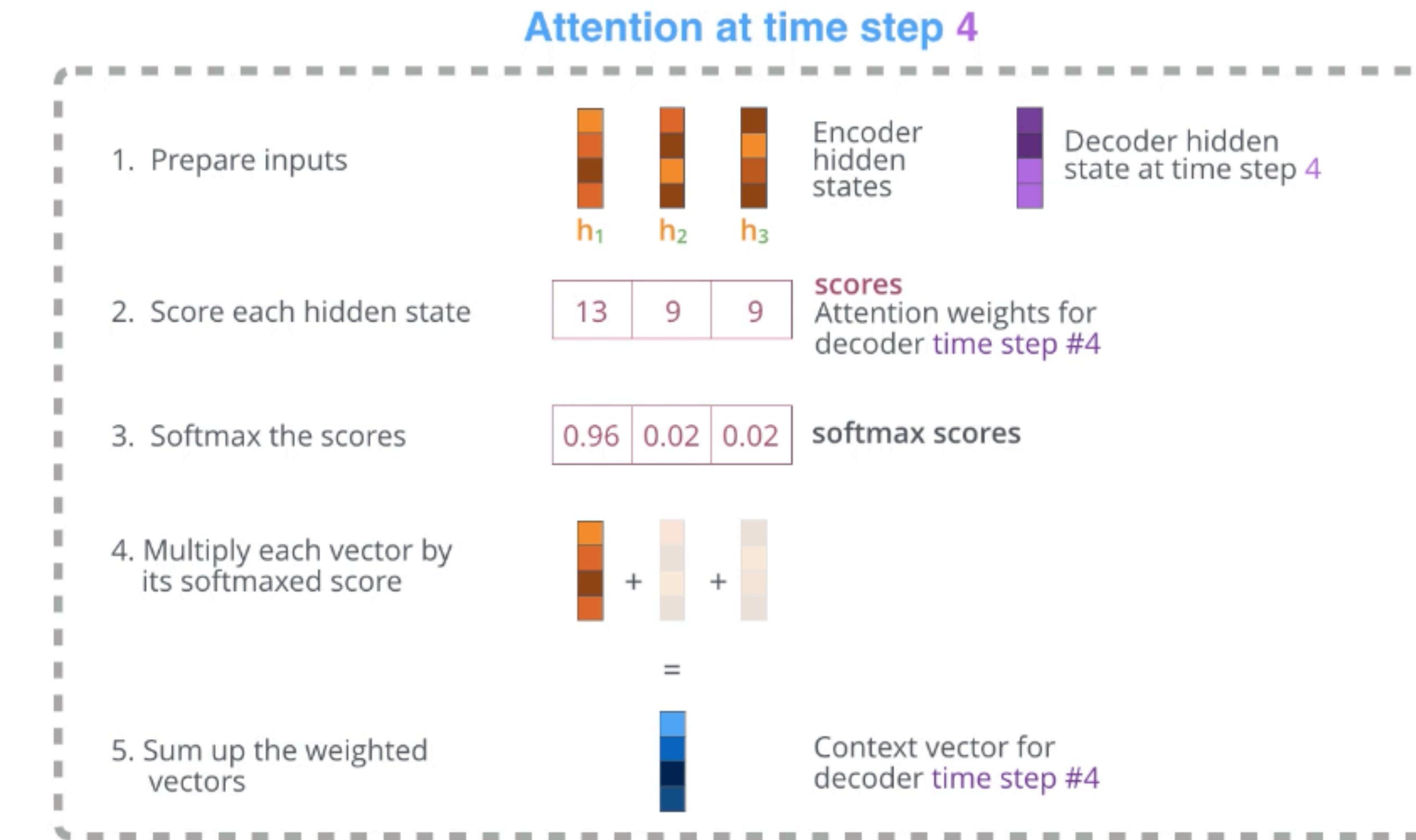
Toma una secuencia de elementos y genera otra secuencia de elementos

Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention



Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention

Proceso completo de atención

1. El decodificador de atención RNN incorpora el token <END> y un **estado oculto del decodificador inicial**.
2. El RNN procesa sus entradas, produciendo una salida y un **nuevo** vector de estado oculto (h_4). La salida se descarta.
3. Paso de atención: utilizamos los **estados ocultos del codificador** y el vector h_4 para calcular un vector de contexto (C_4) para este paso de tiempo.
4. Concatenamos h_4 y C_4 en un vector.
5. Pasamos este vector a través de una **red neuronal feedforward** (una entrenada conjuntamente con el modelo).
6. La **salida** de las redes neuronales de avance indica la palabra de salida de este paso de tiempo.
7. Repita los pasos la próxima vez.

Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

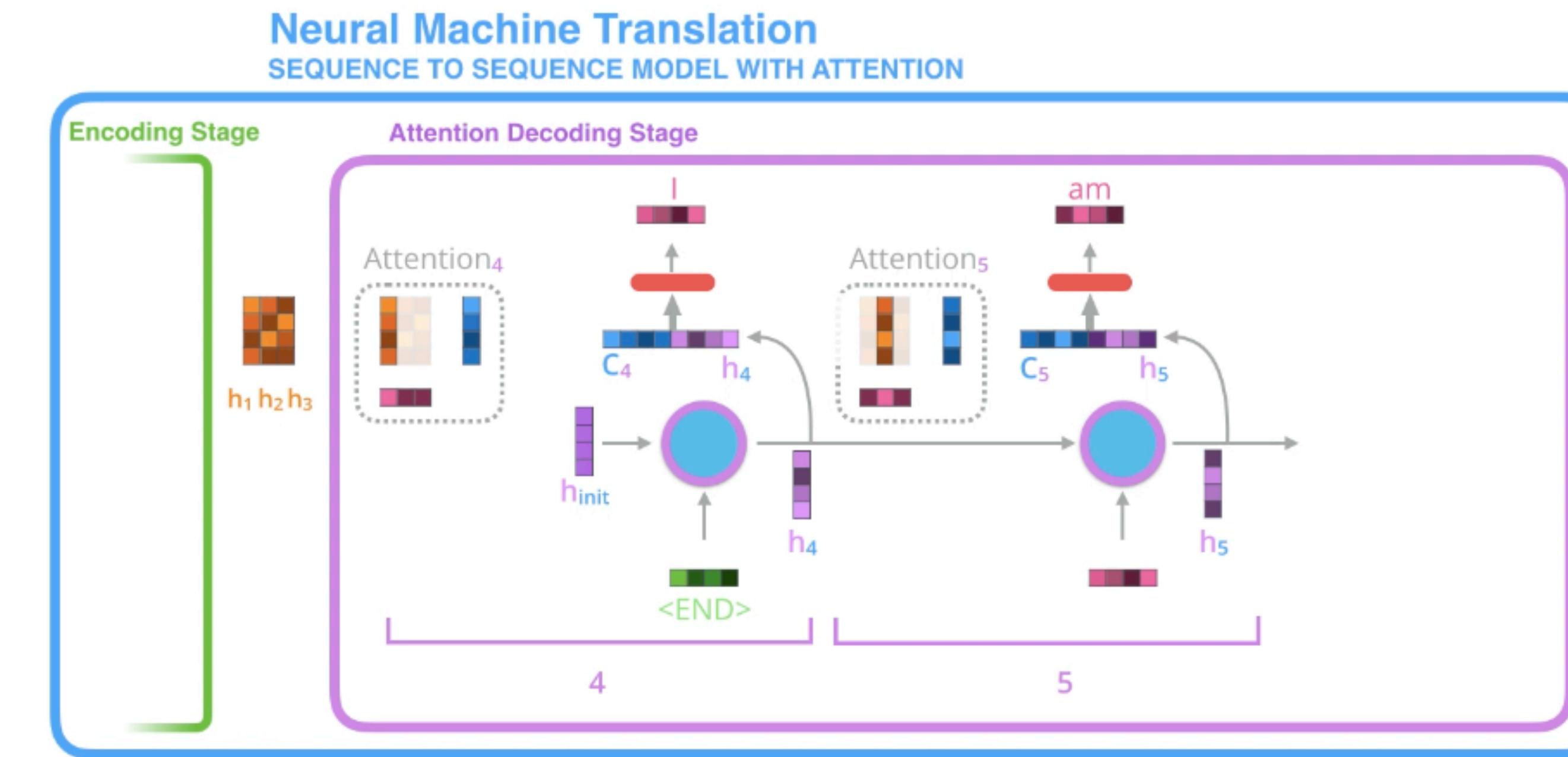
Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention

Proceso completo de atención



Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

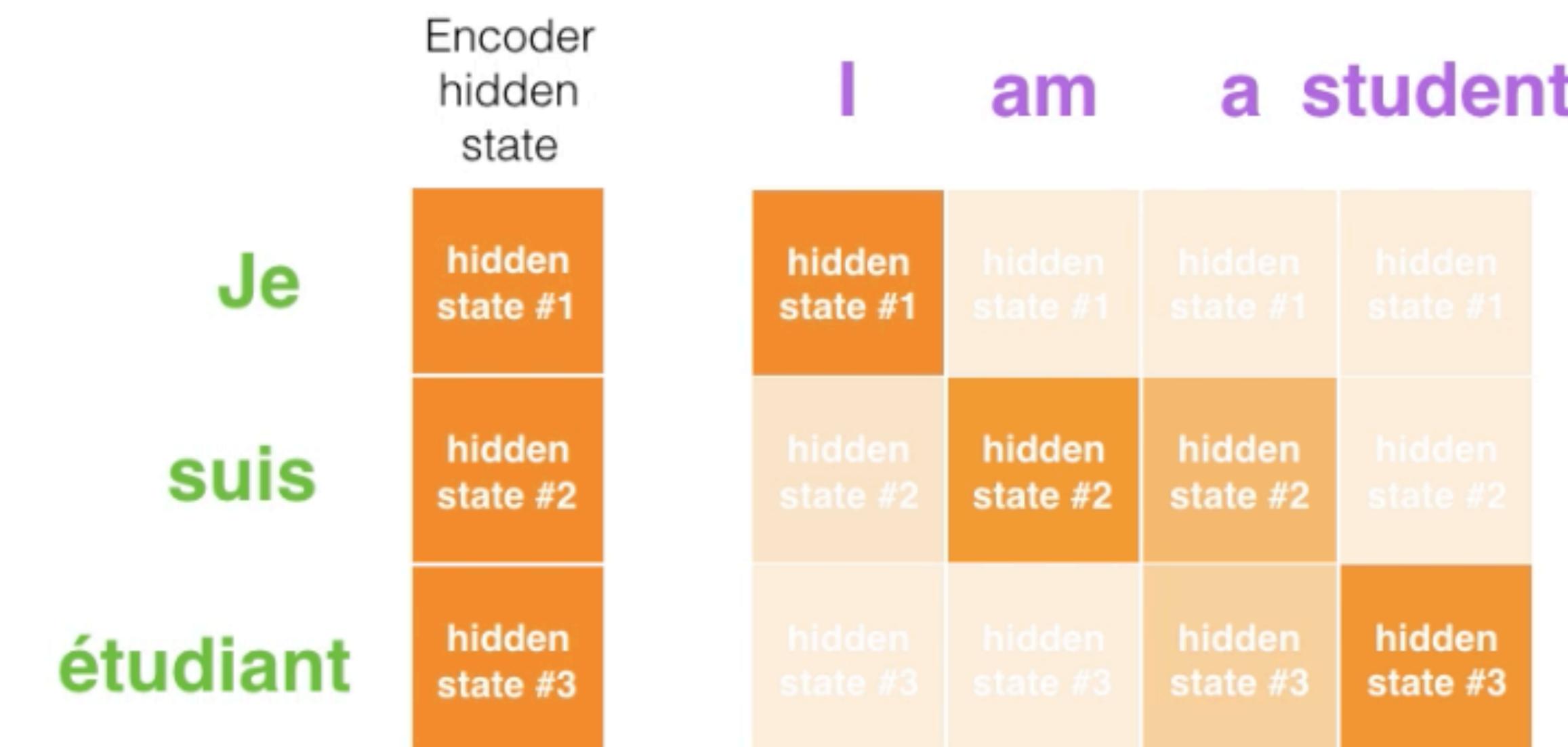
Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention

Proceso completo de atención



Encoder-Decoder

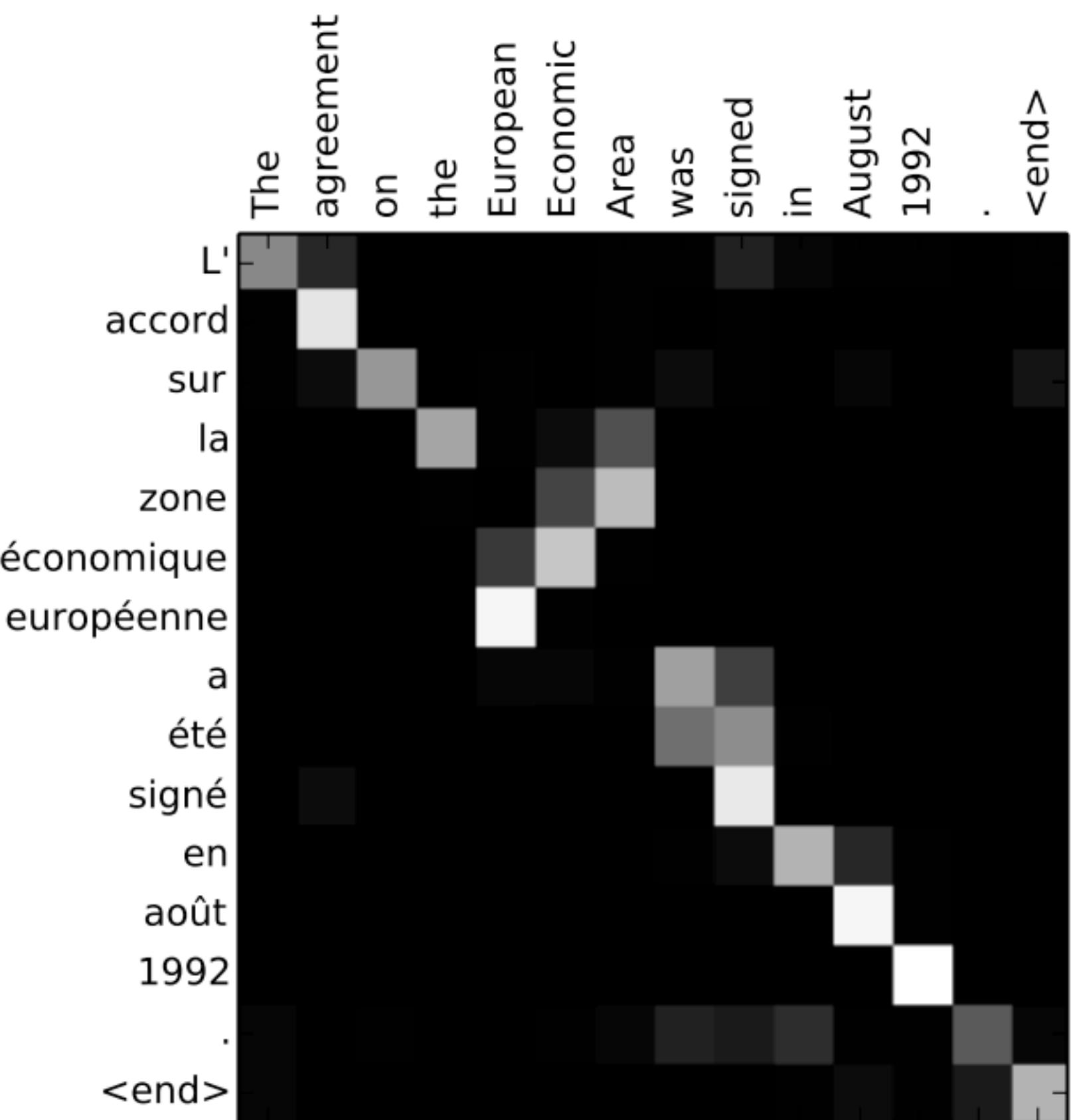
SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)
<i>Toma una secuencia de elementos y genera otra secuencia de elementos</i>
<i>Diferencias:</i>
SeqToSeq
Vs.
SeqToSeq + Attention

Proceso completo de atención

Traduce *indistinto del
orden de las palabras.*

'Espacio Económico
Europeo'



Encoder-Decoder

SeqToSeq with Attention

SeqtoSeq Attention (Neural Machine Translation)

Toma una secuencia de elementos y genera otra secuencia de elementos

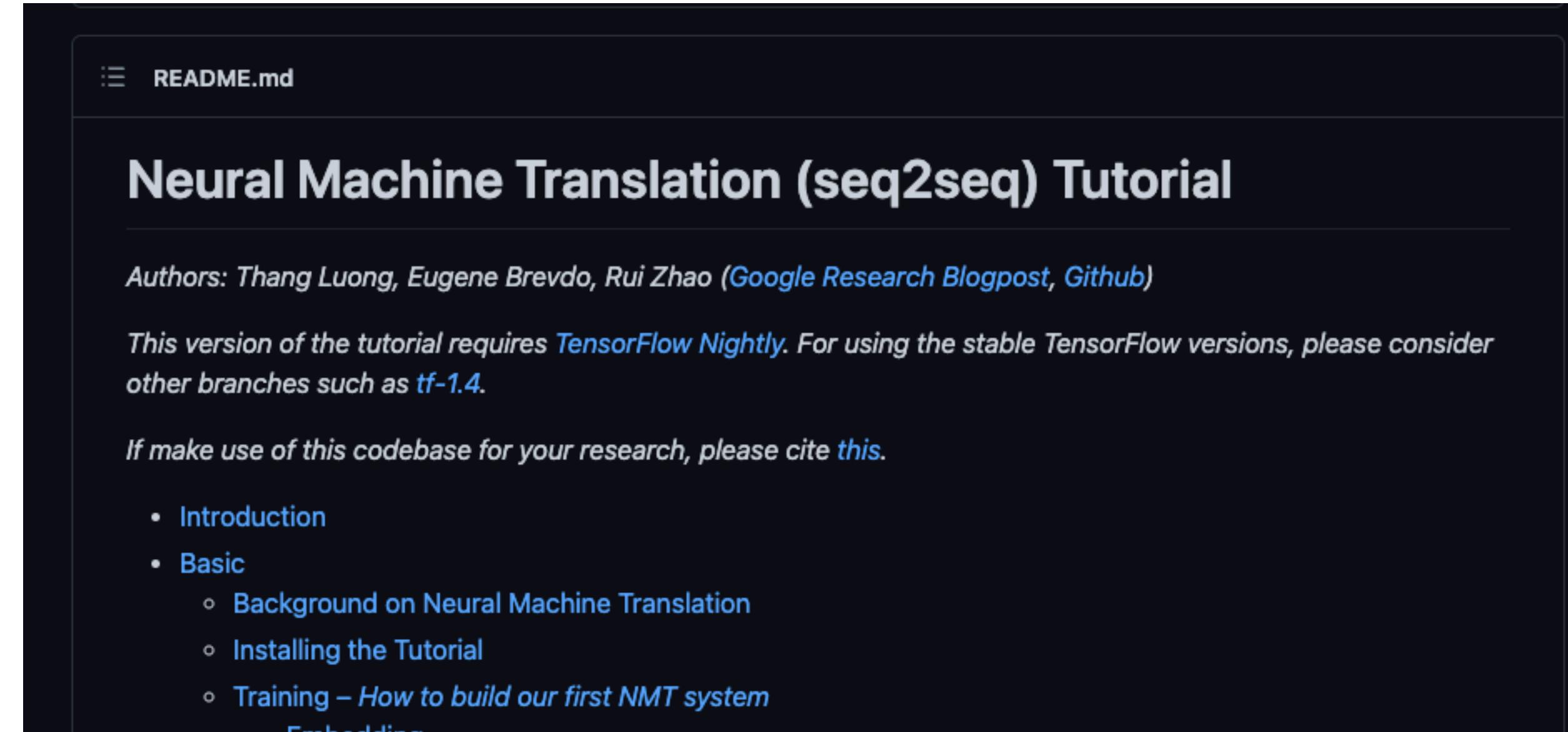
Diferencias:

SeqToSeq

Vs.

SeqToSeq +
Attention

TUTORIAL



The screenshot shows the first few lines of a README.md file. It starts with a header section titled "Neural Machine Translation (seq2seq) Tutorial". Below the title, it mentions the authors: Thang Luong, Eugene Brevdo, Rui Zhao, and provides links to a Google Research Blogpost and a GitHub repository. It also notes that the version shown requires TensorFlow Nightly and suggests using stable versions like tf-1.4. A citation section at the bottom encourages users to cite a specific paper if they use the codebase. A table of contents follows, listing sections such as Introduction, Basic, Background on Neural Machine Translation, Installing the Tutorial, and Training – How to build our first NMT system.

README.md

Neural Machine Translation (seq2seq) Tutorial

Authors: Thang Luong, Eugene Brevdo, Rui Zhao ([Google Research Blogpost](#), [Github](#))

This version of the tutorial requires [TensorFlow Nightly](#). For using the stable TensorFlow versions, please consider other branches such as [tf-1.4](#).

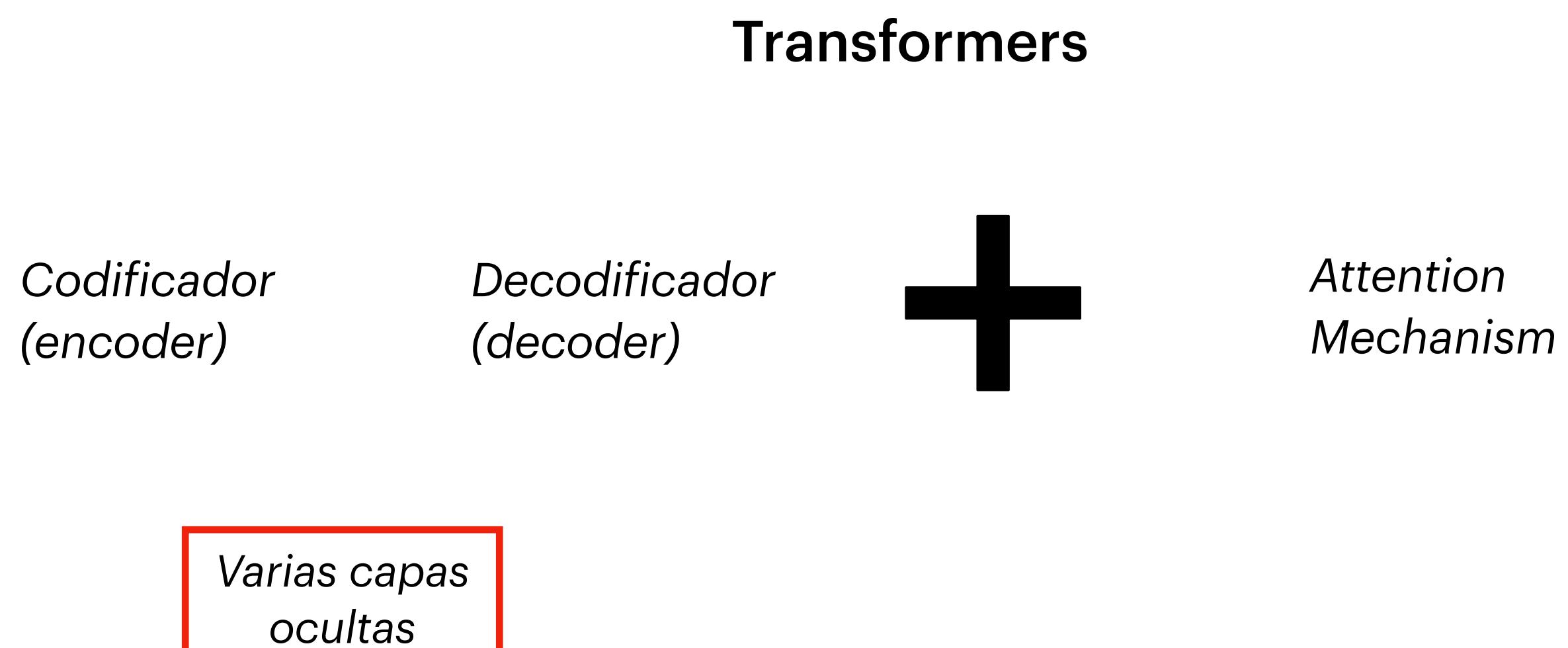
If make use of this codebase for your research, please cite [this](#).

- [Introduction](#)
- [Basic](#)
 - [Background on Neural Machine Translation](#)
 - [Installing the Tutorial](#)
 - [Training – How to build our first NMT system](#)
 - [Embedding](#)

Transformers and Attention Mechanisms

Transformers

Attention Is All You Need



Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

Attention Is All You Need

Ashish Vaswani* Noam Shazeer* Niki Parmar* Jakob Uszkoreit*
Google Brain Google Brain Google Research Google Research
avaswani@google.com noam@google.com nikip@google.com usz@google.com

Llion Jones* Aidan N. Gomez* † Lukasz Kaiser*
Google Research University of Toronto Google Brain
llion@google.com aidan@cs.toronto.edu lukaszkaiser@google.com

Ilia Polosukhin* ‡
ilia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

* Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Ilia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and implementing tensor2tensor, replacing our earlier codebase, greatly improving results and massively accelerating our research.

† Work performed while at Google Brain.

‡ Work performed while at Google Research.

Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática



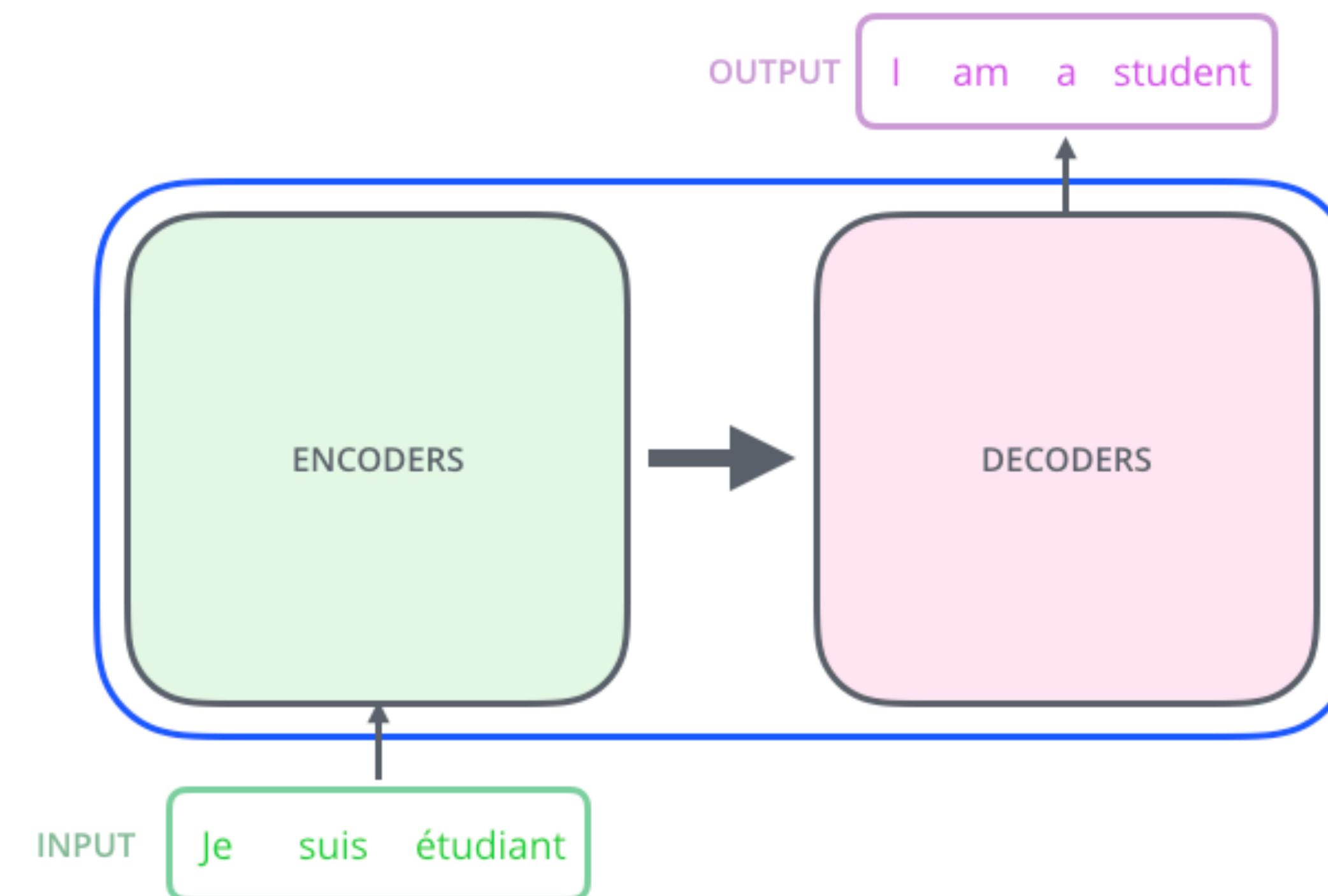
Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática



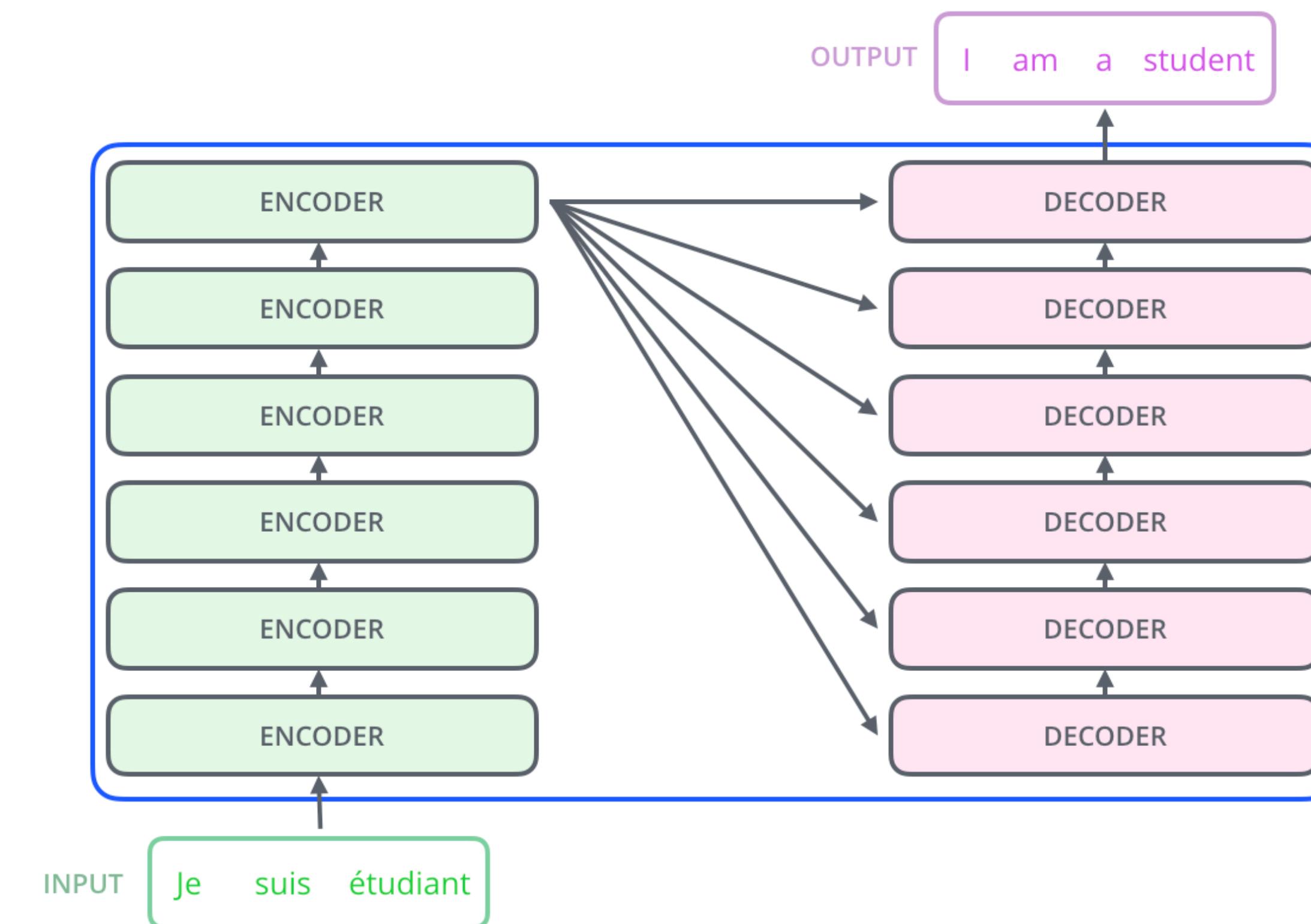
Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática



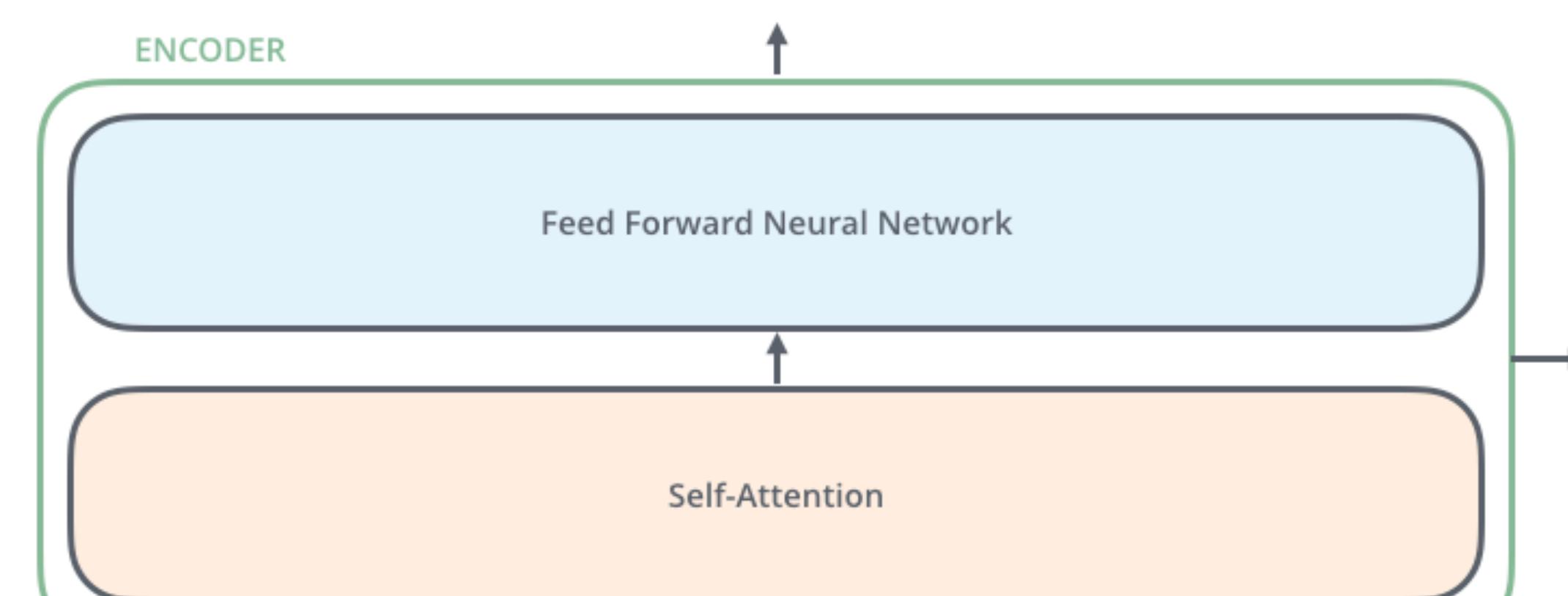
Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática



Cada encoder es igual → Pero no comparten pesos

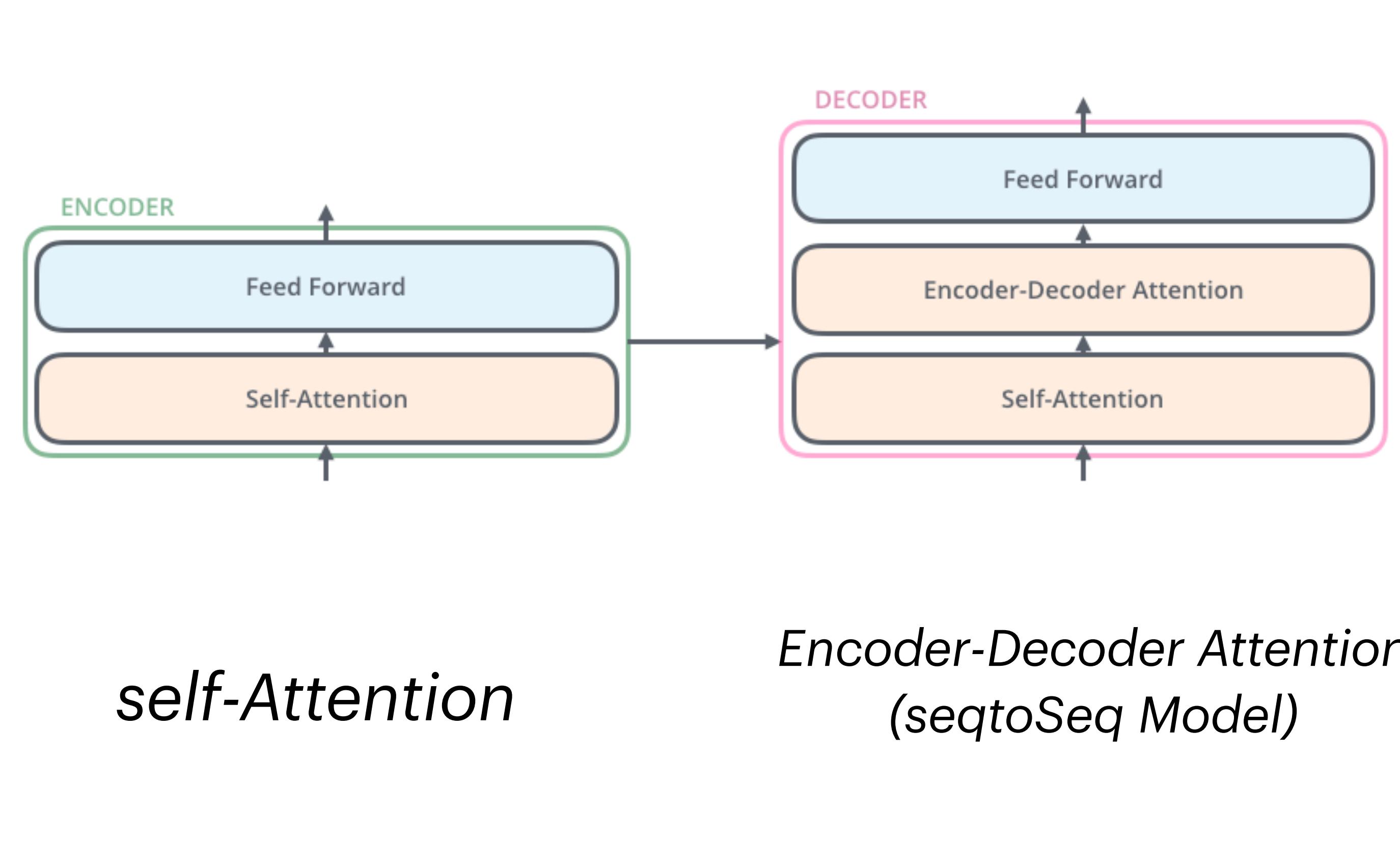
Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

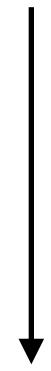


Encoder-Decoder

Self-Attention

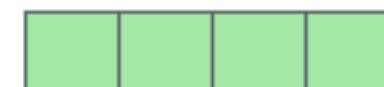
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

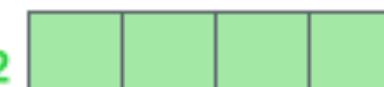


Ej.: Traducción automática

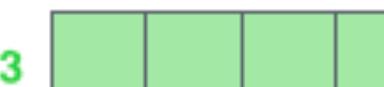
Word Embedding

x_1 

Je

x_2 

suis

x_3 

étudiant

- Embedding de 512 dimensiones
- Longitud de la oración más larga en nuestro conjunto de datos

Encoder-Decoder

Self-Attention

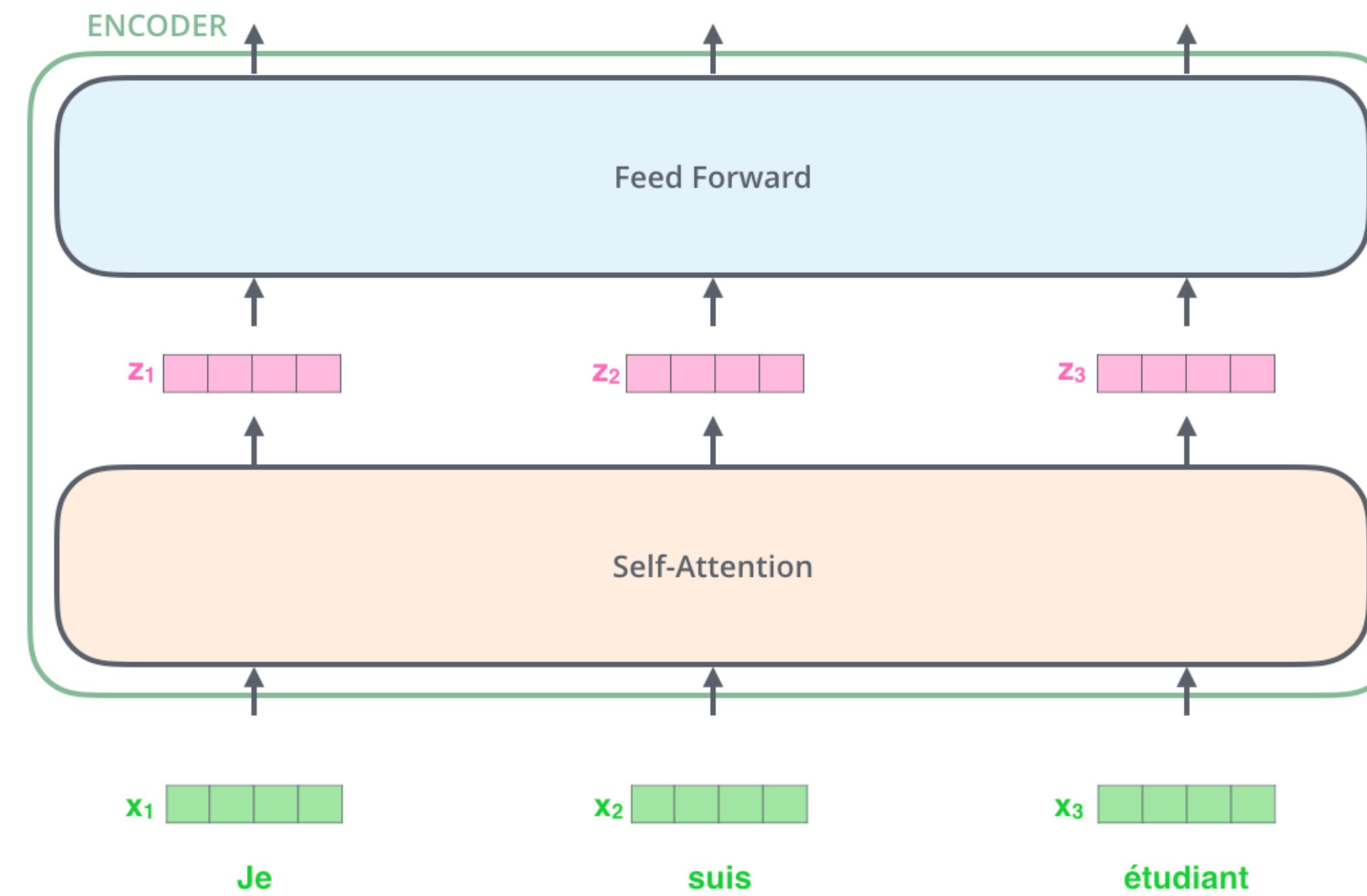
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Word Embedding

Fluyen en paralelo (clave)



Encoder-Decoder

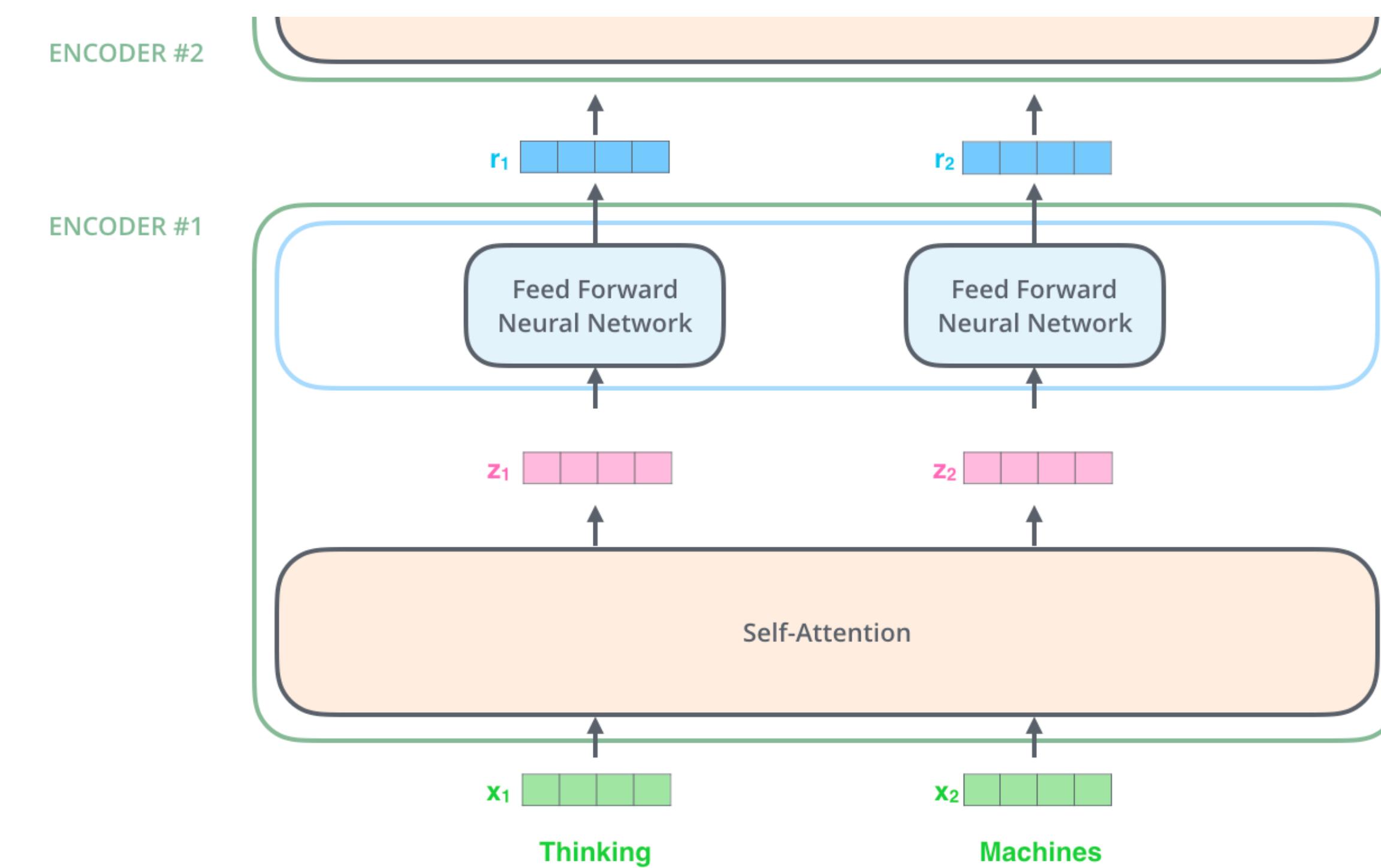
Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Now We're Encoding!



Encoder-Decoder

Self-Attention

Transformers

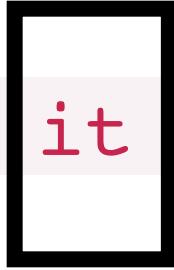
Atención para aumentar la velocidad en entrenamiento de modelos.



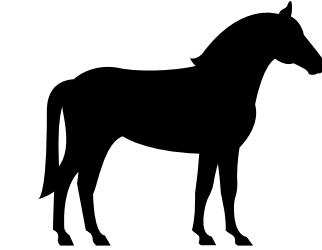
Ej.: Traducción automática

Self-Attention at a High Level

The animal didn't cross the street because **it** was too tired



¿A qué nos referimos con 'it'?



Encoder-Decoder

Self-Attention

Transformers

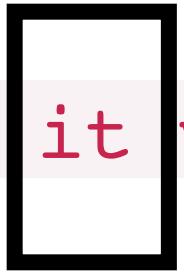
Atención para aumentar la velocidad en entrenamiento de modelos.



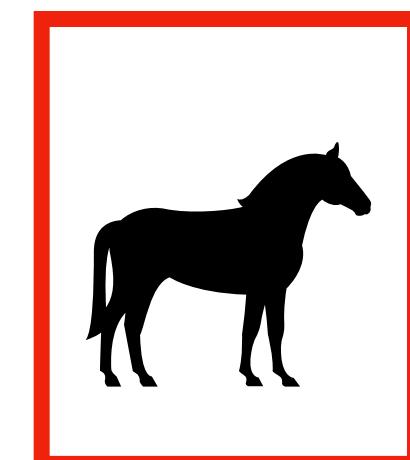
Ej.: Traducción automática

Self-Attention at a High Level

The animal didn't cross the street because it was too tired



¿A qué nos referimos con 'it'?

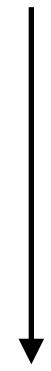


Encoder-Decoder

Self-Attention

Transformers

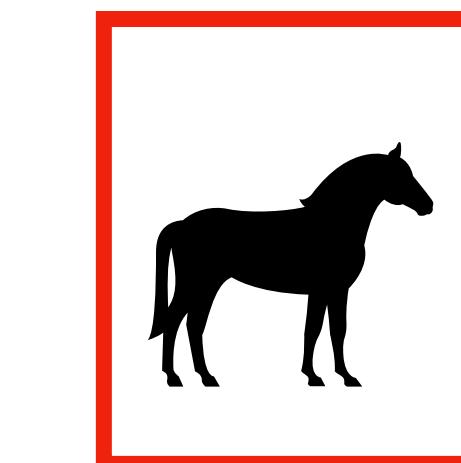
Atención para aumentar la velocidad en entrenamiento de modelos.



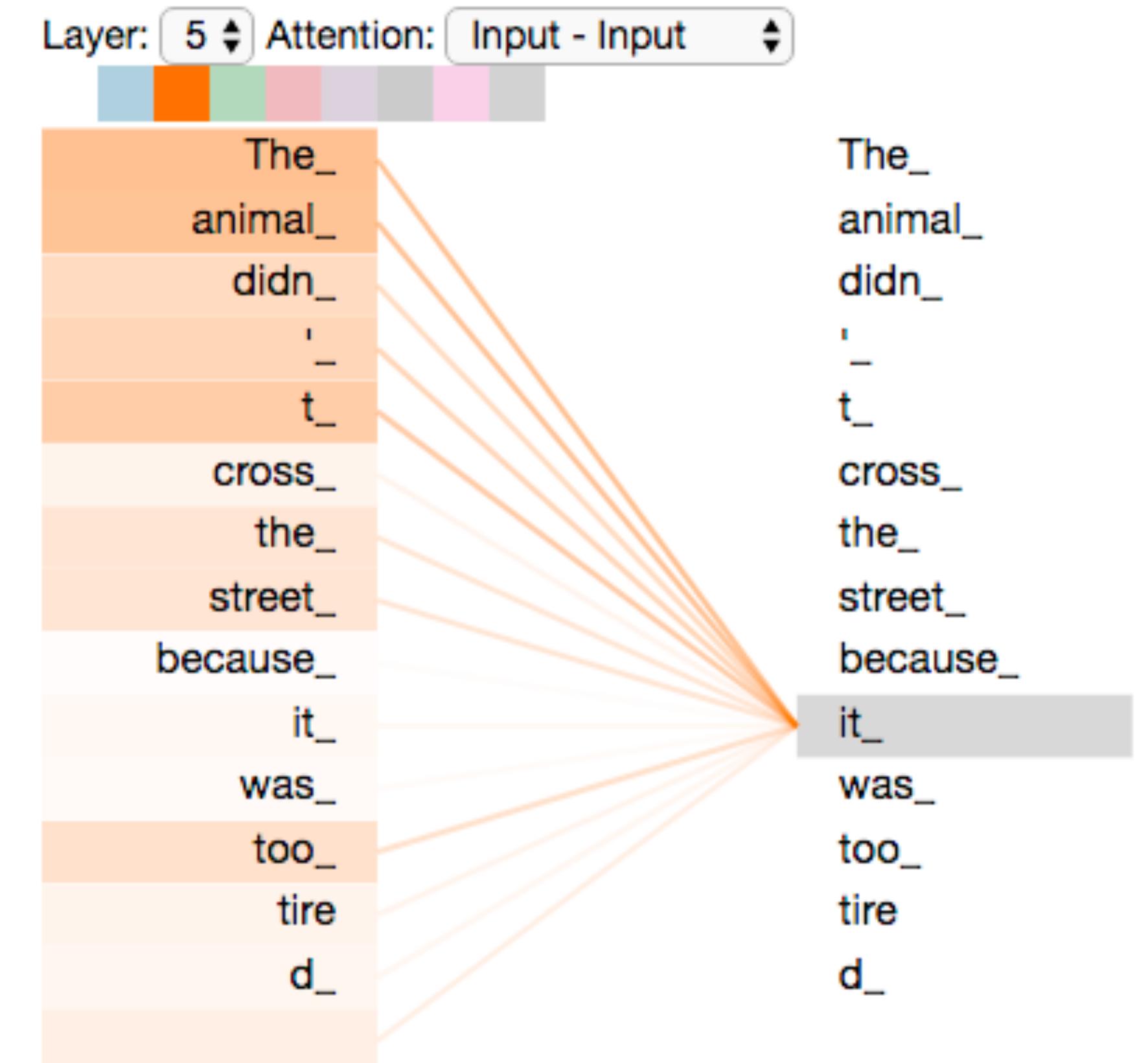
Ej.: Traducción automática



¿A qué nos referimos con 'it'?



Self-Attention at a High Level



Encoder-Decoder

Self-Attention

Calculo vectorial

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.



Ej.: Traducción automática

$q_1, k_1, v_1 \in \text{palabras}$

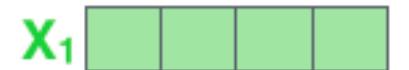
Self-Attention in Detail Implementación de vectores

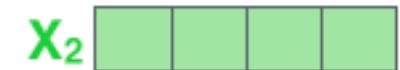
Input

Thinking

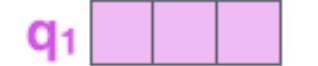
Machines

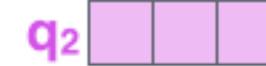
Embedding

X_1 

X_2 

Queries

q_1 

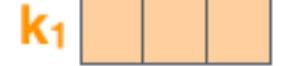
q_2 

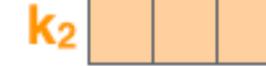
$$q_1 = X_1 \cdot W_Q$$

Vectores consulta, clave y valor

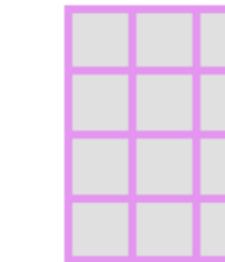
Vectores con dimensión de 64

Keys

k_1 

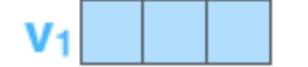
k_2 

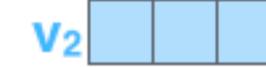
$$k_1 = X_1 \cdot W_K$$



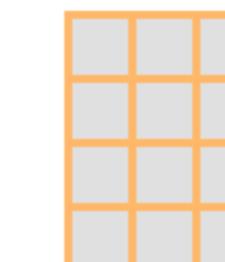
W_Q

Values

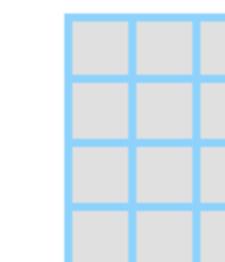
v_1 

v_2 

$$v_1 = X_1 \cdot W_V$$



W_K



W_V

Encoder-Decoder

Self-Attention

Calculo vectorial

Transformers

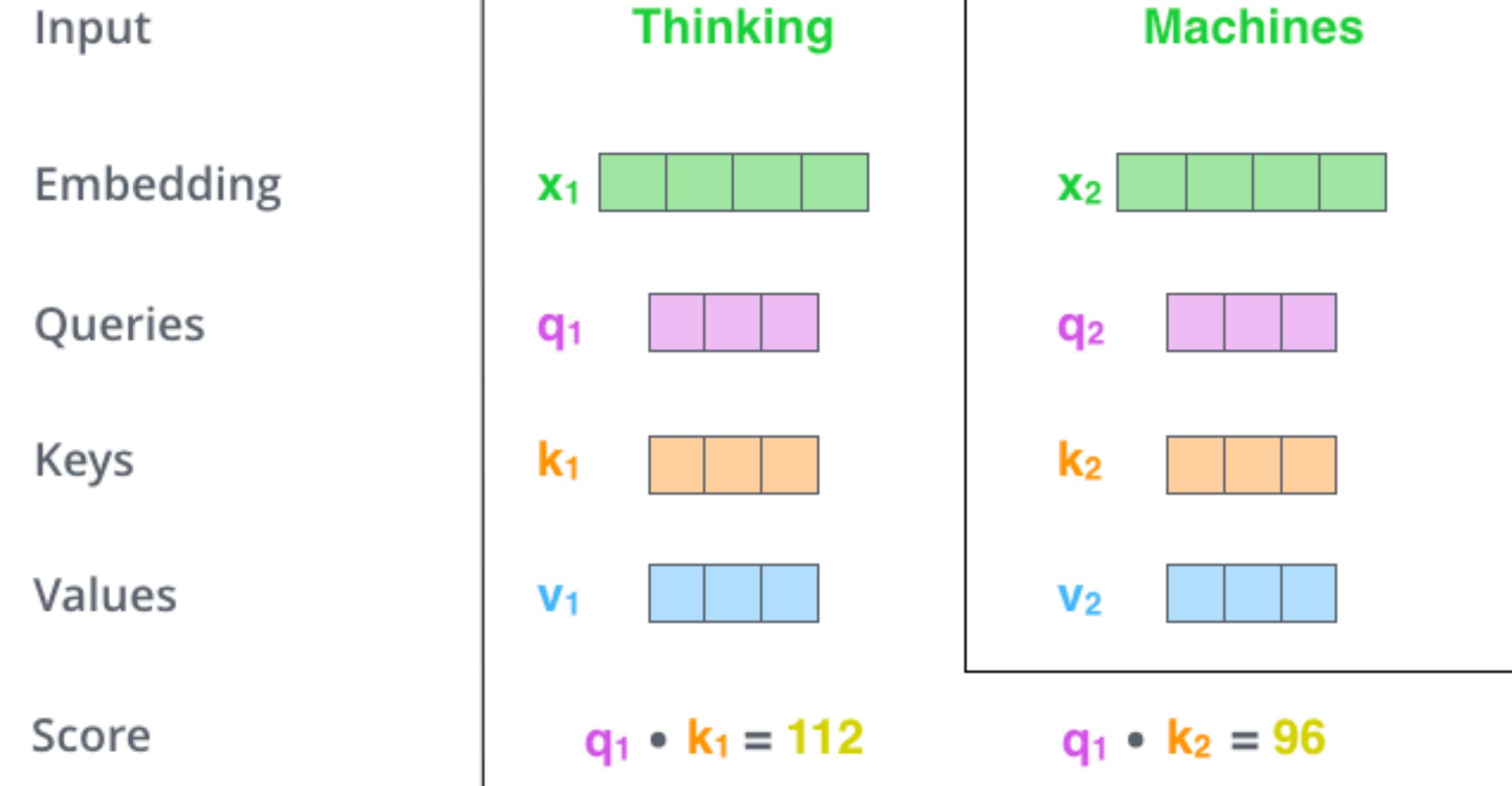
Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

$q_1, k_1, v_1 \in \text{palabras}$

Self-Attention in Detail Calcular puntuación

(vectores consulta, clave y valor)



Encoder-Decoder

Self-Attention

Calculo vectorial

Transformers

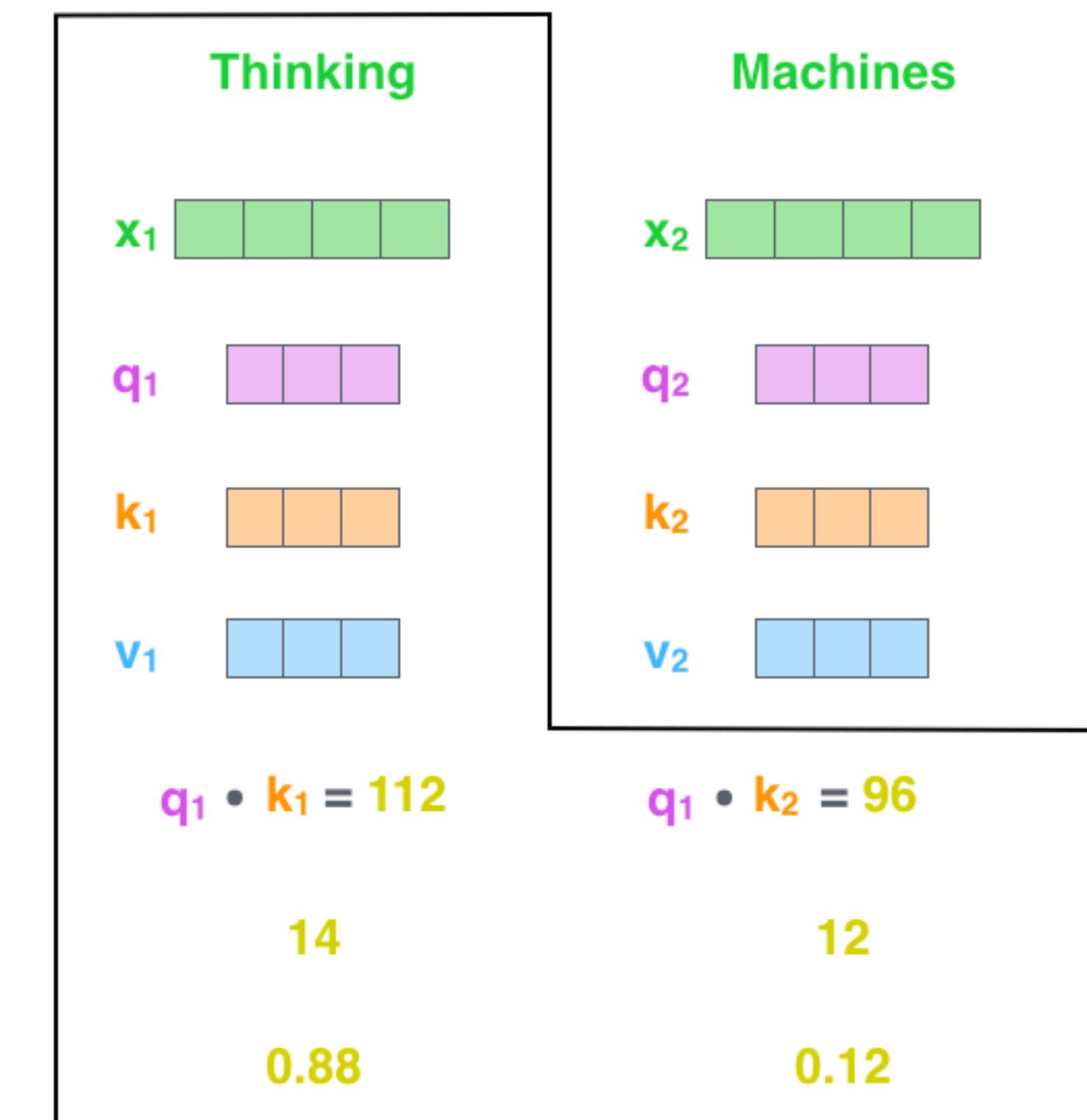
Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

$q_1, k_1, v_1 \in \text{palabras}$

Self-Attention in Detail Normalizar puntuación

Input
Embedding
Queries
Keys
Values
Score
Divide by 8 ($\sqrt{d_k}$)
Softmax



Gradientes estables (sqrt -> 8)

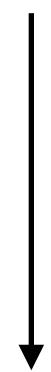
Encoder-Decoder

Self-Attention

Calculo vectorial

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

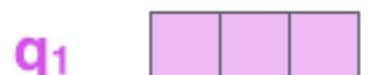


Ej.: Traducción automática

$q_1, k_1, v_1 \in \text{palabras}$

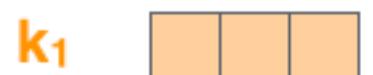
Self-Attention in Detail Preparar para resumir

Queries

q_1 

q_2 

Keys

k_1 

k_2 

Values

v_1 

v_2 

Score

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

Divide by 8 ($\sqrt{d_k}$)

14

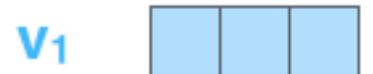
12

Softmax

0.88

0.12

Softmax X Value

v_1 

v_2 

Sum

z_1 

z_2 

Se mantienen intactos los valores de palabras donde se quiere centrarse

Se suman los vectores de valores ponderados (self-attention vector)

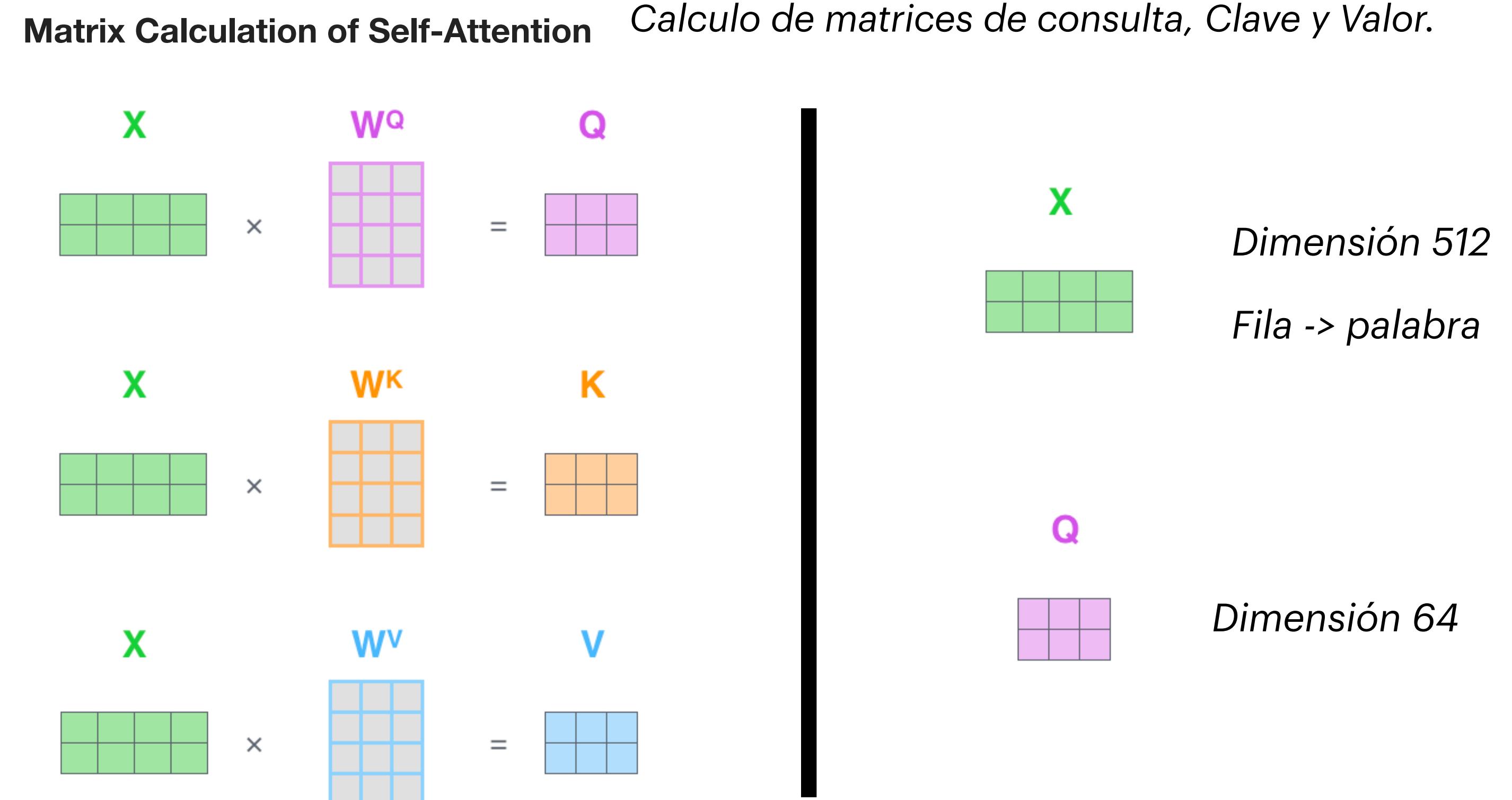
Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática



Encoder-Decoder

Self-Attention

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Matrix Calculation of Self-Attention

$$\text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V} = \mathbf{Z}$$

The diagram illustrates the matrix calculation of self-attention. It shows three square matrices: Q (purple), K^T (orange), and V (blue). The Q matrix has 3 columns and 3 rows. The K^T matrix has 3 columns and 3 rows. The V matrix has 3 columns and 3 rows. A multiplication symbol (×) is placed between Q and K^T. A division symbol (÷) with the square root of d_k is placed between the multiplication and the right parenthesis of the softmax function. An equals sign (=) is followed by a Z matrix (pink) with 3 columns and 3 rows.

Cálculo de la autoatención en forma matricial

Encoder-Decoder

Multi-Head Attention

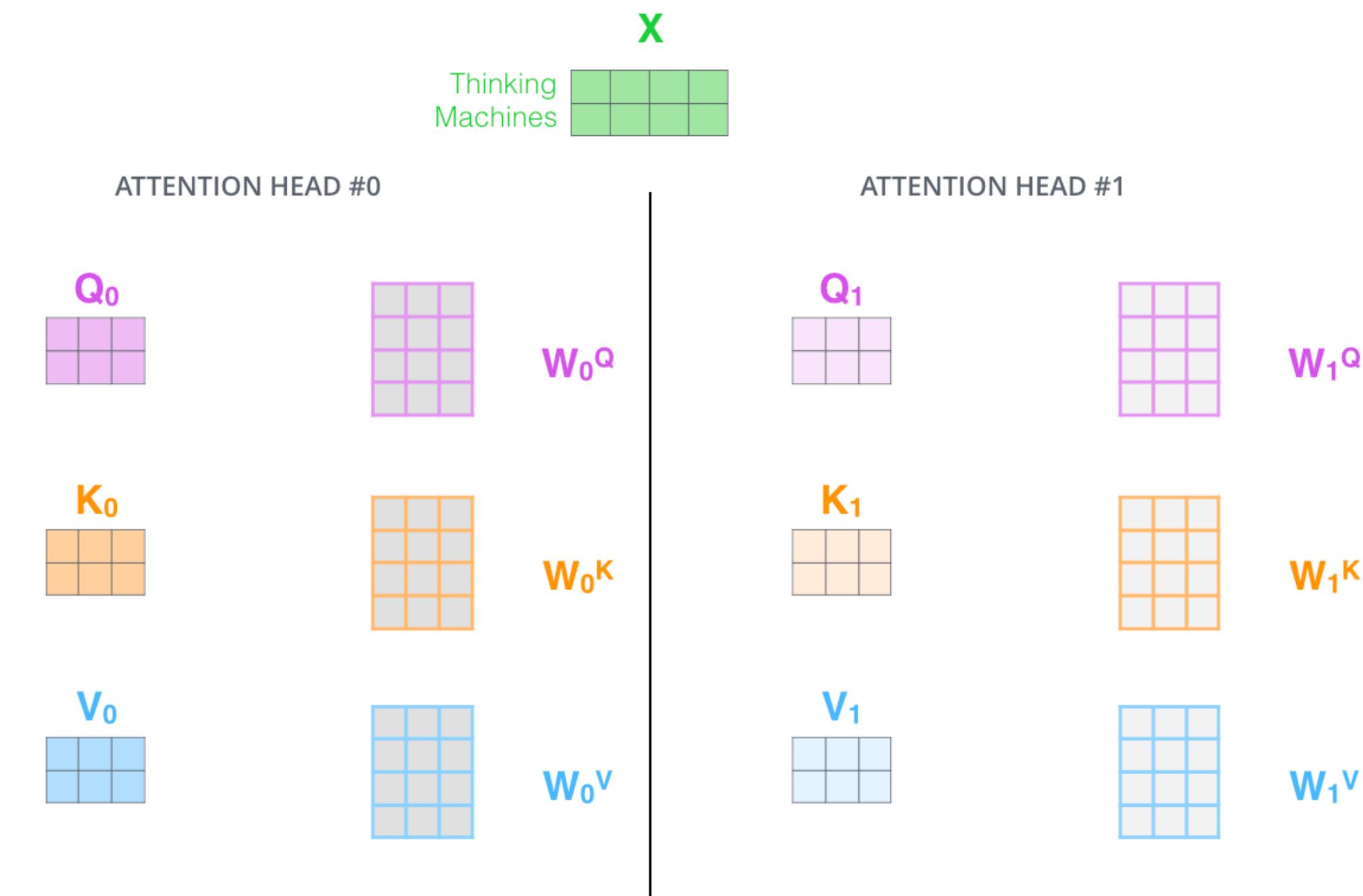
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Atención “multidireccional”

The Beast With Many Heads Se usan 8 cabezas de atención por cada codificador/decodificador



Encoder-Decoder

Multi-Head Attention

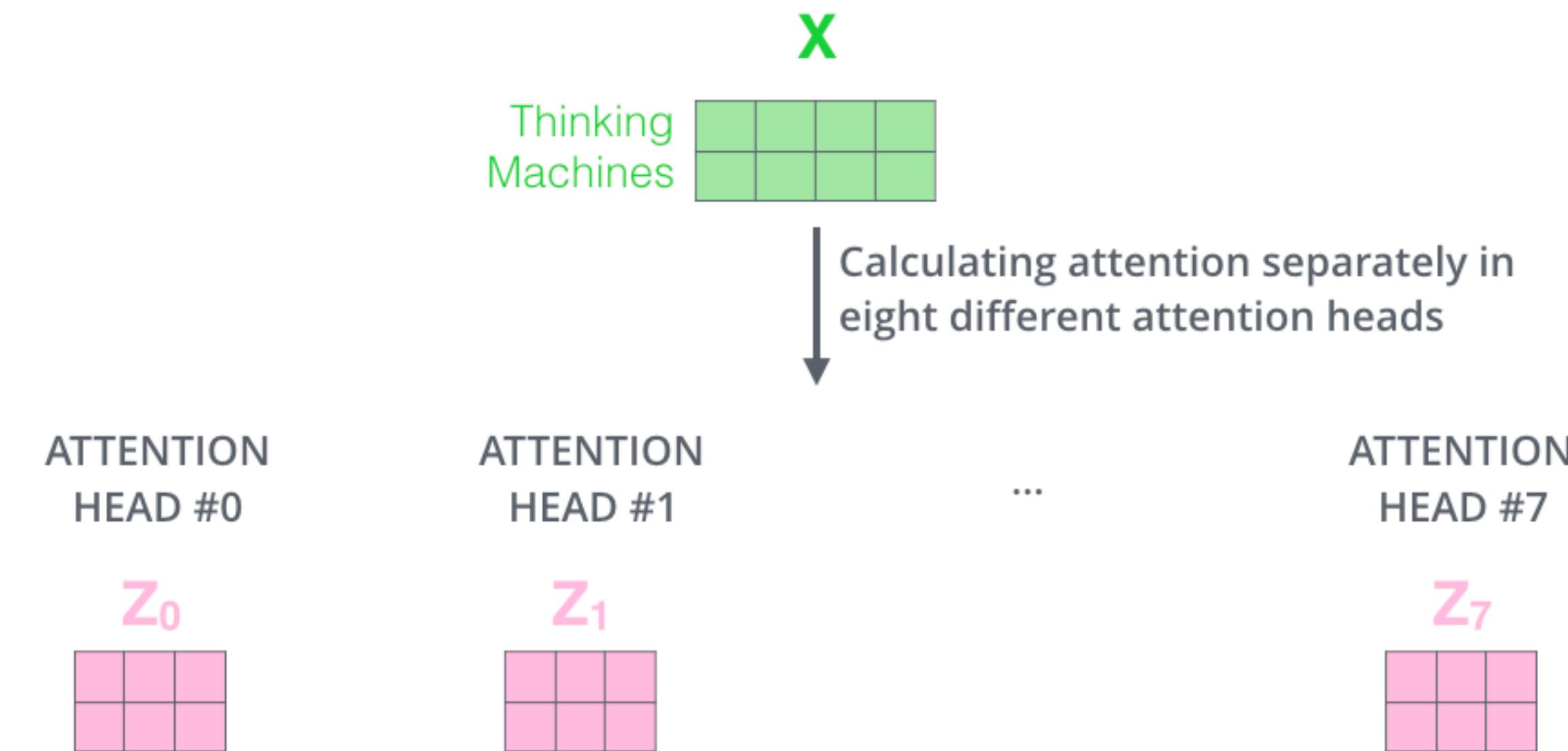
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Atención “multidireccional”

The Beast With Many Heads Se usan 8 cabezas de atención por cada codificador/decodificador



Encoder-Decoder

Multi-Head Attention

Transformers

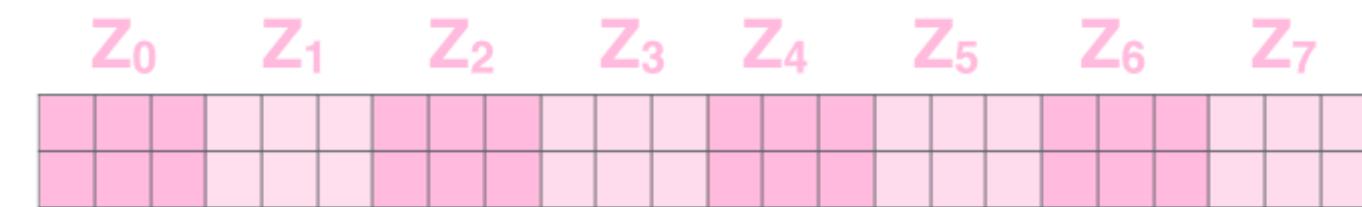
Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Atención “multidireccional”

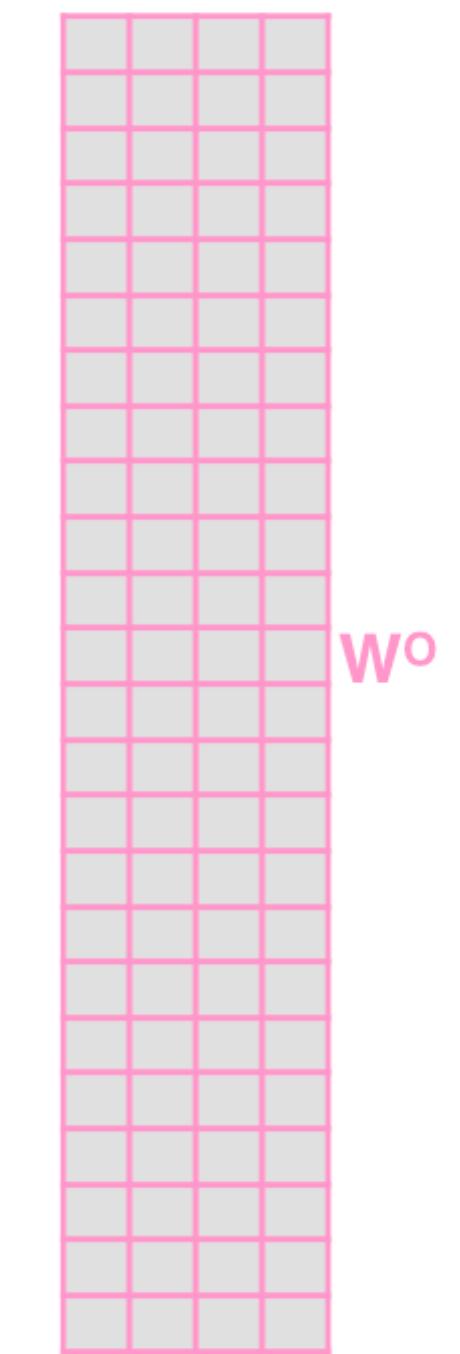
The Beast With Many Heads Se usan 8 cabezas de atención por cada codificador/decodificador

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \hline \end{matrix}$$

Encoder-Decoder

Multi-Head Attention

Transformers

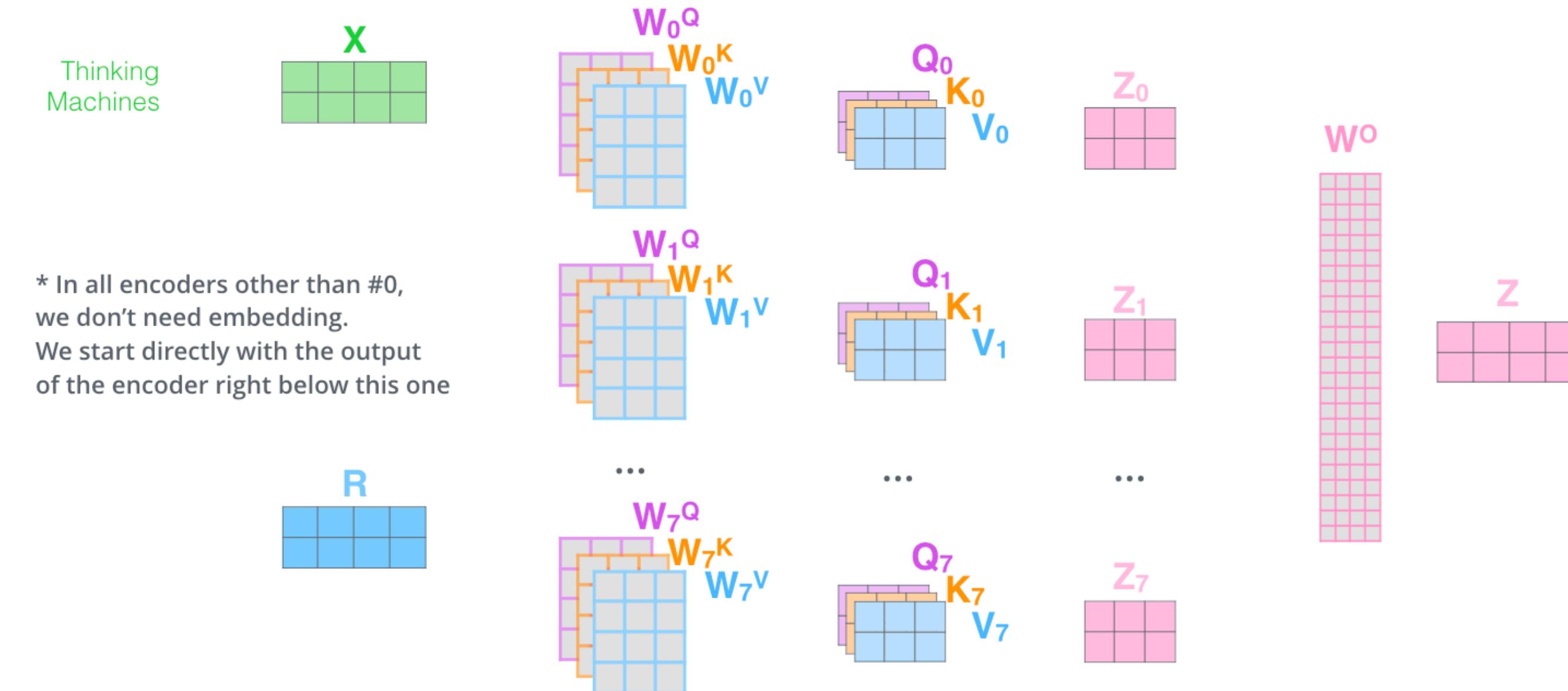
Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Atención “multidireccional”

The Beast With Many Heads Se usan 8 cabezas de atención por cada codificador/decodificador

- 1) This is our input sentence* each word*
- 2) We embed
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



Encoder-Decoder

Multi-Head Attention

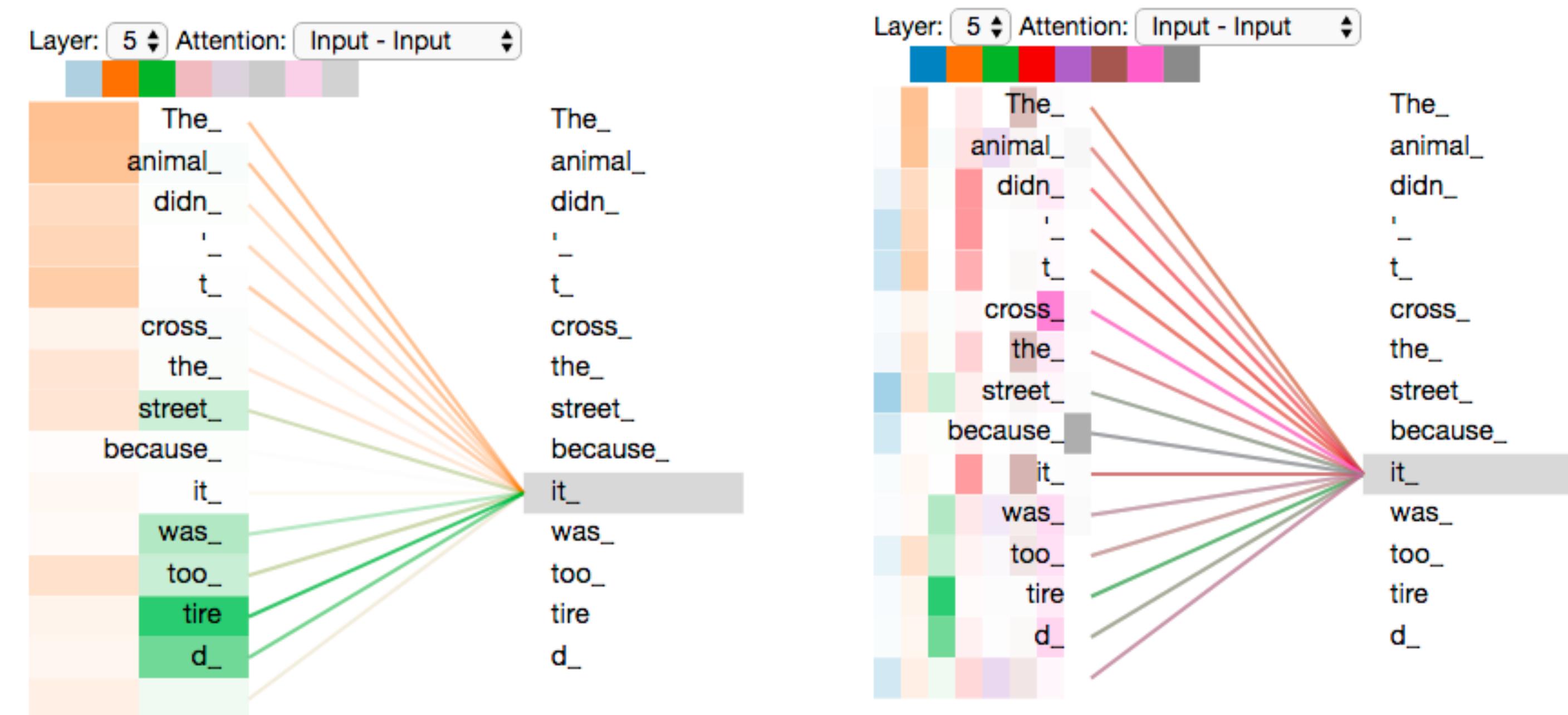
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Atención “multidireccional”

The Beast With Many Heads Se usan 8 cabezas de atención por cada codificador/decodificador



Encoder-Decoder

Positional Encoding

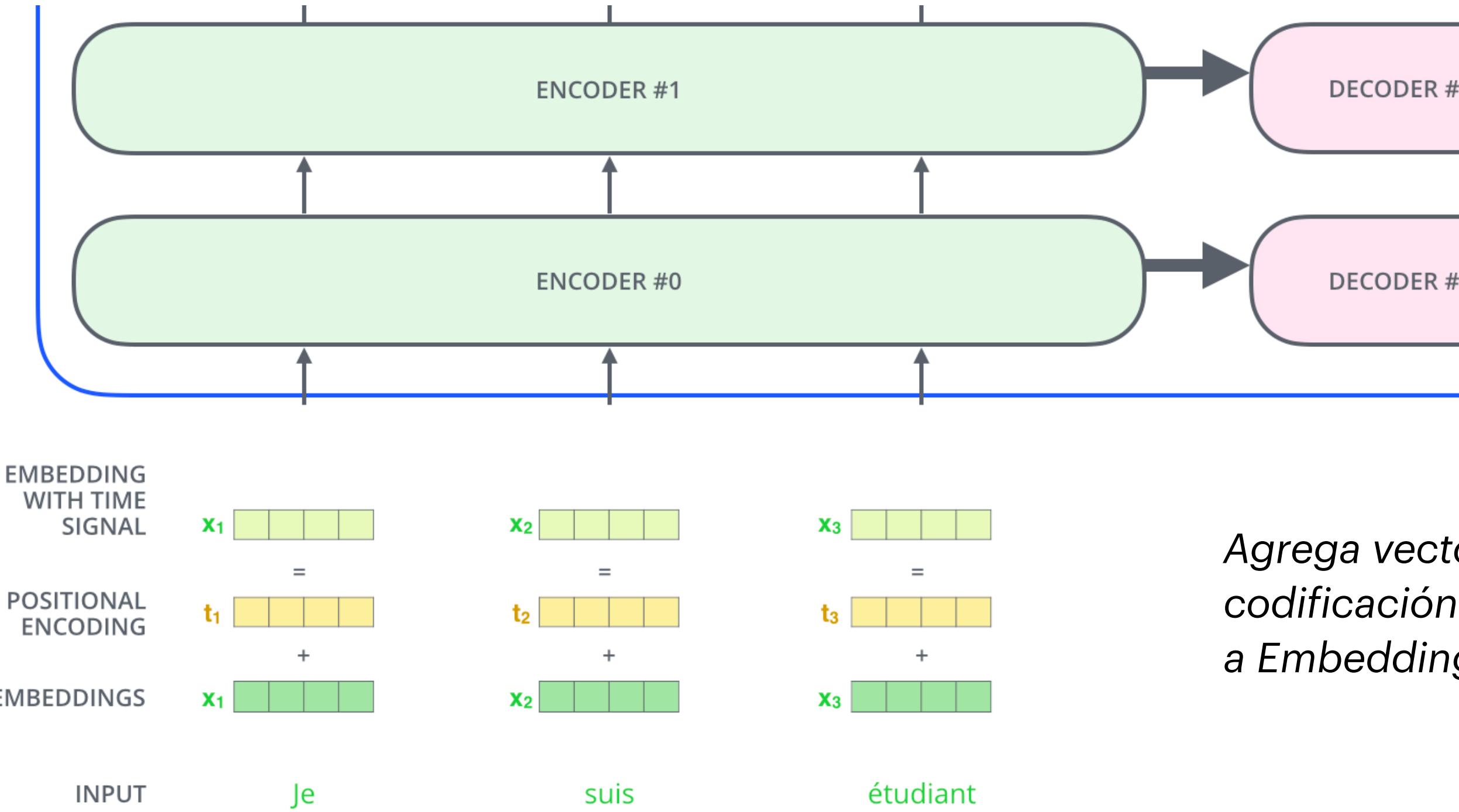
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

'Idea de orden en las palabras'

The Order of The Sequence Using Positional Encoding

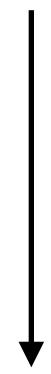


Encoder-Decoder

Positional Encoding

Transformers

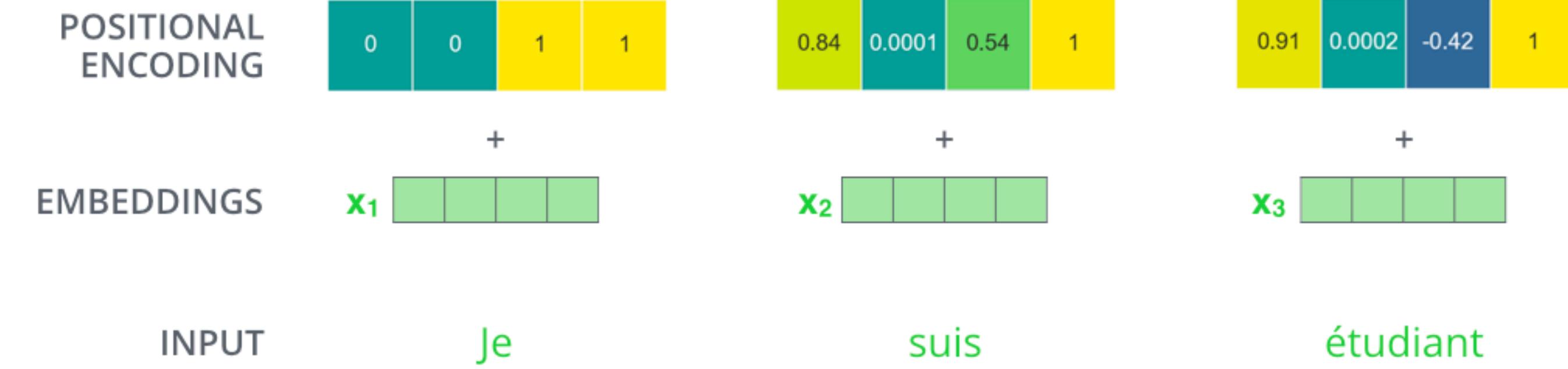
Atención para aumentar la velocidad en entrenamiento de modelos.



Ej.: Traducción automática

'Idea de orden en las palabras'

The Order of The Sequence Using Positional Encoding



Encoder-Decoder

Positional Encoding

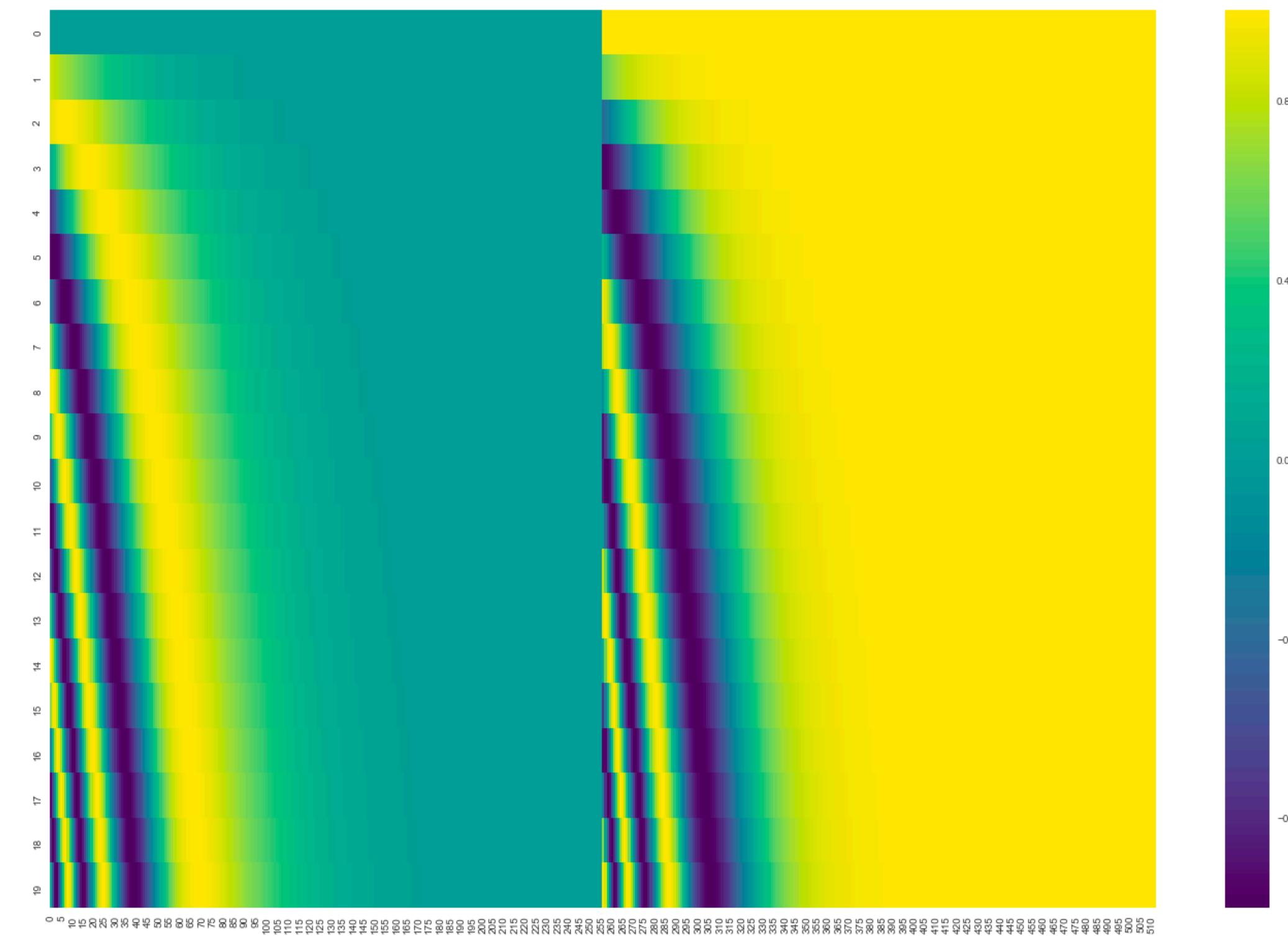
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

'Idea de orden en las palabras'

The Order of The Sequence Using Positional Encoding



Función
seno

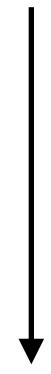
Función
coseno

Encoder-Decoder

Positional Encoding

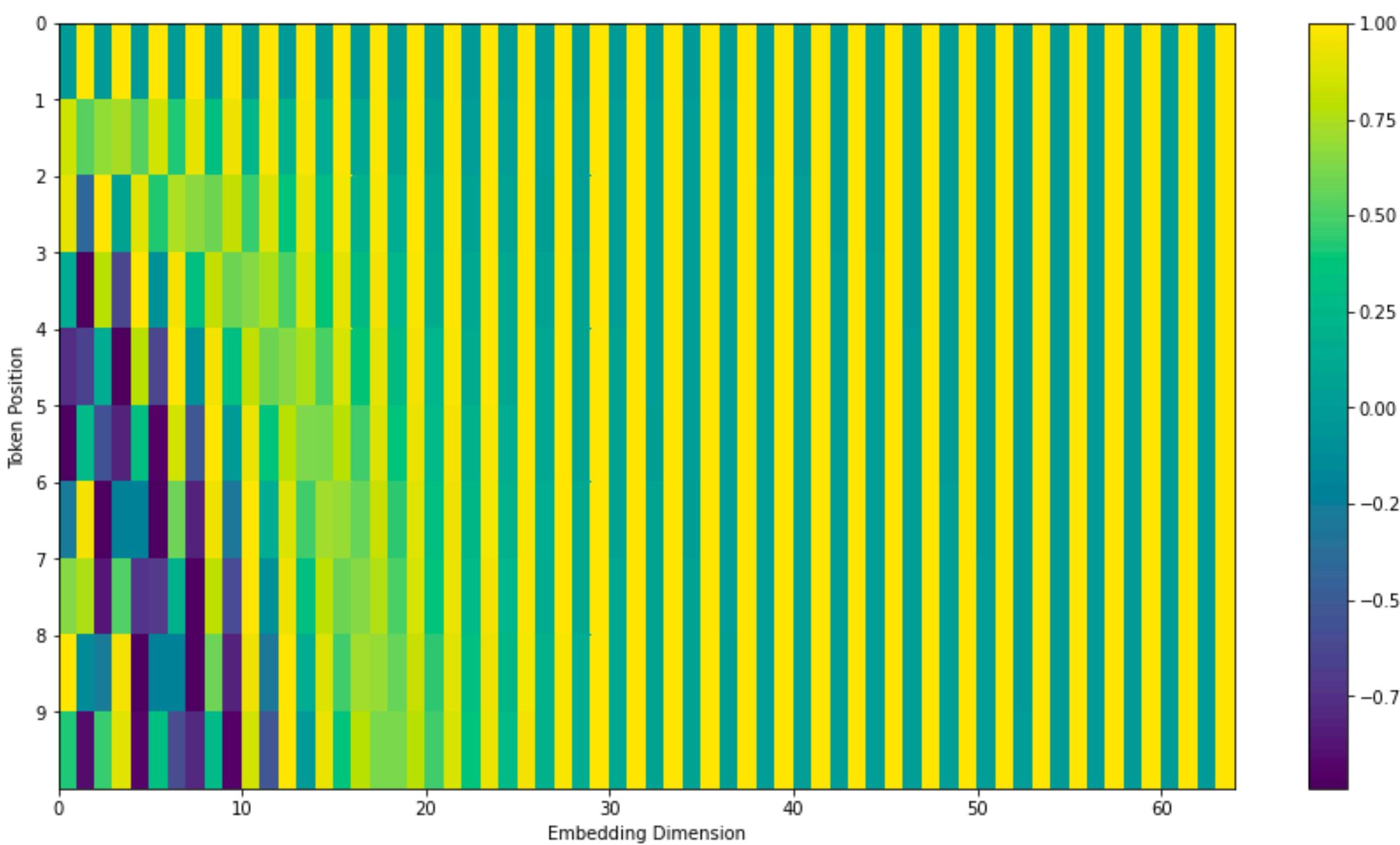
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.



Ej.: Traducción automática
'Idea de orden en las palabras'

The Order of The Sequence Using Positional Encoding



20 palabras (filas)

512 dimensiones embedding

Valores entre -1 y 1

Entrelazando las dos señales

Función seno

Función coseno

Encoder-Decoder

The Residuals

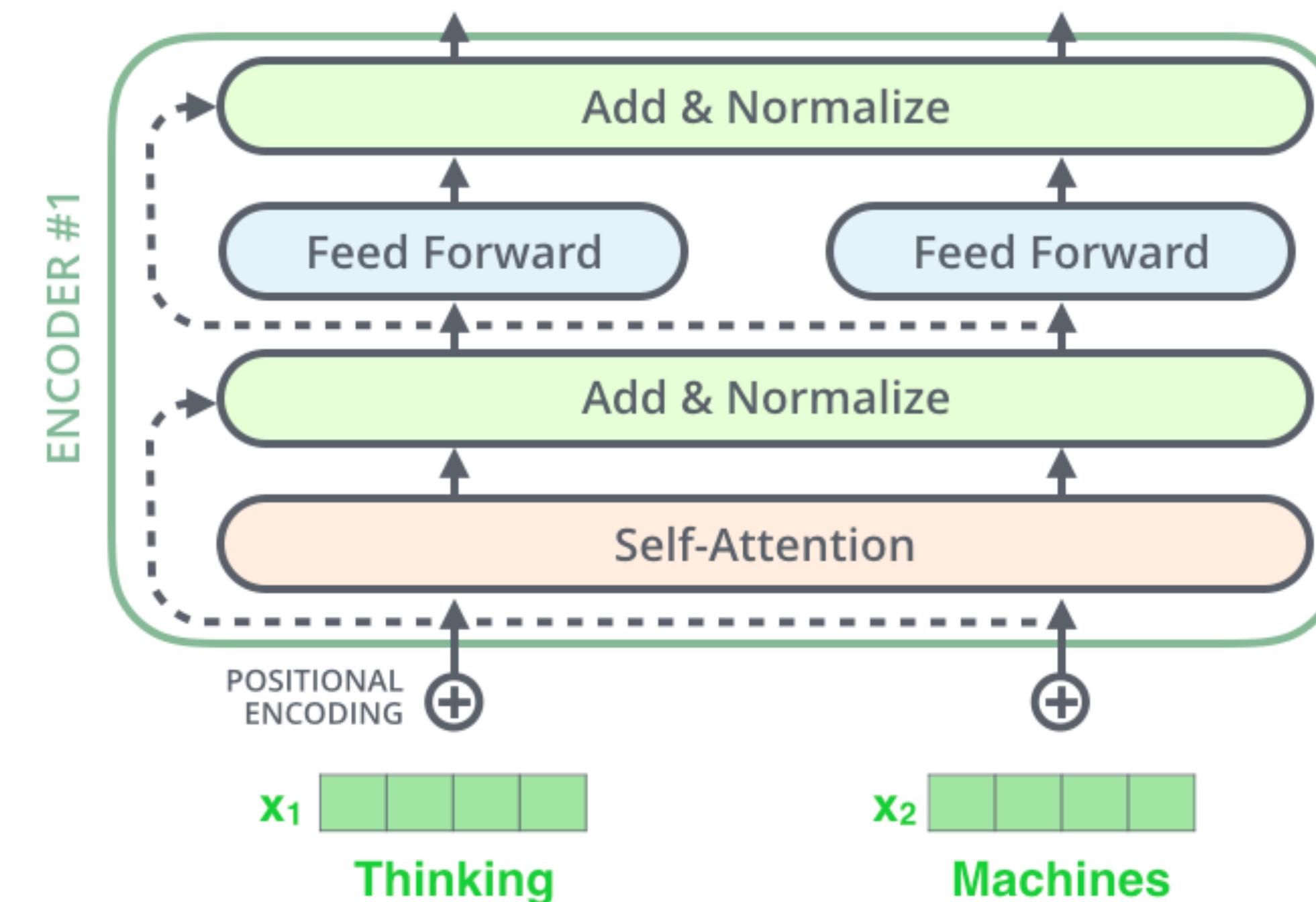
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Normalización entre atención, ffnn

The Order of The Sequence Using Positional Encoding



Conexión residual a una capa de normalización

Encoder-Decoder

The Residuals

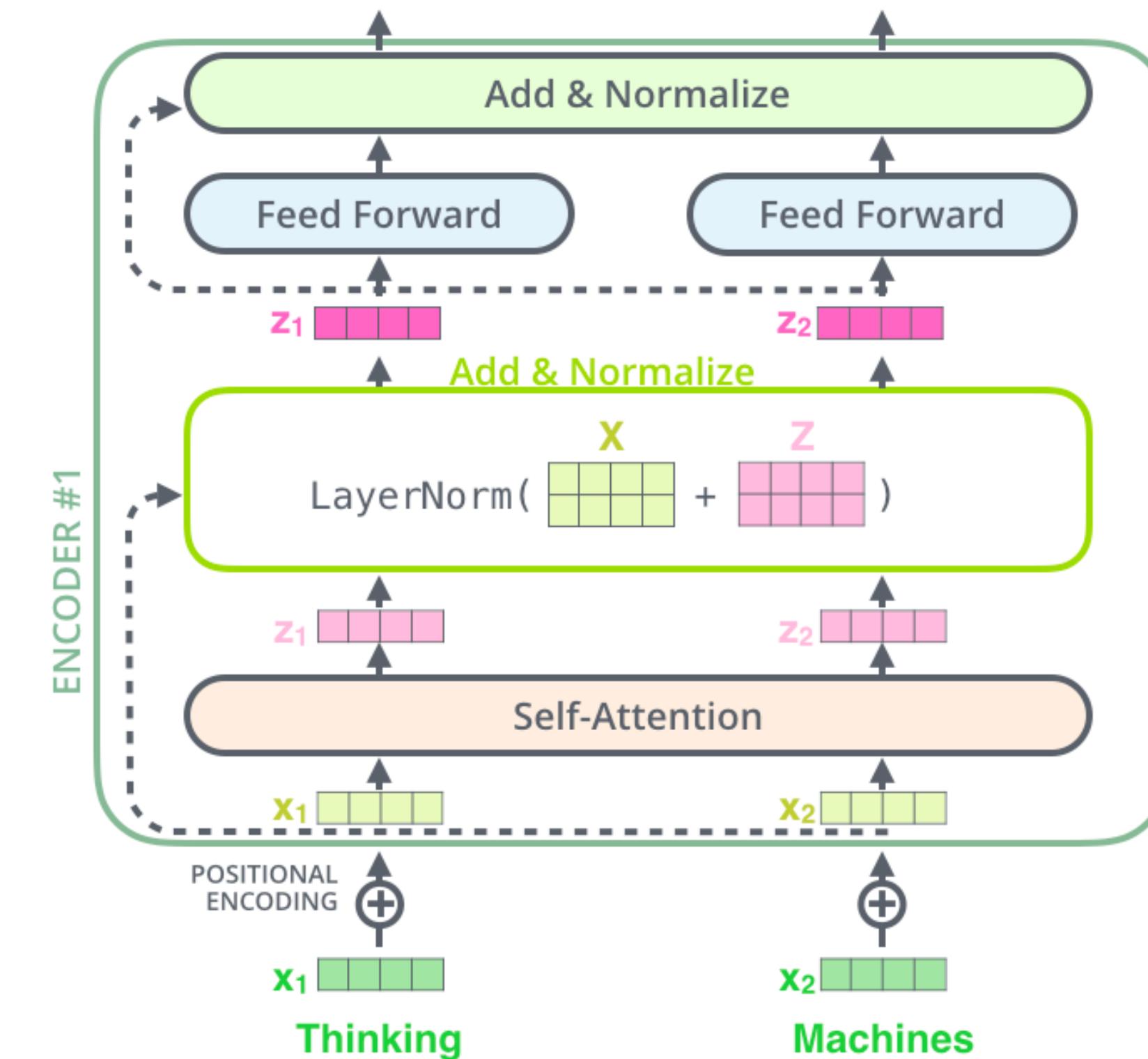
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Normalización entre atención, ffnn

The Order of The Sequence Using Positional Encoding



Visualización de parámetros en capa de norma

Esto se extiende a las subcapas del decodificador

Encoder-Decoder

The Residuals

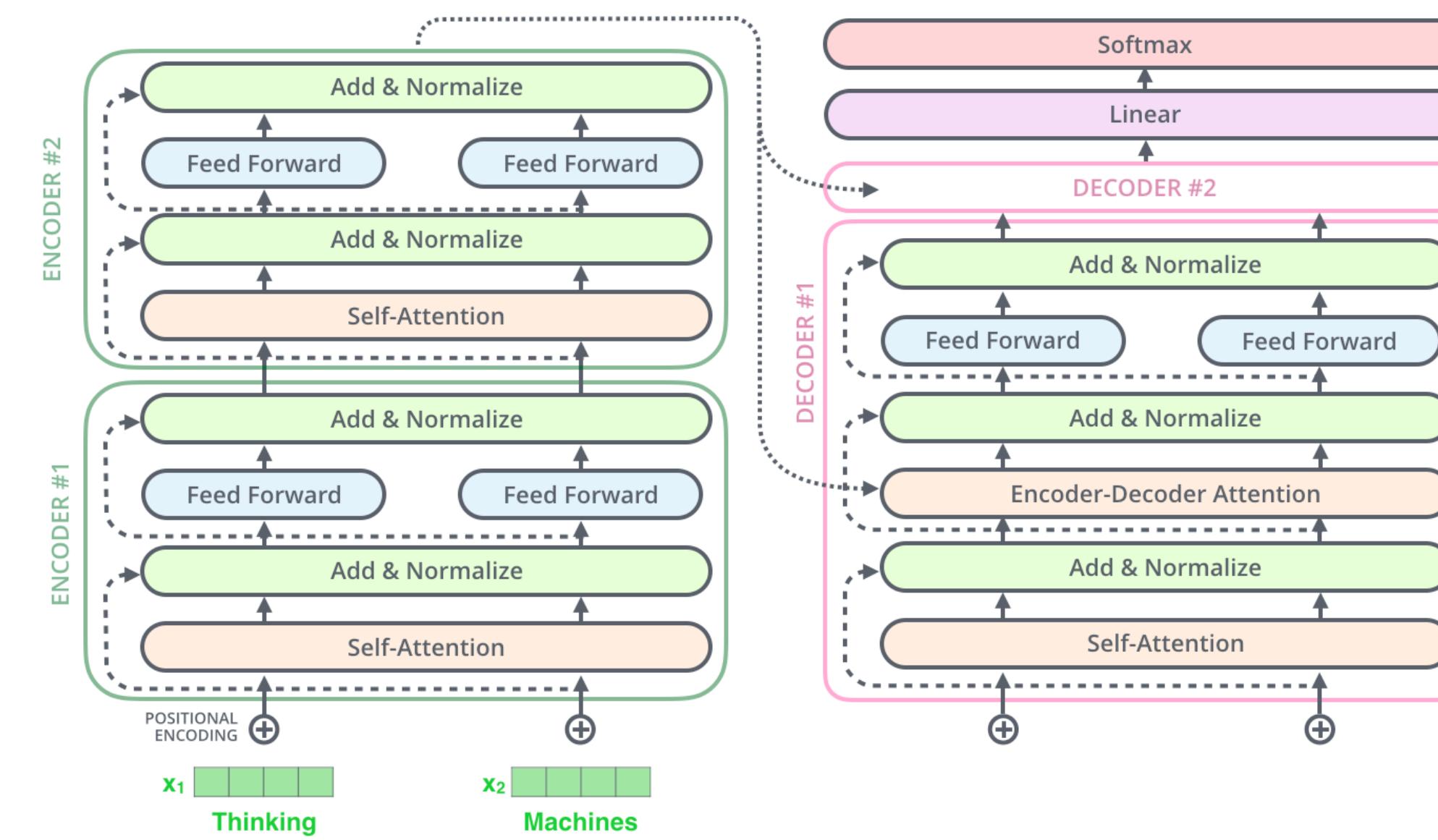
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Normalización entre atención, ffnn

The Order of The Sequence Using Positional Encoding



Visualización de parámetros en capa de norma

Esto se extiende a las subcapas del decodificador

Encoder-Decoder

The Decoder Side

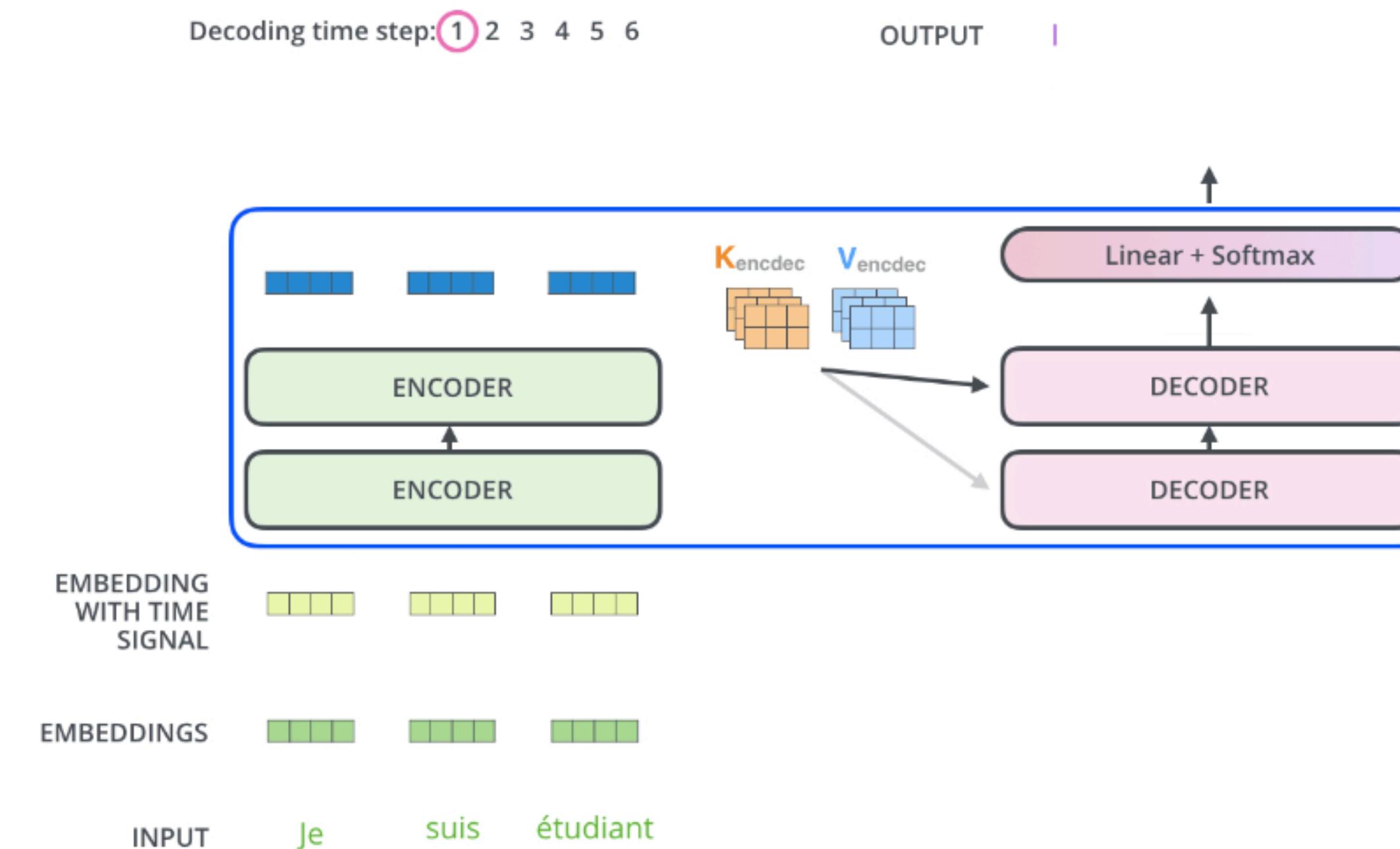
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Normalización entre atención, ffnn

The Decoder



Encoder-Decoder

The Decoder Side

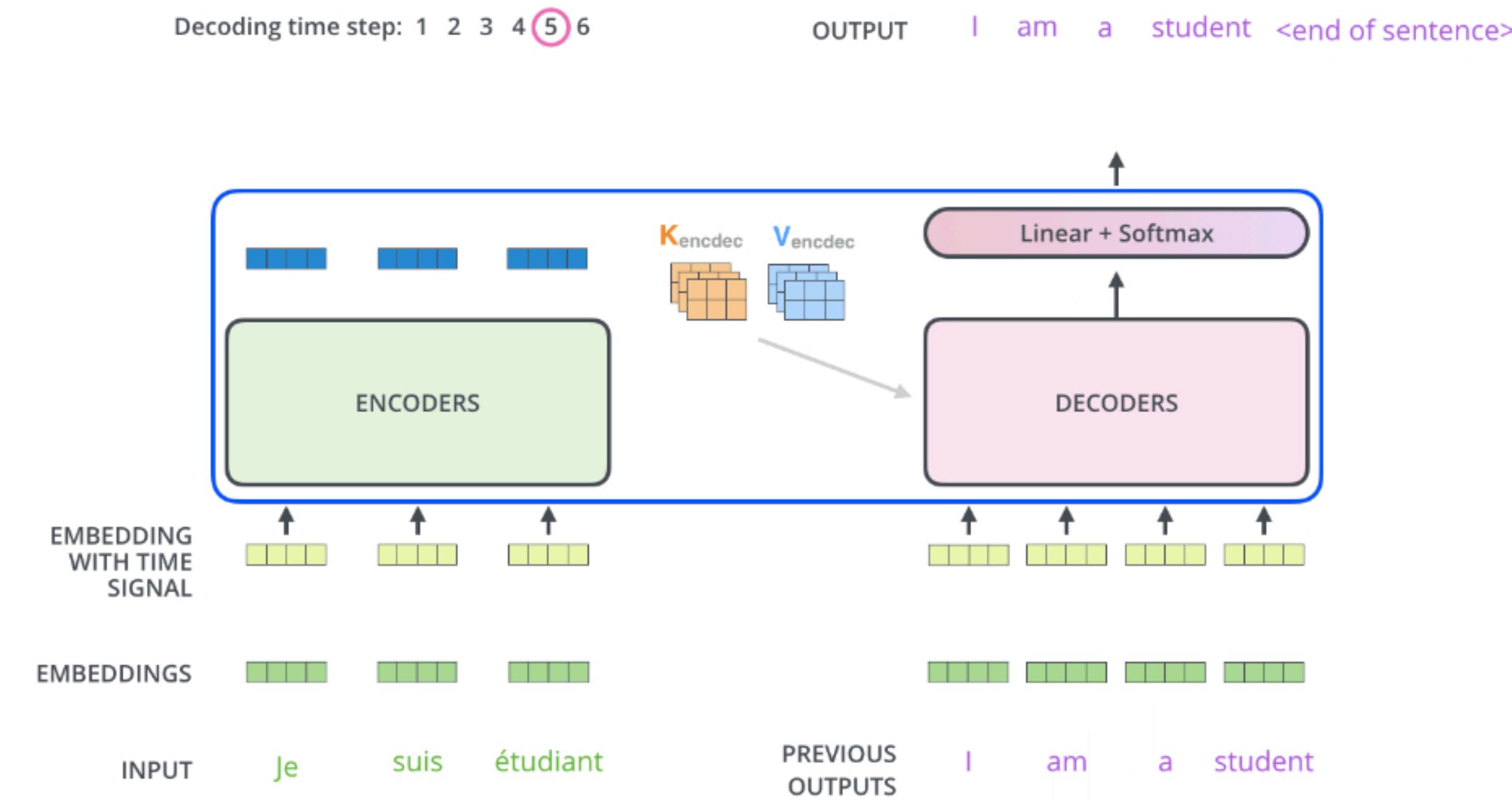
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Normalización entre atención, ffnn

The Decoder



Encoder-Decoder

The Final

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

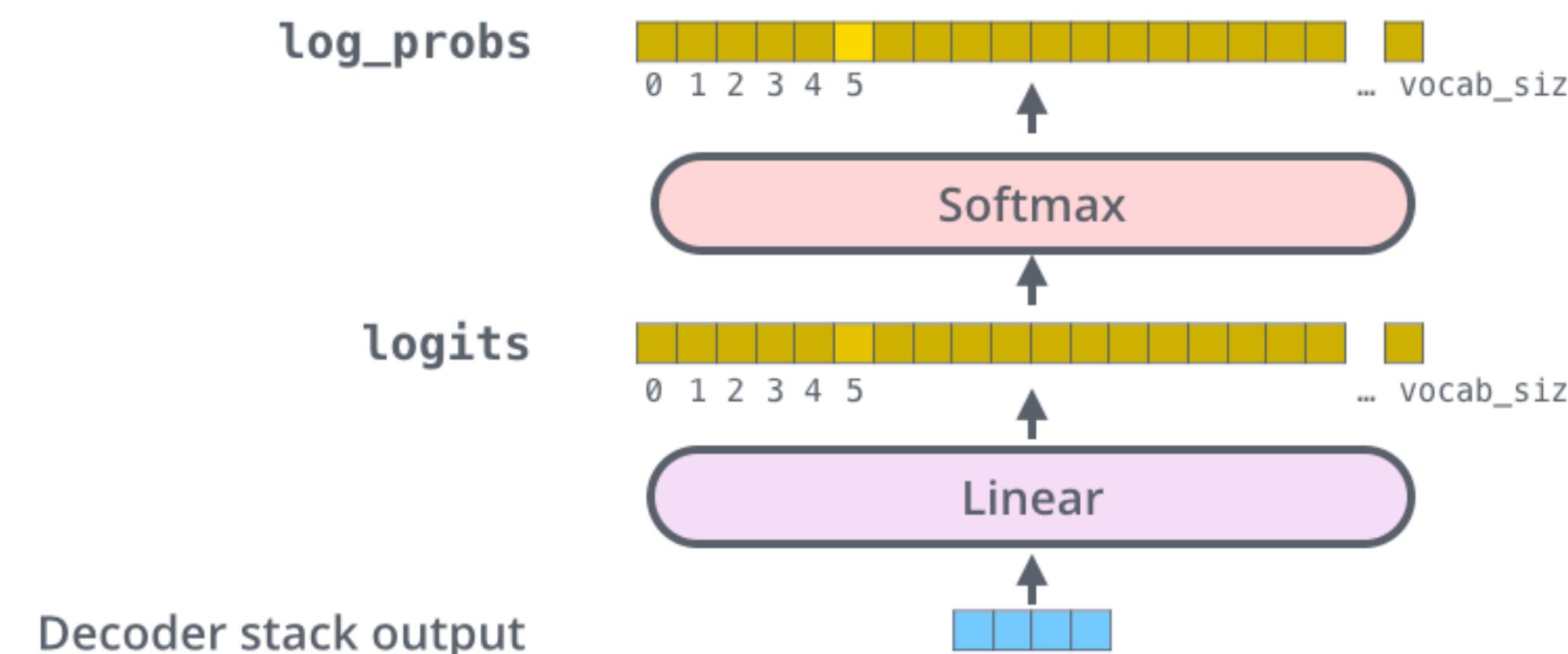
The Final Linear and Softmax Layer

Which word in our vocabulary is associated with this index?

Get the index of the cell with the highest value (`argmax`)

am

5



Encoder-Decoder

Recap of Training

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.



Ej.: Traducción automática

Output vocabulary *El vocabulario de salida del modelo, se hace en la fase de preprocesamiento*

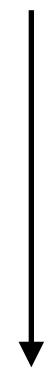
Output Vocabulary						
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

Encoder-Decoder

Recap of Training

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

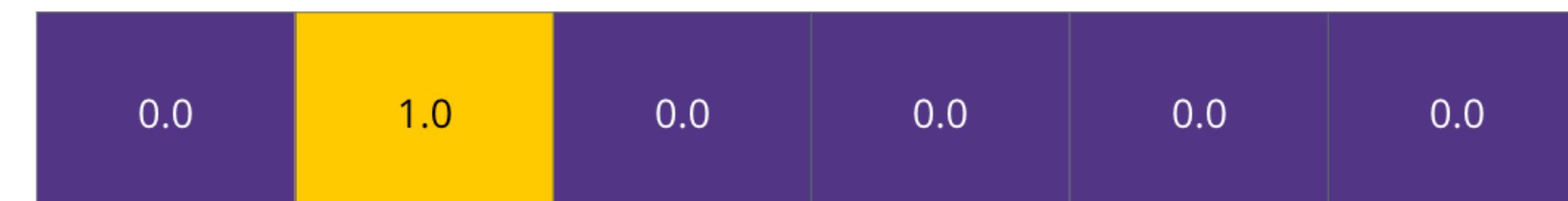


Ej.: Traducción automática

Output vocabulary *El vocabulario de salida del modelo, se hace en la fase de preprocesamiento*

Output Vocabulary						
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

One-hot encoding of the word "am"



Encoder-Decoder

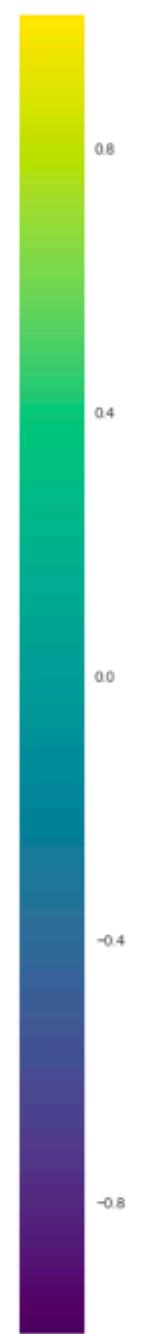
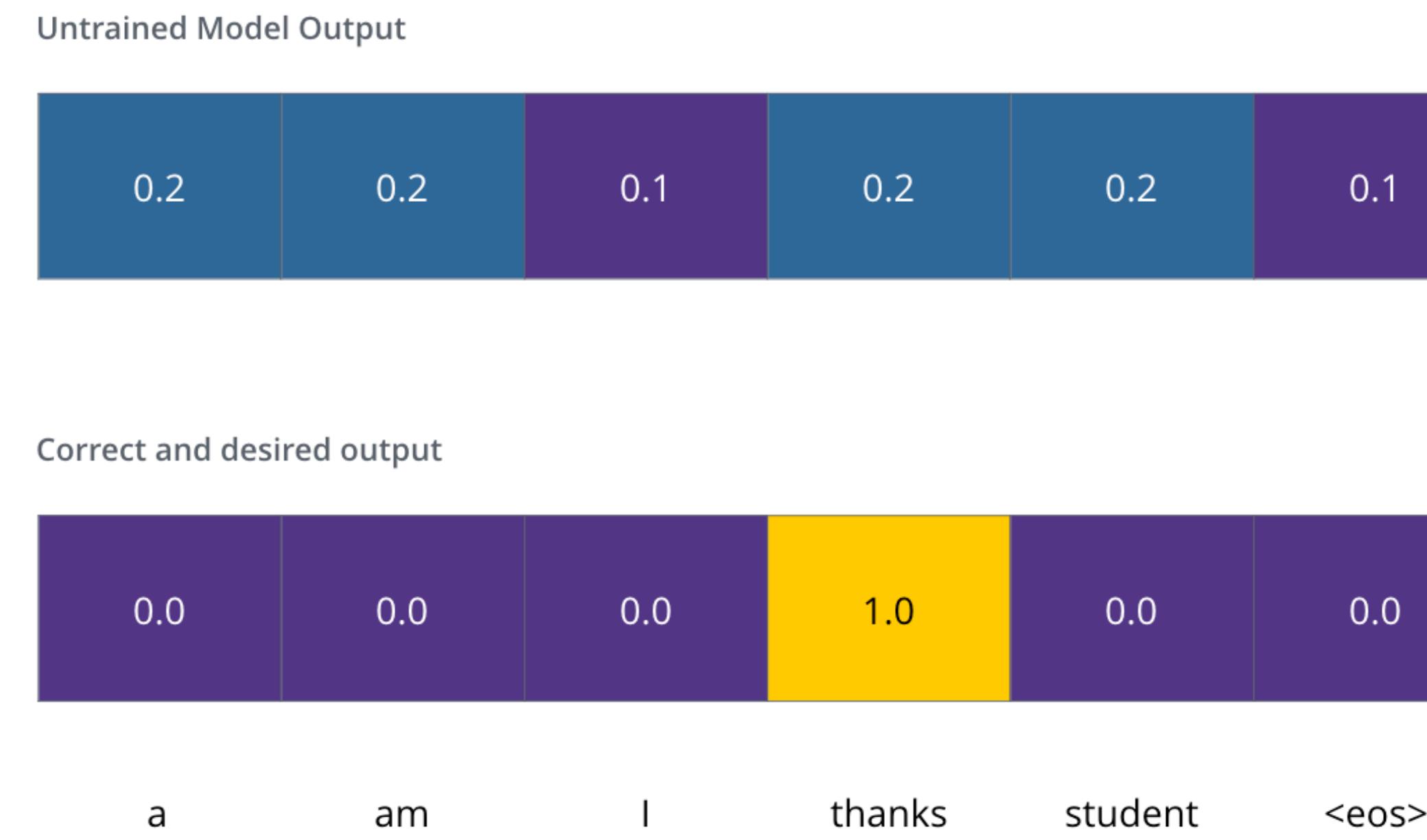
Recap of Training

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Loss function



Encoder-Decoder

Recap of Training

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

Loss function

Target Model Outputs

Output Vocabulary: a am I thanks student <eos>

position #1	0.0	0.0	1.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #2	0.0	1.0	0.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #3	1.0	0.0	0.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #4	0.0	0.0	0.0	0.0	1.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #5	0.0	0.0	0.0	0.0	0.0	1.0
-------------	-----	-----	-----	-----	-----	-----

a am I thanks student <eos>



Encoder-Decoder

Recap of Training

Transformers

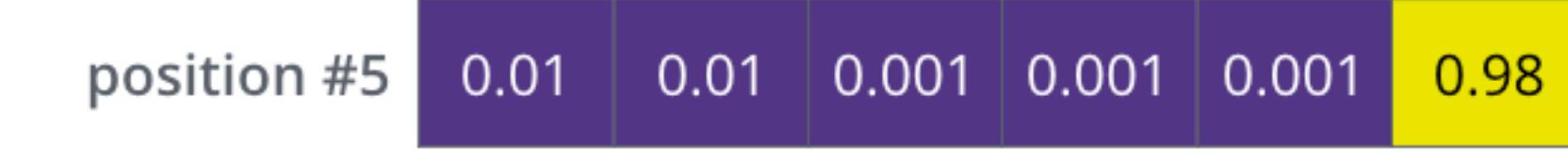
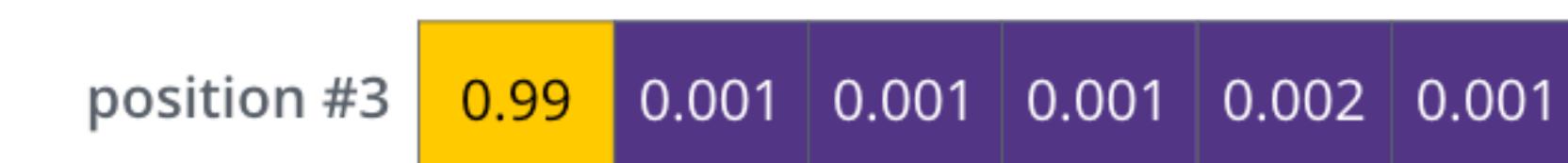
Atención para aumentar la velocidad en entrenamiento de modelos.

Ej.: Traducción automática

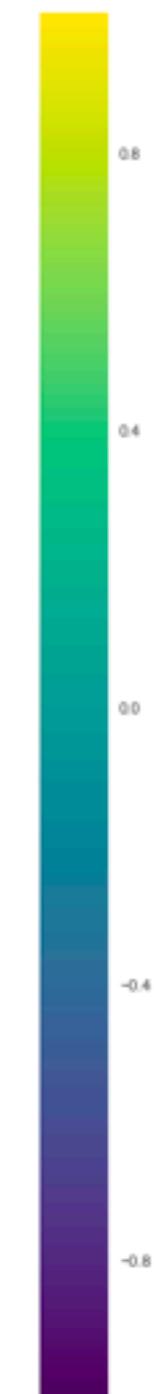
Loss function

Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>



a am I thanks student <eos>



Encoder-Decoder

Recap of Training

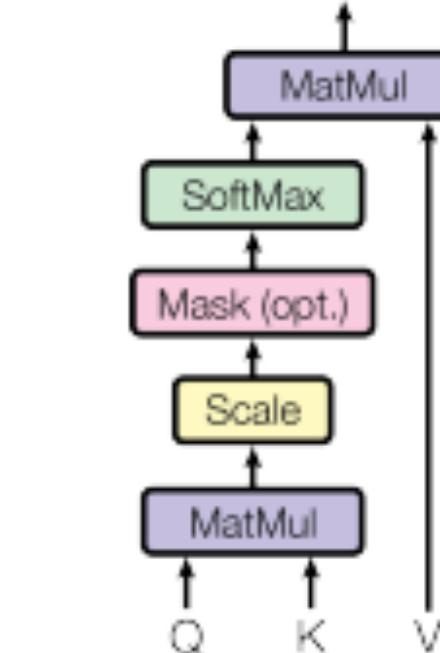
Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.

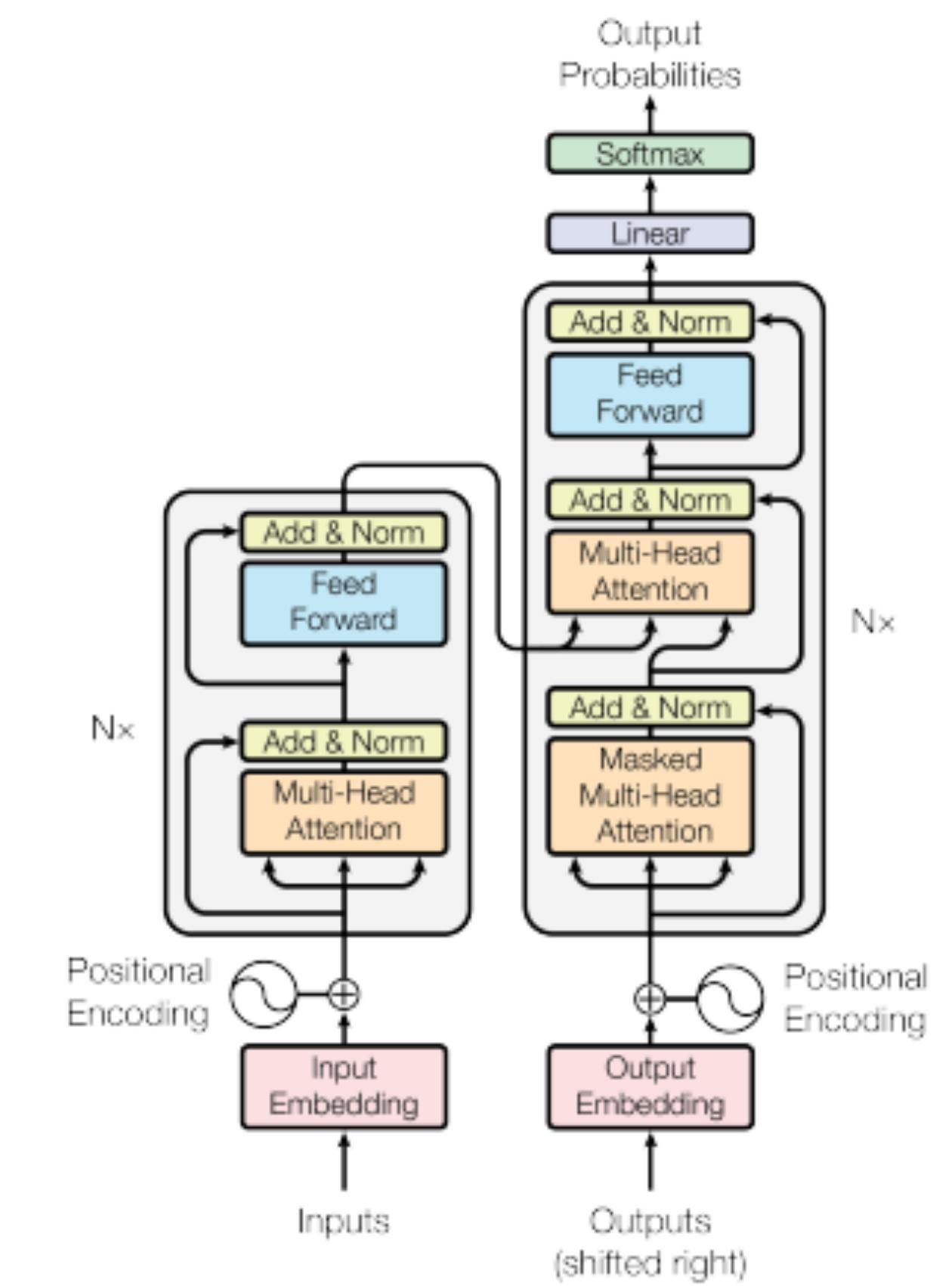
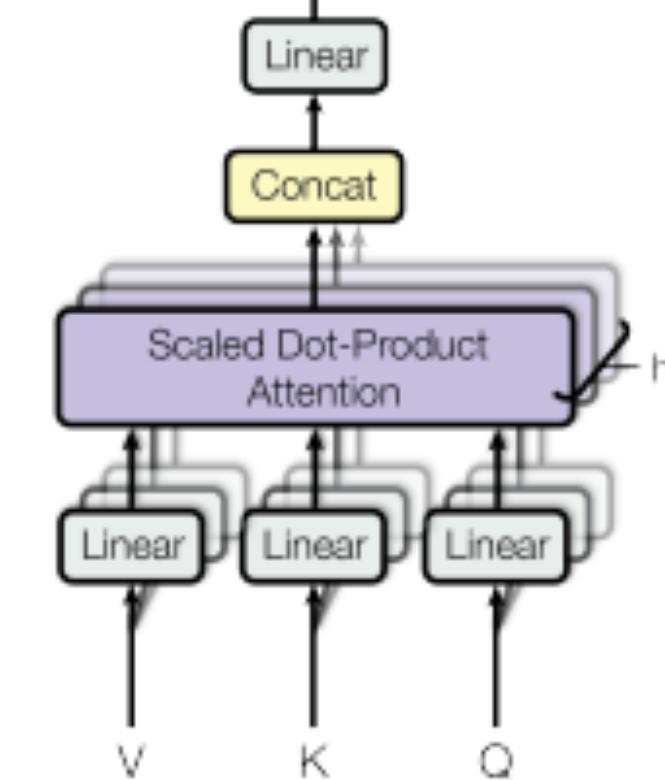
Ej.: Traducción automática

The Transformer - model architecture.

Scaled Dot-Product Attention



Multi-Head Attention

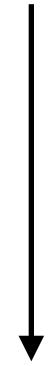


Encoder-Decoder

Recap of Training

Transformers

Atención para aumentar la velocidad en entrenamiento de modelos.



Ej.: Traducción automática

Loss function

Kullback-Leibler Divergence Explained

Entropía Cruzada