

Procesamiento de Lenguaje Natural

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

Sergio Alberto Mora Pardo - H2 2023

BERT:

**Bidirectional Encoder Representations
from Transformers**

BERT

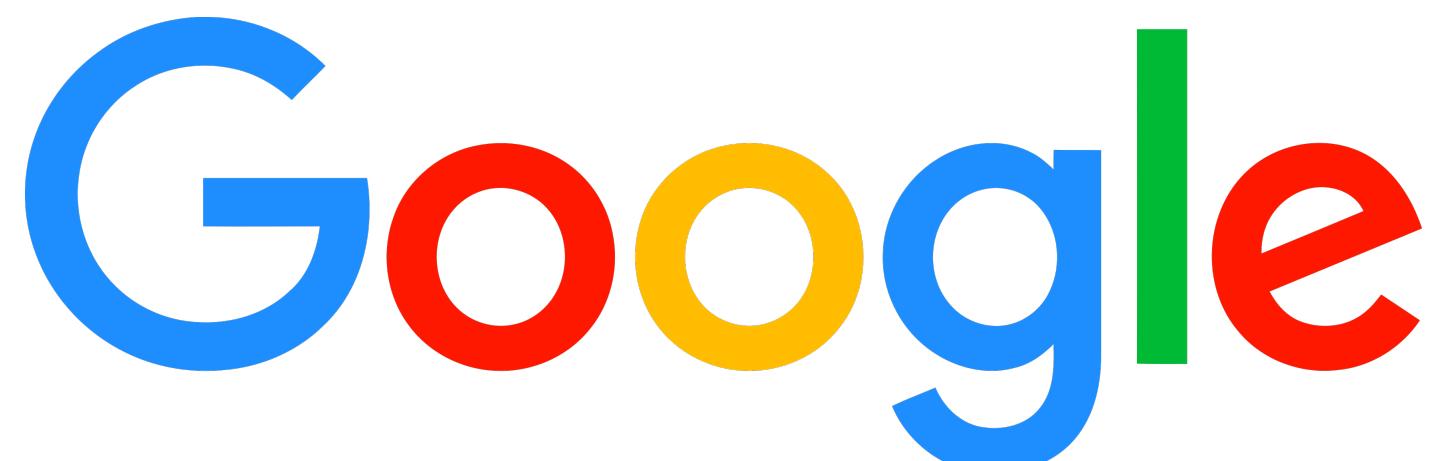
Bidirectional Encoder Representation for Transformers (2018)

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP



- La predicción del modelo Masked-language (Mask LM). (15%)
- La predicción de la siguiente oración (NSP). (50%)

BERT

Bidirectional Encoder Representation for Transformers

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Masked-Language

El modelo BERT toma como entrada dos oraciones separadas por un token especial ([SEP])

Input: the man went to the [MASK1] . he bought a [MASK2] of milk.

Labels:

- [MASK1] = store
- [MASK2] = gallon

BERT

Bidirectional Encoder Representation for Transformers

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Next Sentence Prediction

El modelo BERT toma como entrada dos oraciones separadas por un token especial ([SEP])

Sentence A: the man went to the store .
Sentence B: he bought a gallon of milk .

Label: **IsNextSentence**

Sentence A: the man went to the store .
Sentence B: penguins are flightless.

Label: **NotNextSentence**

BERT

Bidirectional Encoder Representation for Transformers

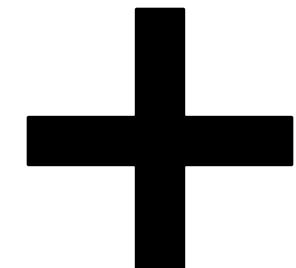
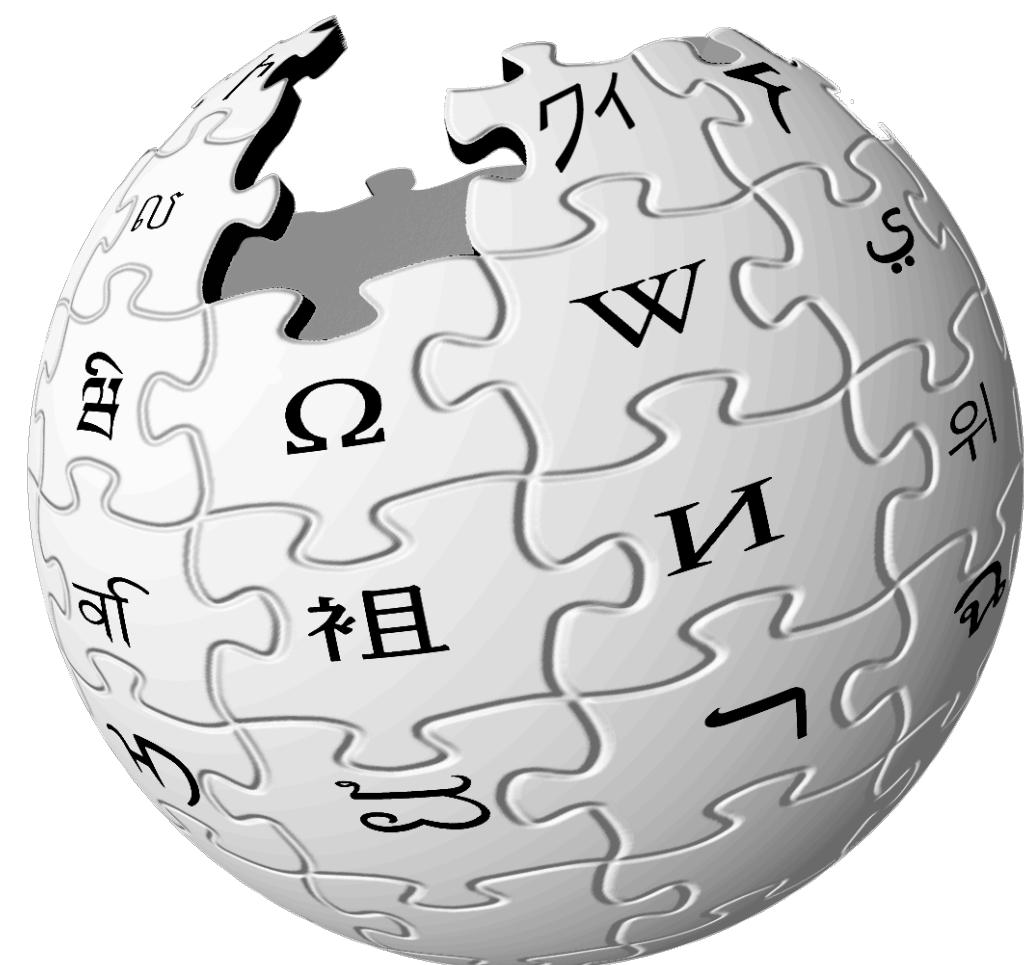
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Wikipedia



BookCorpus

**Aligning Books and Movies: Towards Story-like Visual Explanations by
Watching Movies and Reading Books**

Yukun Zhu ^{*,1} Ryan Kiros^{*,1} Richard Zemel¹ Ruslan Salakhutdinov¹
Raquel Urtasun¹ Antonio Torralba² Sanja Fidler¹
¹University of Toronto ²Massachusetts Institute of Technology
 {yukun,rkiros,zemel,rsalaku,urtasun,fidler}@cs.toronto.edu, torralba@csail.mit.edu

Abstract

Books are a rich source of both fine-grained information, how a character, an object or a scene looks like, as well as high-level semantics, what someone is thinking, feeling and how these states evolve through a story. This paper aims to align books to their movie releases in order to provide rich descriptive explanations for visual content that go semantically far beyond the captions available in current datasets. To align movies and books we exploit a neural sentence embedding that is trained in an unsupervised way from a large corpus of books, as well as a video-text neural embedding for computing similarities between movie clips and sentences in the book. We propose a context-aware CNN to combine information from multiple sources. We demonstrate good quantitative performance for movie/book alignment and show several qualitative examples that showcase the diversity of tasks our model can be used for.

1. Introduction

A truly intelligent machine needs to not only parse the surrounding 3D environment, but also understand whv ne-

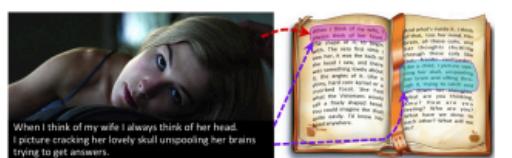


Figure 1: Shot from the movie *Gone Girl*, along with the book. We reason about the visual and dialog (text) alignment between the movie and a book.

Books provide us with very rich, descriptive text that conveys both fine-grained visual details (how people or scenes look like) as well as high-level semantics (what people think and feel, and how their states evolve through a story). This source of knowledge, however, does not come with associated visual information that would enable us to ground it with descriptions. Grounding descriptions in books to vision would allow us to get textual explanations or stories behind visual information rather than simplistic captions available in current datasets. It can also provide us with extremely large amount of data (with tens of thousands books available online).

**Alinear libros y películas: hacia
explicaciones visuales parecidas a
historias viendo películas y leyendo
libros**

BERT

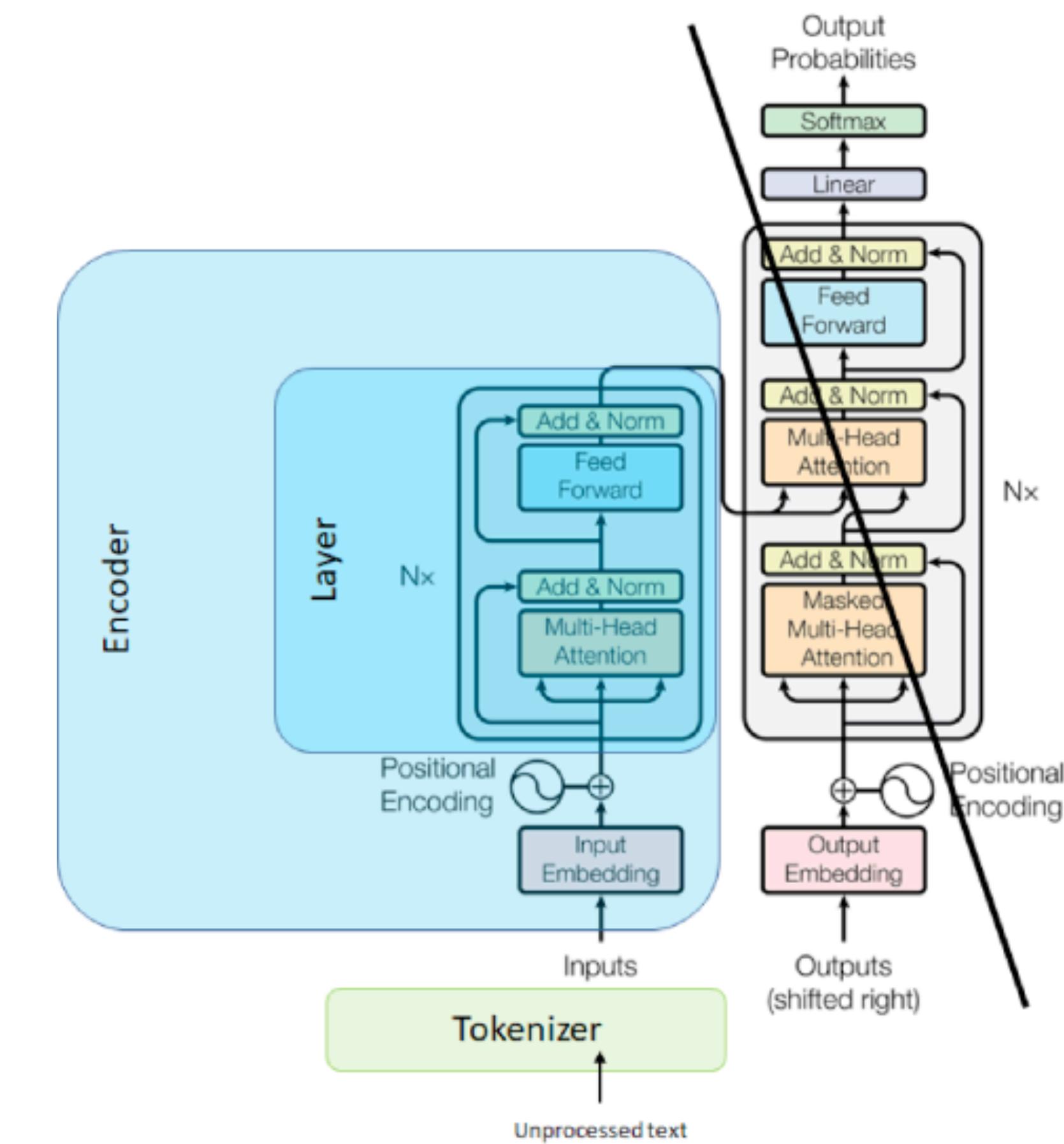
Bidirectional Encoder Representation for Transformers

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP



BERT

Bidirectional Encoder Representation for Transformers

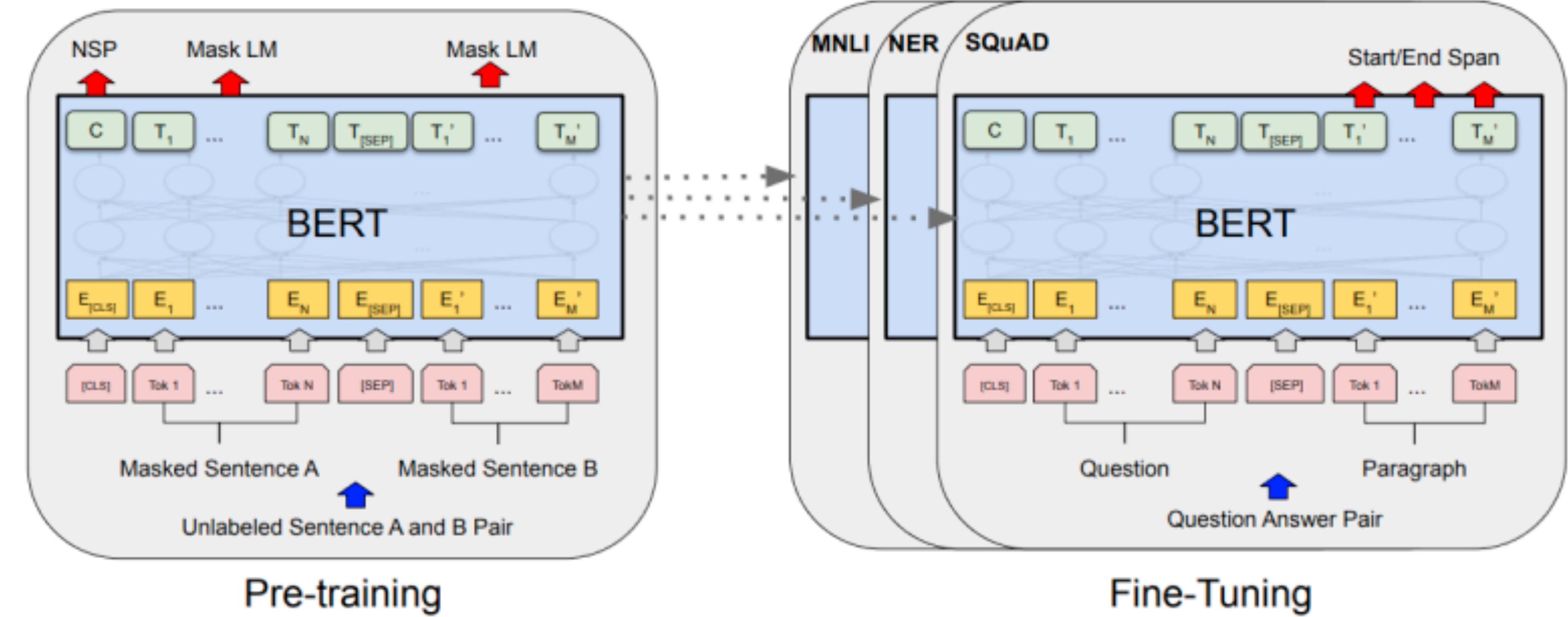
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

BERT es utilizado a menudo como un modelo *encoder* de lenguaje.



Puede entenderse como un punto de chequeo (checkpoint) pre-entrenado el cual es extendido con capas adicionales que son específicas a la tarea de interés

BERT

Bidirectional Encoder Representation for Transformers

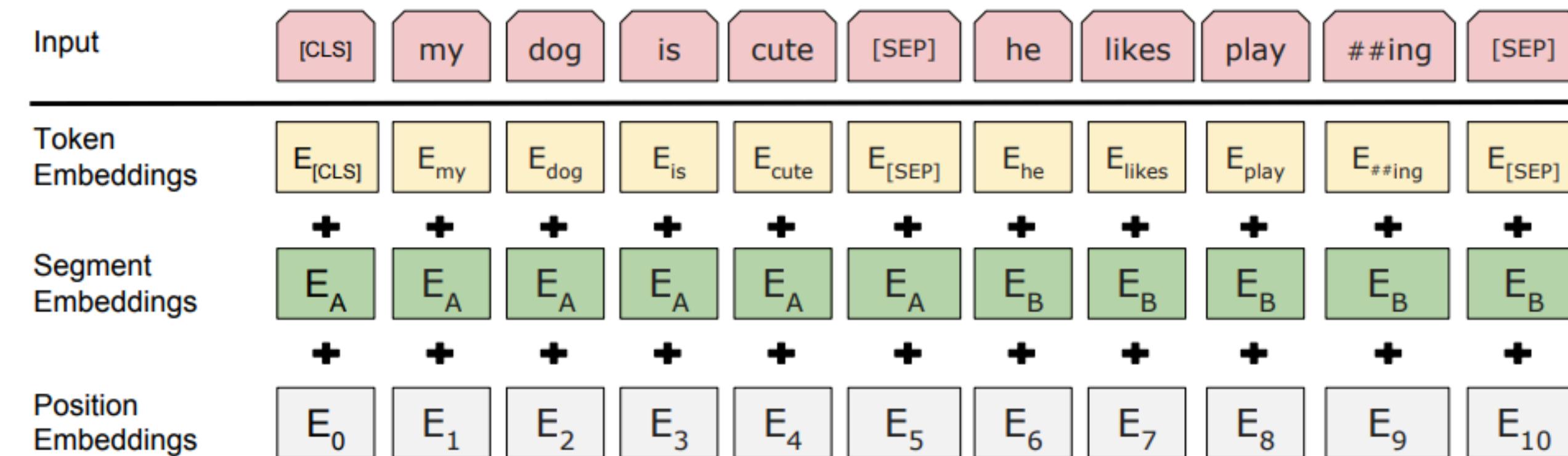
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Input representation



Usa los tokens [SEP] y [CLS] para cada oración y [MASK] para pronosticar esa misma palabra.

BERT

Bidirectional Encoder Representation for Transformers

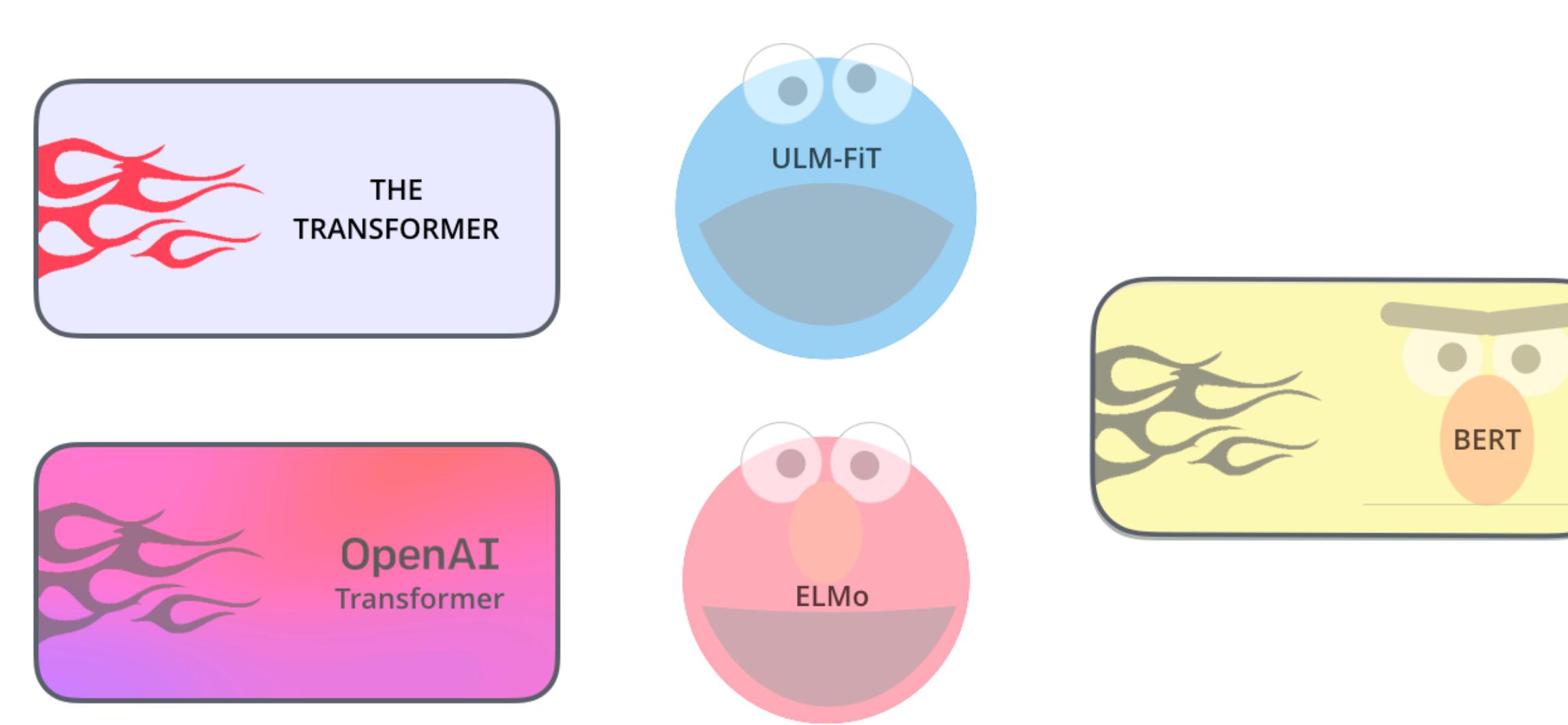
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

BERT family... (RoBERT, RoBERTA, alBERT, Hasta un tal BETO existe)



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

BERT

Bidirectional Encoder Representation for Transformers

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

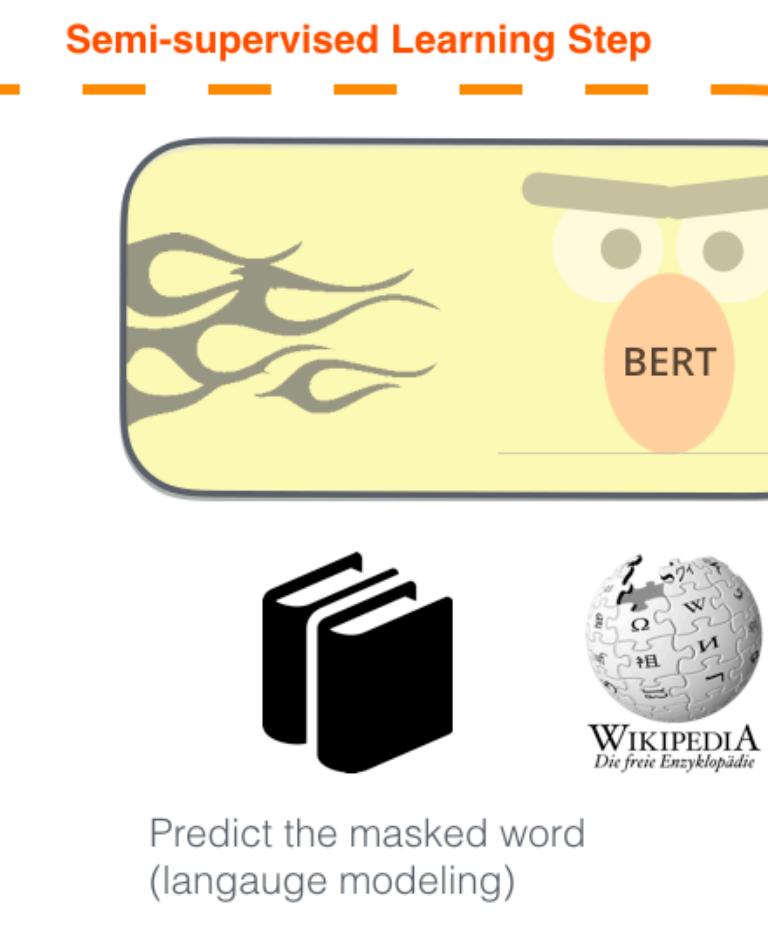
Función de perdida:

- Mask ML
- NSP

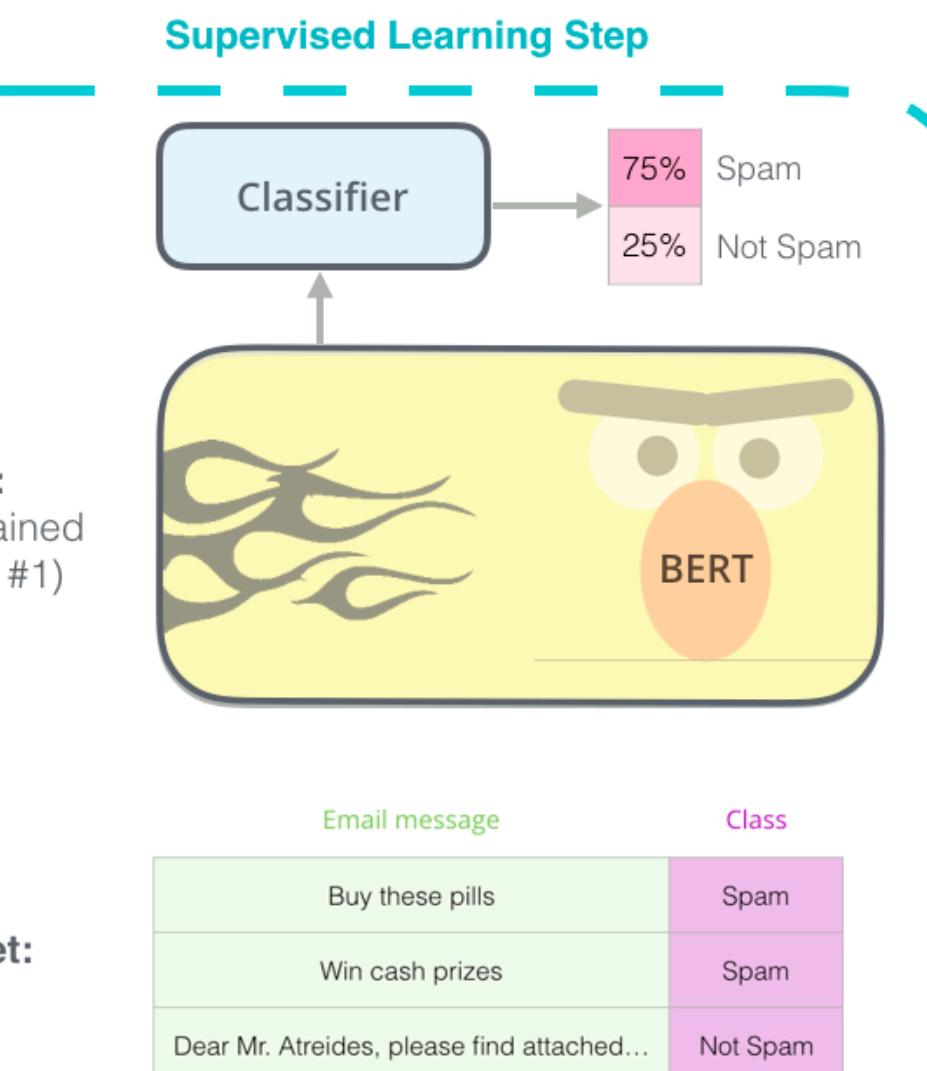
Steps:

- 1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



- 2 - **Supervised** training on a specific task with a labeled dataset.



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

BERT

Bidirectional Encoder Representation for Transformers

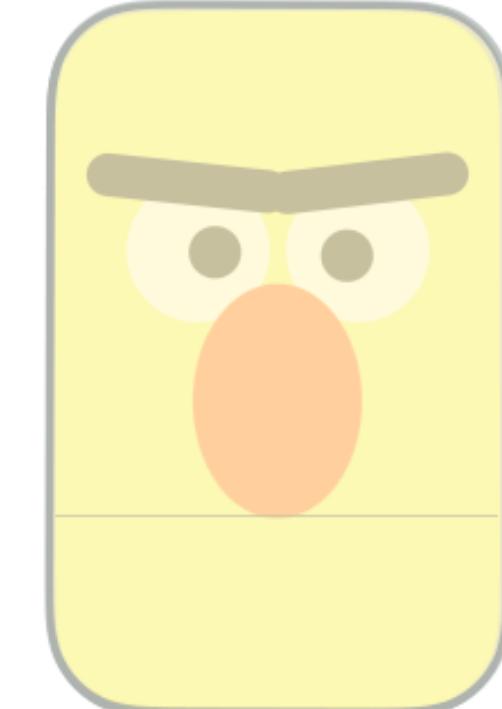
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

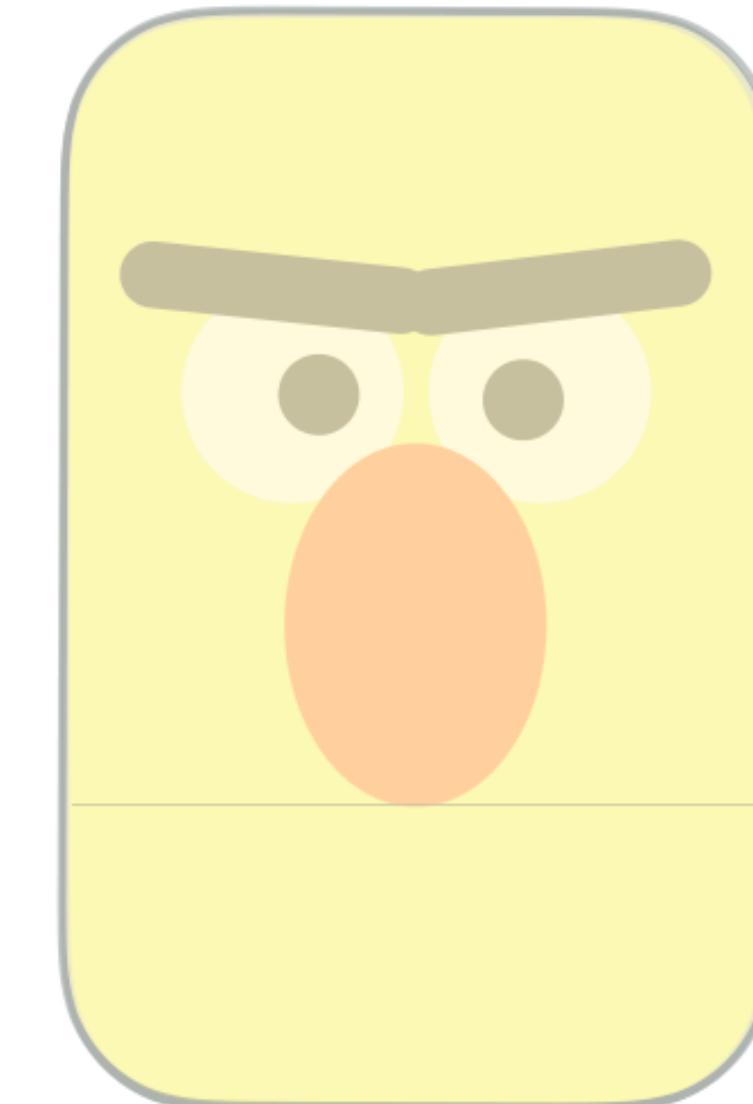
Función de perdida:

- Mask ML
- NSP

Model arquitectura:



BERT_{BASE}



BERT_{LARGE}

BERT-base de tamaño similar al OpenAI Transformer

BERT

Bidirectional Encoder Representation for Transformers

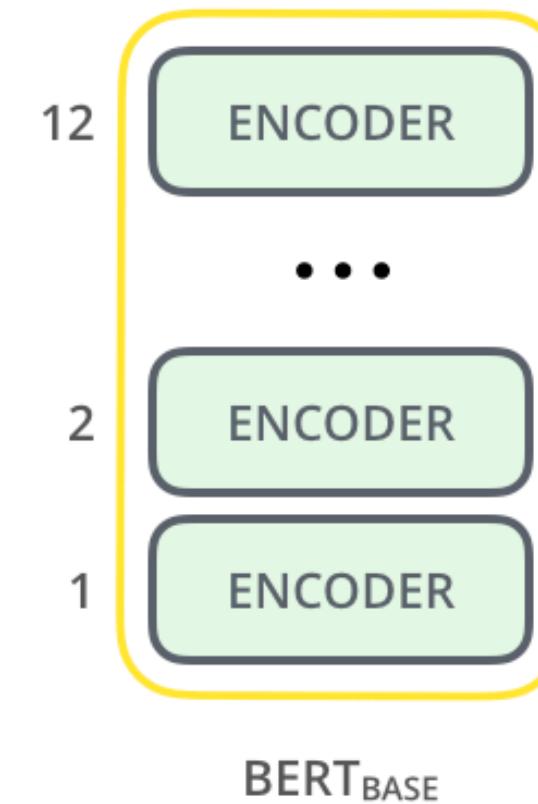
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

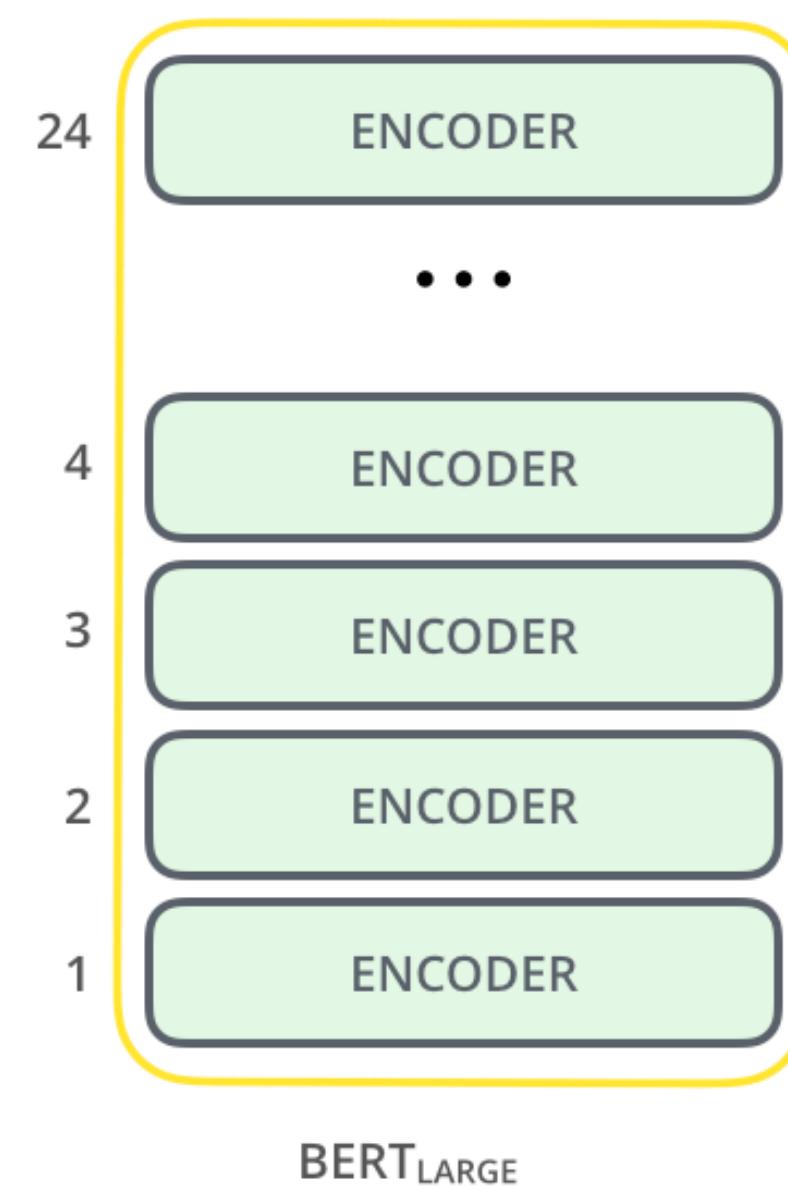
- Mask ML
- NSP

Model arquitectura:



BERT_{BASE}

768 dimensiones
12 cabezas de atención



BERT_{LARGE}

1024 unidades ocultas (dimensiones)
16 cabezas de atención

Transformer:
6 capas de encoder
512 dimensiones
8 cabezas de atención

BERT

Bidirectional Encoder Representation for Transformers

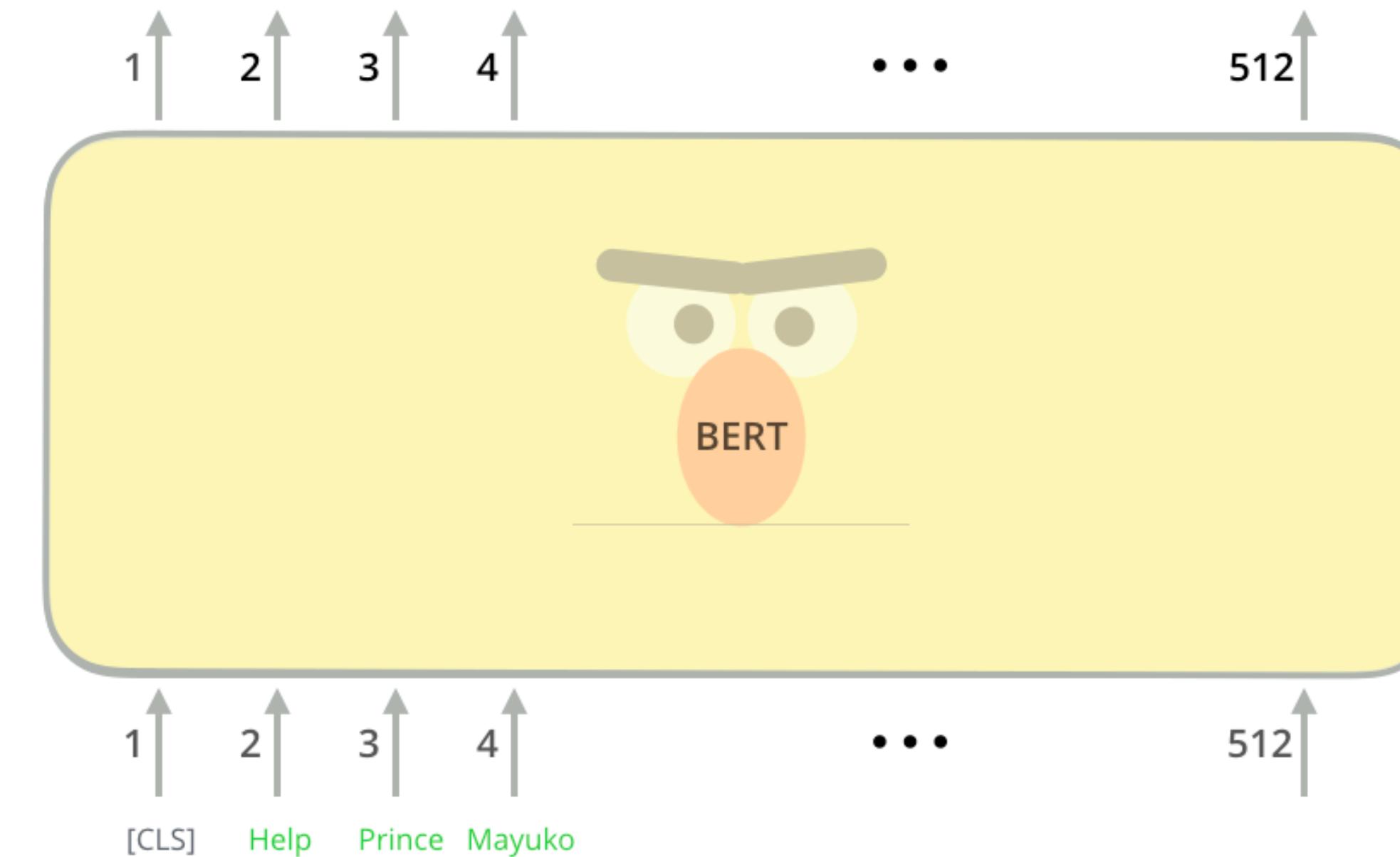
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Entradas del modelo



El primer token es [CLS] significa clasificación

BERT

Bidirectional Encoder Representation for Transformers

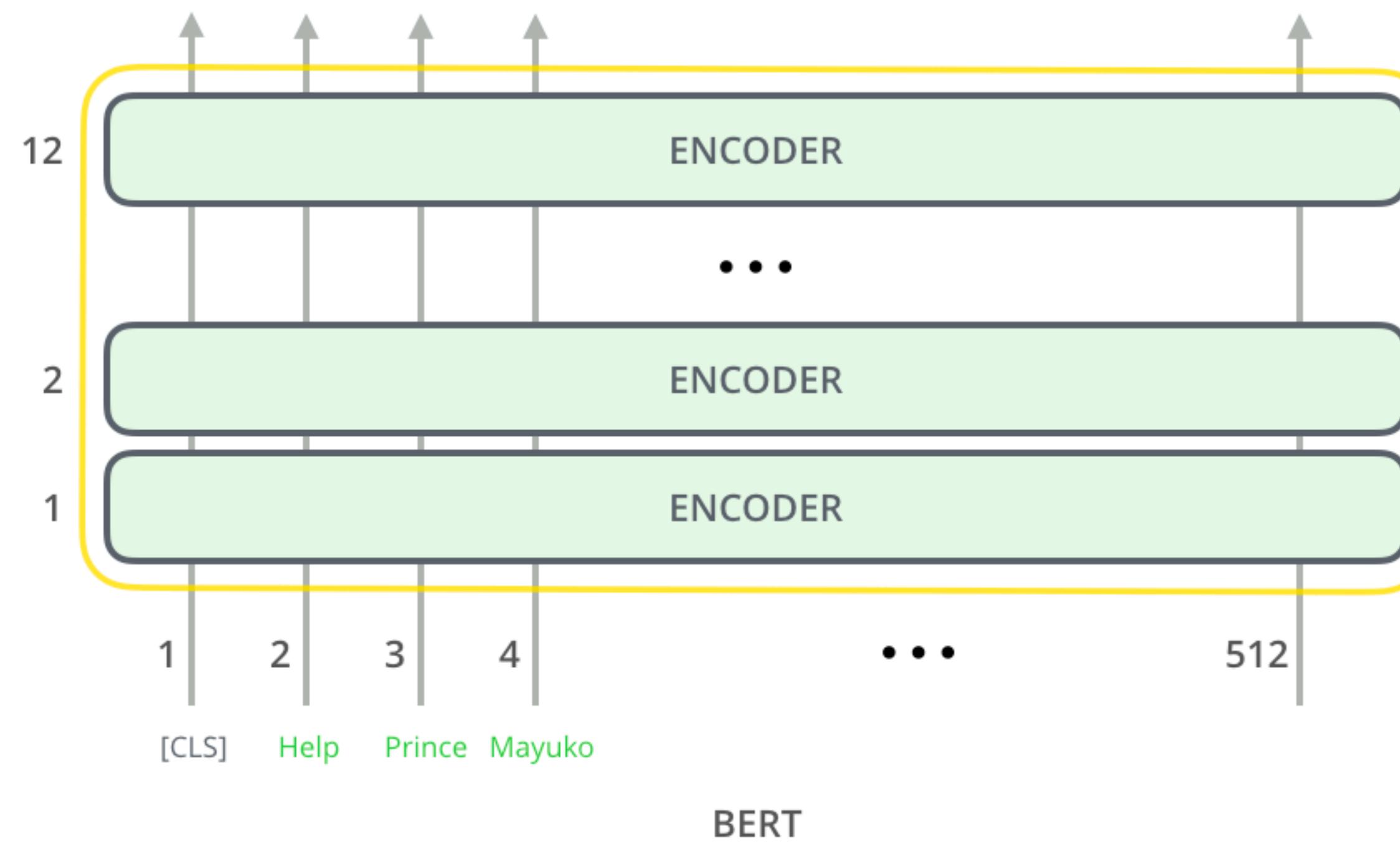
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Entradas del modelo



La arquitectura es idéntica al Transformer (hasta ahora)

BERT

Bidirectional Encoder Representation for Transformers

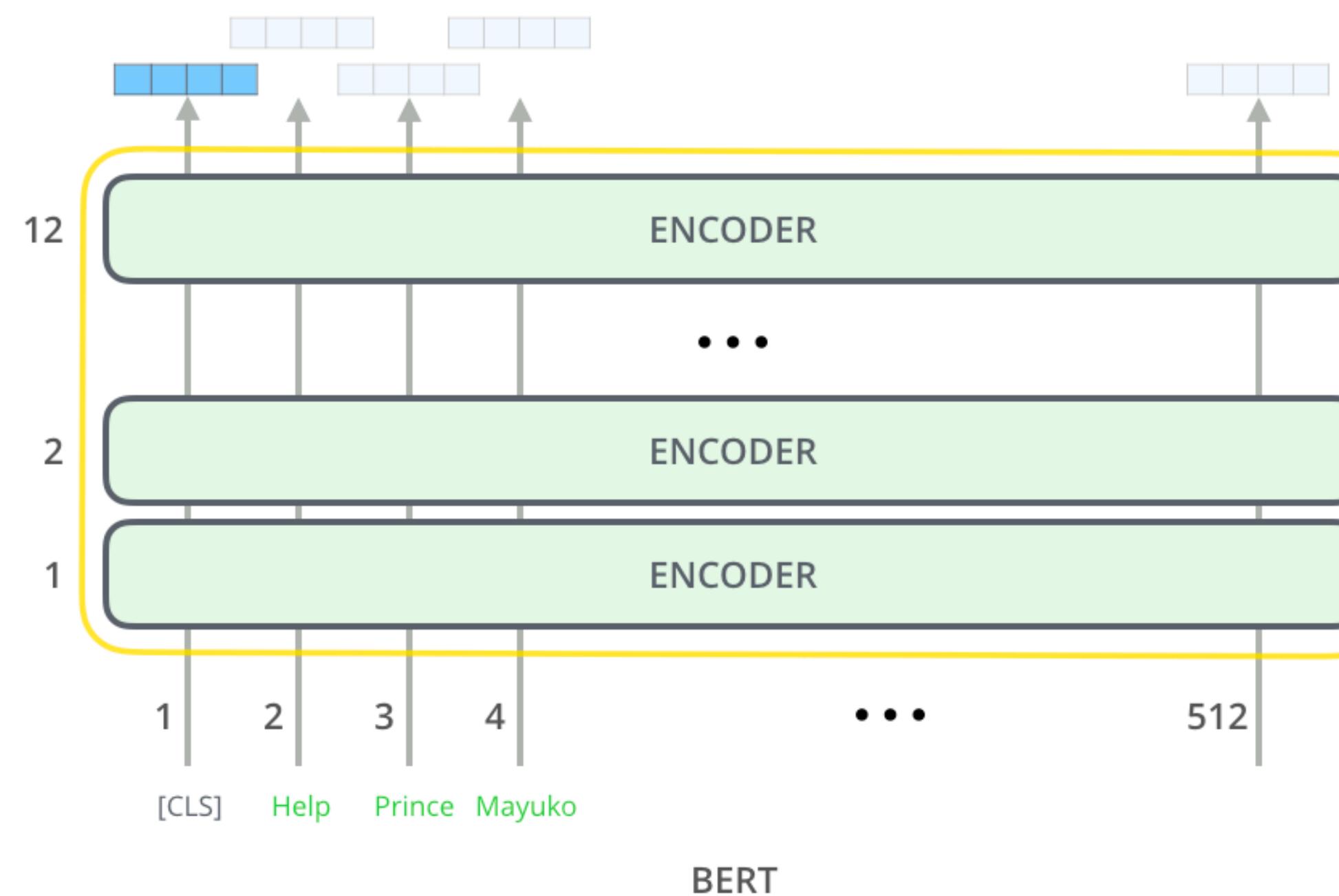
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Resultados del modelo



Cada posición genera
un vector de tamaño
768 (BERT base)

Nos centramos en el
vector con el token
especial [CLS]

BERT

Bidirectional Encoder Representation for Transformers

BERT

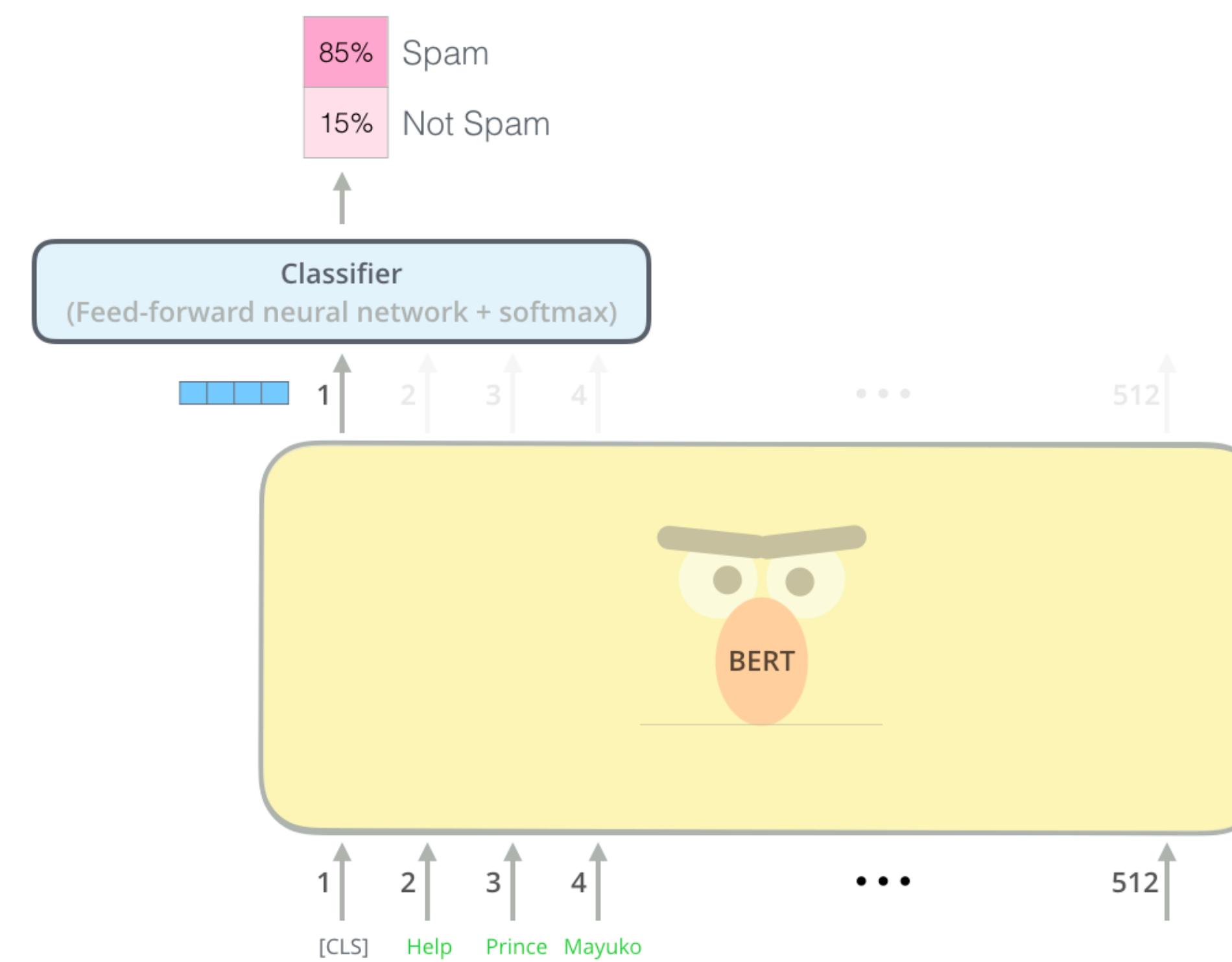
*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Resultados del modelo

85% Spam
15% Not Spam



modelo de calificación
Ej.: spam-social-promoción
Ej.: spam

vector como entrada

BERT

Bidirectional Encoder Representation for Transformers

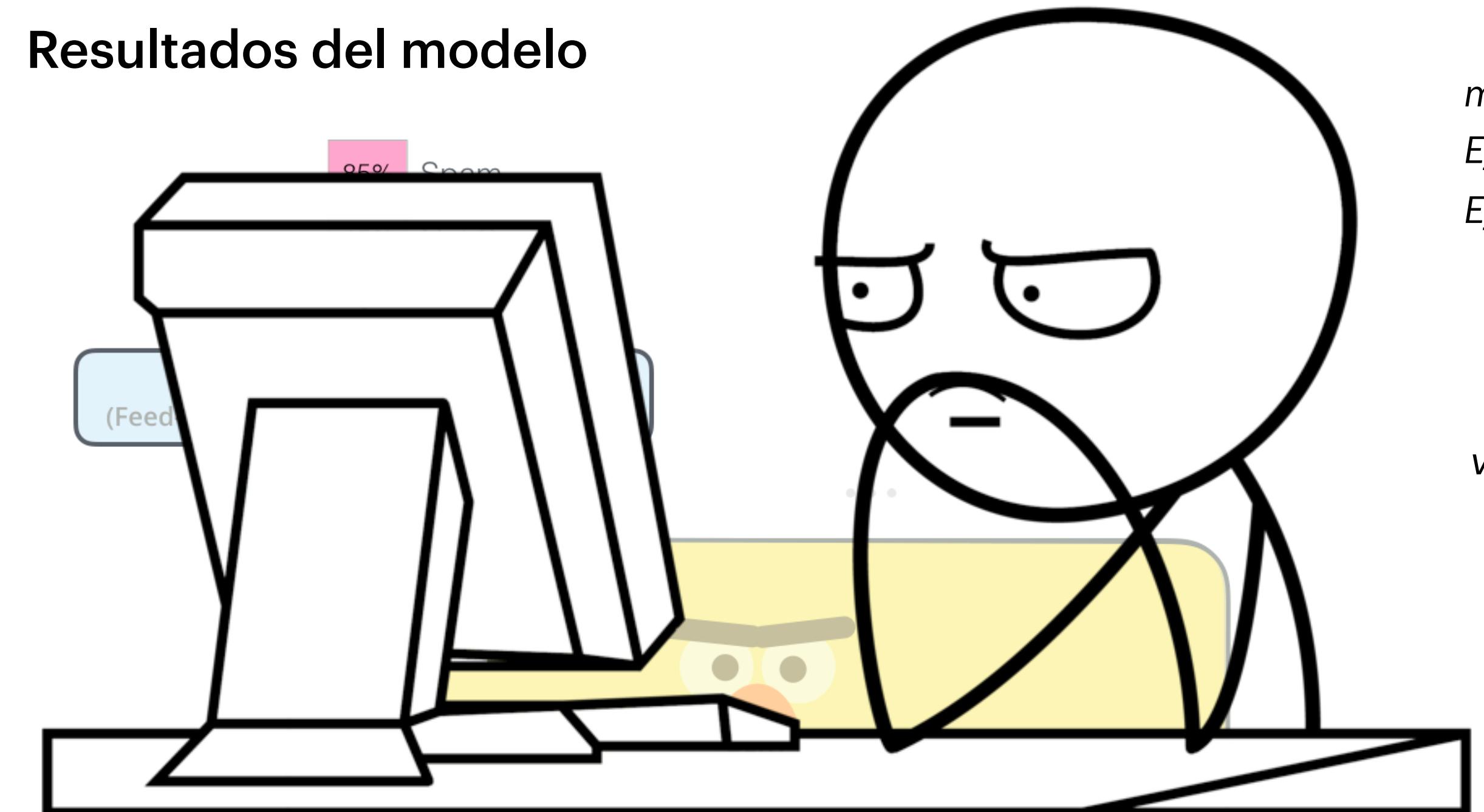
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

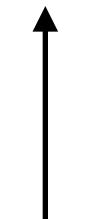
Función de perdida:

- Mask ML
- NSP

Resultados del modelo



modelo de calificación
Ej.: spam-social-promoción
Ej.: spam



vector como entrada

¿Entonces que cambio?

BERT

Bidirectional Encoder Representation for Transformers

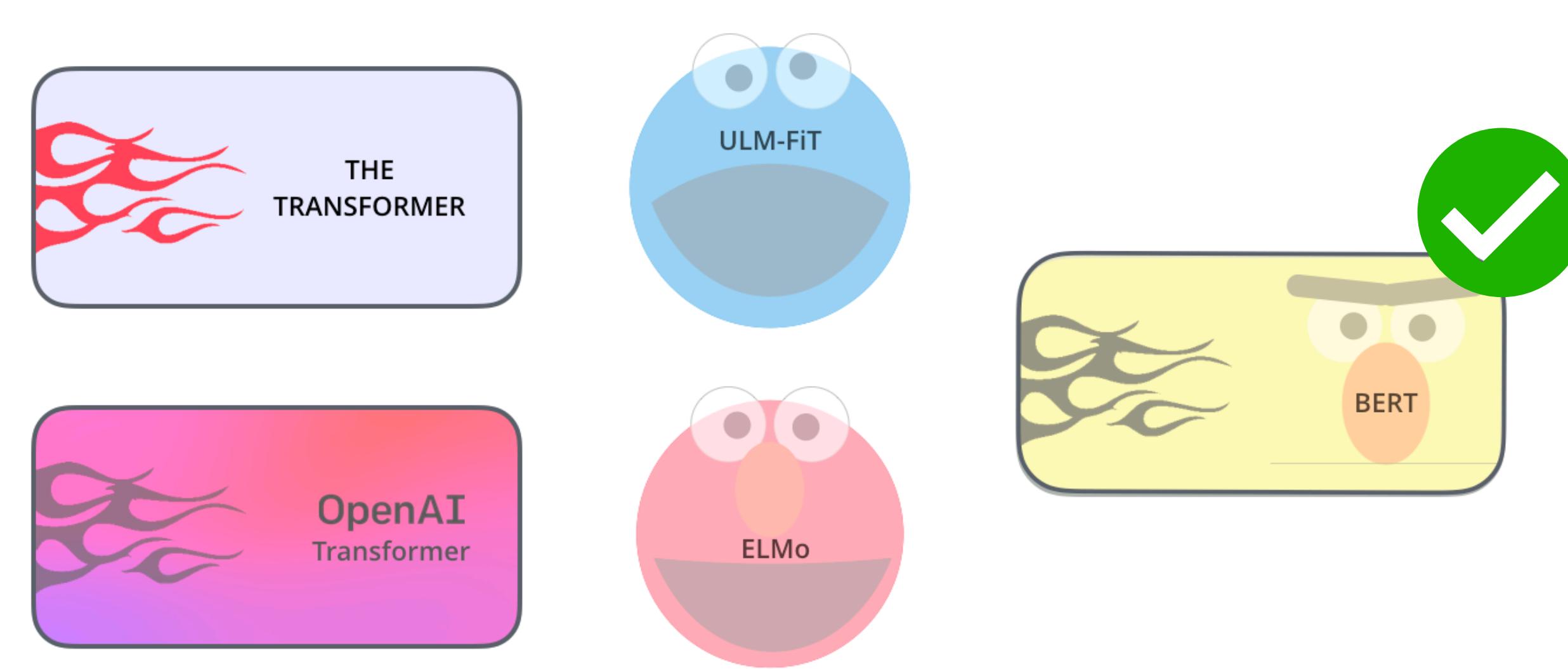
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

BERT family... (RoBERT, RoBERTA, alBERT, Hasta un tal BETO existe)



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

ELMo

Embeddings from Language Models

ELMo

*Embeddings from
Language Models*

LSTM bidireccional



Revisa la oración
completa

Context Matters



Representación semantica

Word2Vec

Glove



Medir distancia/similitud

Compresión en un
espacio vectorial

ELMo

Embeddings from Language Models

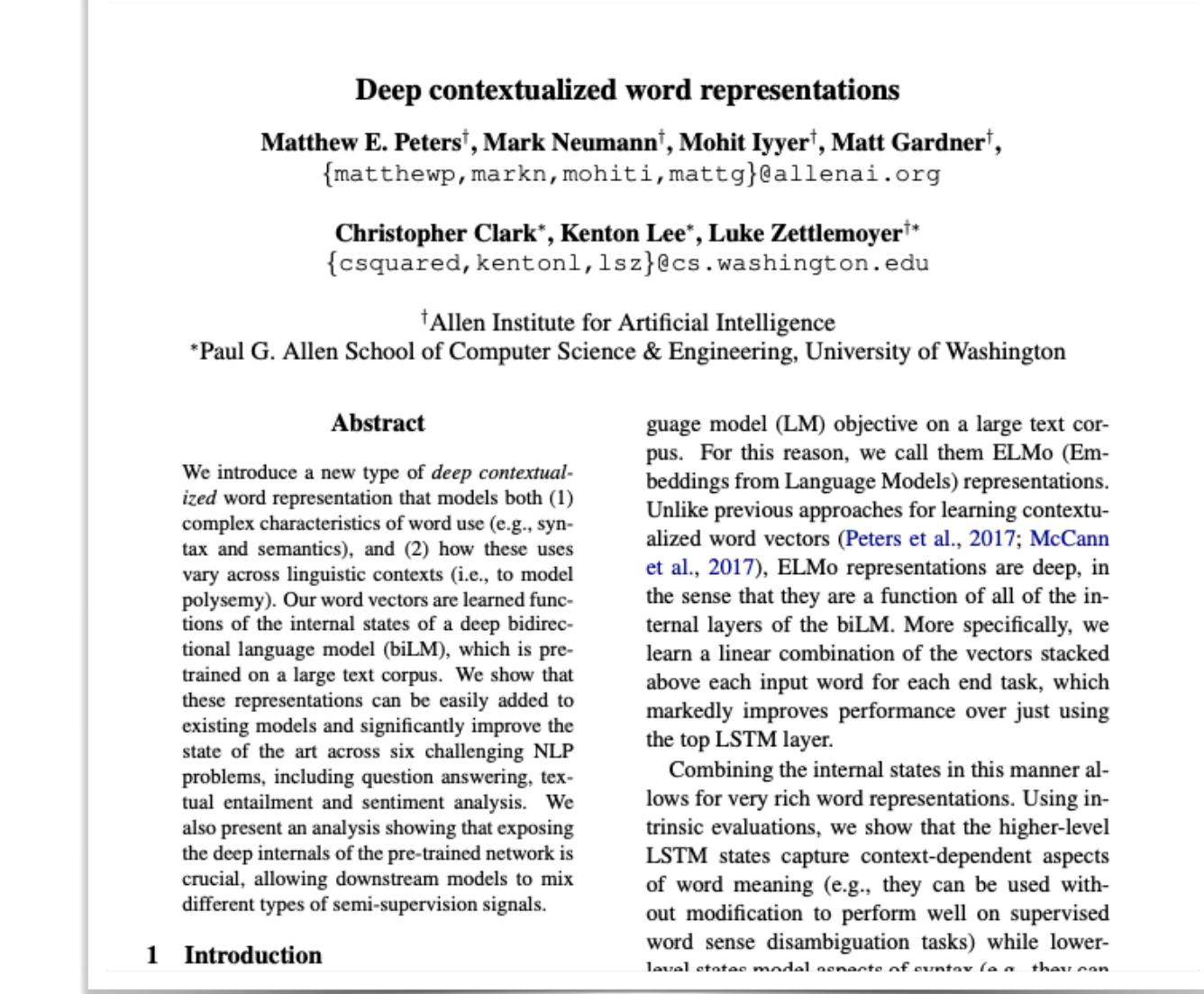
ELMo

*Embeddings from
Language Models*

LSTM bidireccional

↓
Revisa la oración
completa

Context Matters



Peters et. al., 2018 en el artículo de ELMo

Representación semántica

Word2Vec

Glove

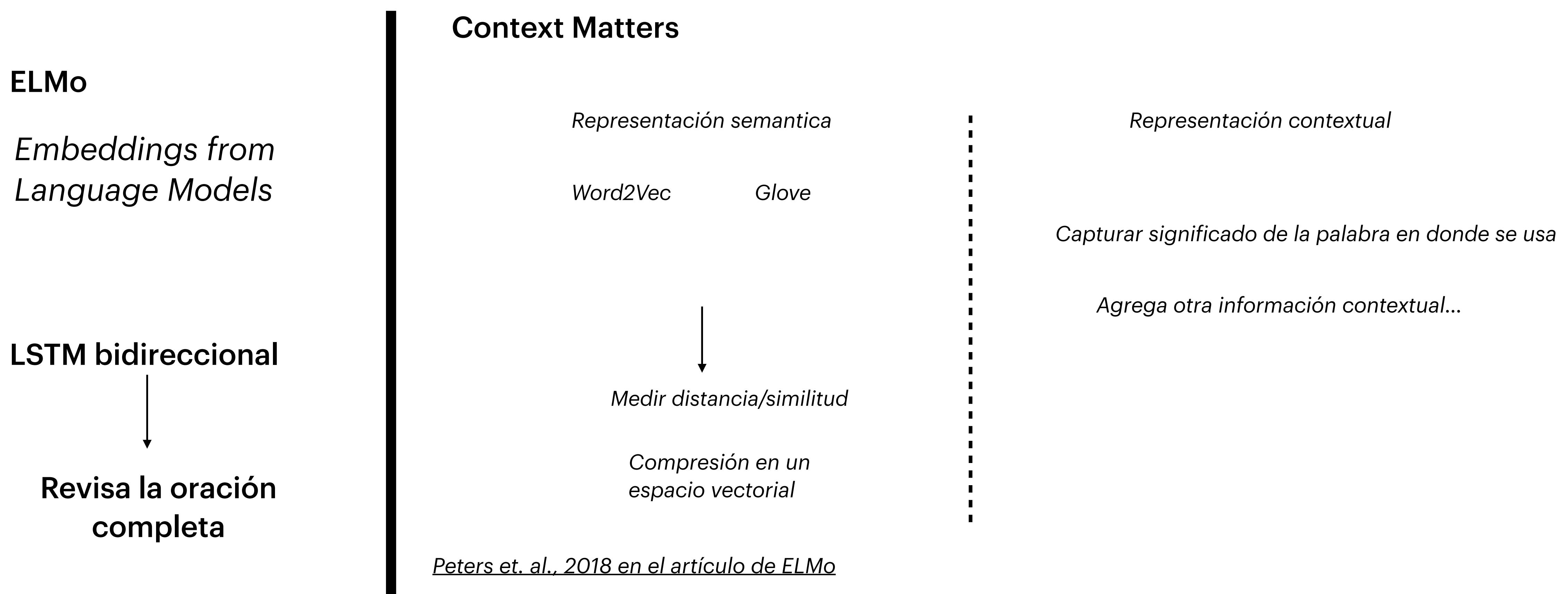


↓
Medir distancia/similitud

Compresión en un
espacio vectorial

ELMo

Embeddings from Language Models



ELMo

Embeddings from Language Models

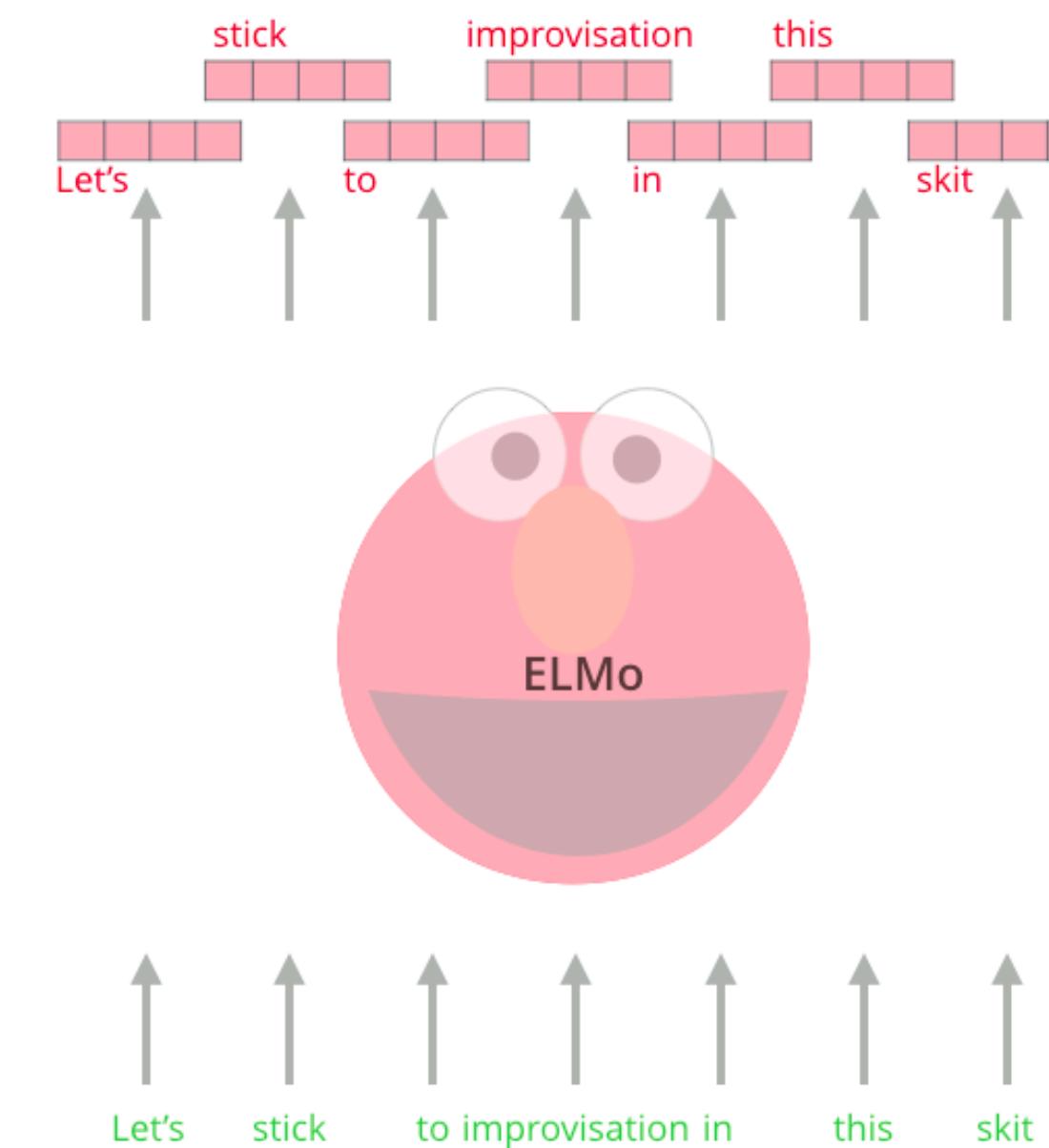
ELMo
*Embeddings from
Language Models*

LSTM bidireccional
↓
**Revisa la oración
completa**

Context Matters

ELMo
Embeddings

Words to embed



Representación contextual
Capturar significado de la
palabra en donde se usa

Peters et. al., 2018 en el artículo de ELMo

ELMo

Embeddings from Language Models

ELMo

*Embeddings from
Language Models*

LSTM bidireccional

Revisa la oración
completa

Context Matters

Modelo de
Lenguaje

Pronostica la
siguiente palabra

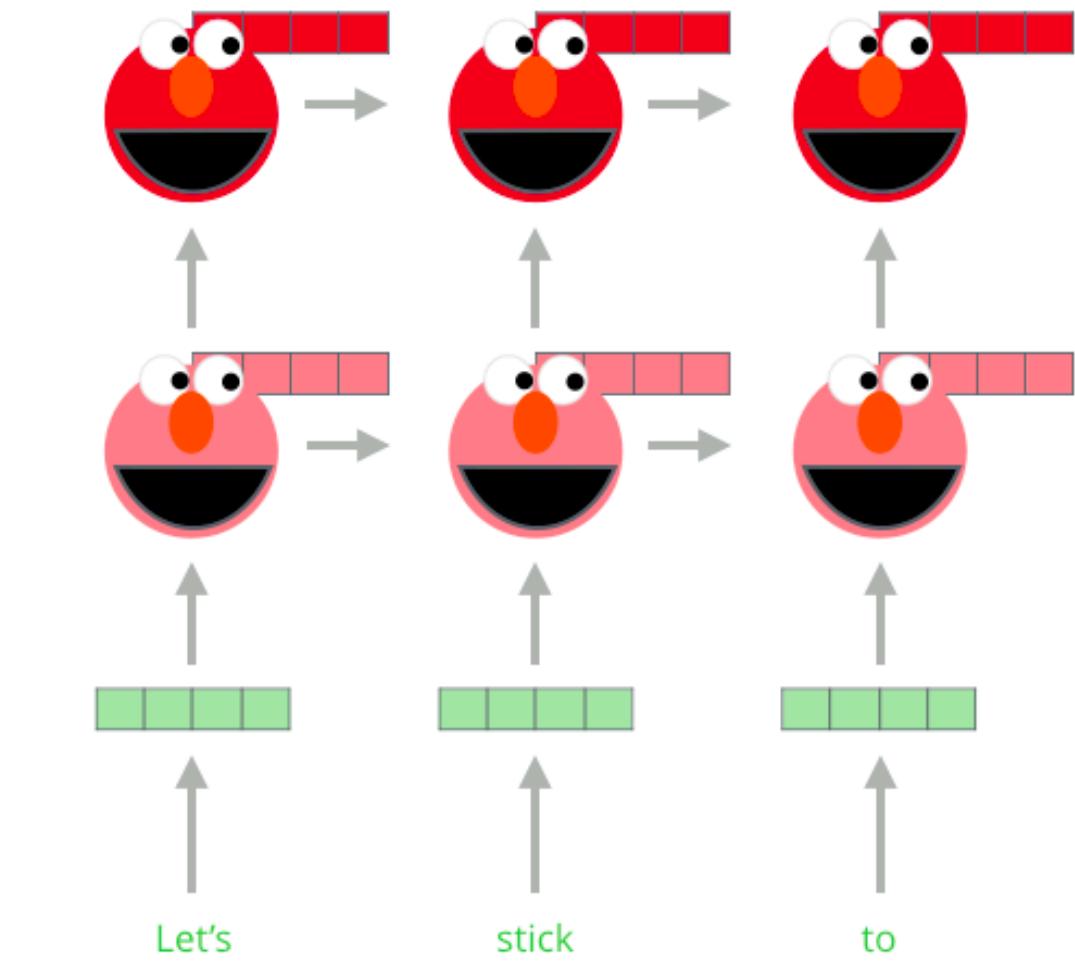
Possible classes:
All English words

Output
Layer

LSTM
Layer #2

LSTM
Layer #1

Embedding



Peters et. al., 2018 en el artículo de ELMo

ELMo

Embeddings from Language Models

ELMo

*Embeddings from
Language Models*

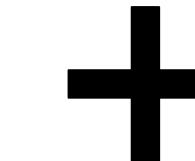
LSTM bidireccional

Revisa la oración
completa

Context Matters

Modelo de
Lenguaje

Pronostica la
siguiente palabra



Pronostica la palabra
anterior

Peters et. al., 2018 en el artículo de ELMo

Possible classes:
All English words

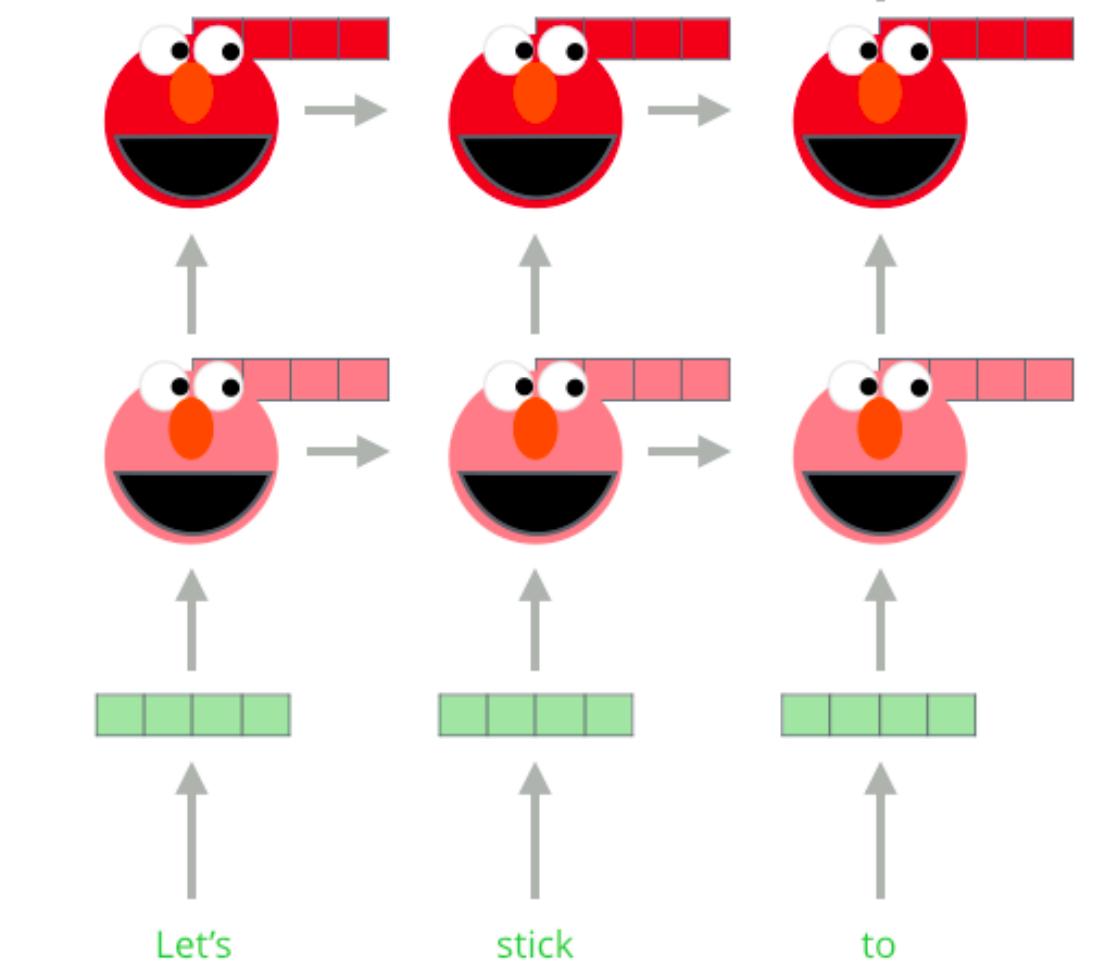


Output
Layer

LSTM
Layer #2

LSTM
Layer #1

Embedding



ELMo

Embeddings from Language Models

ELMo

*Embeddings from
Language Models*

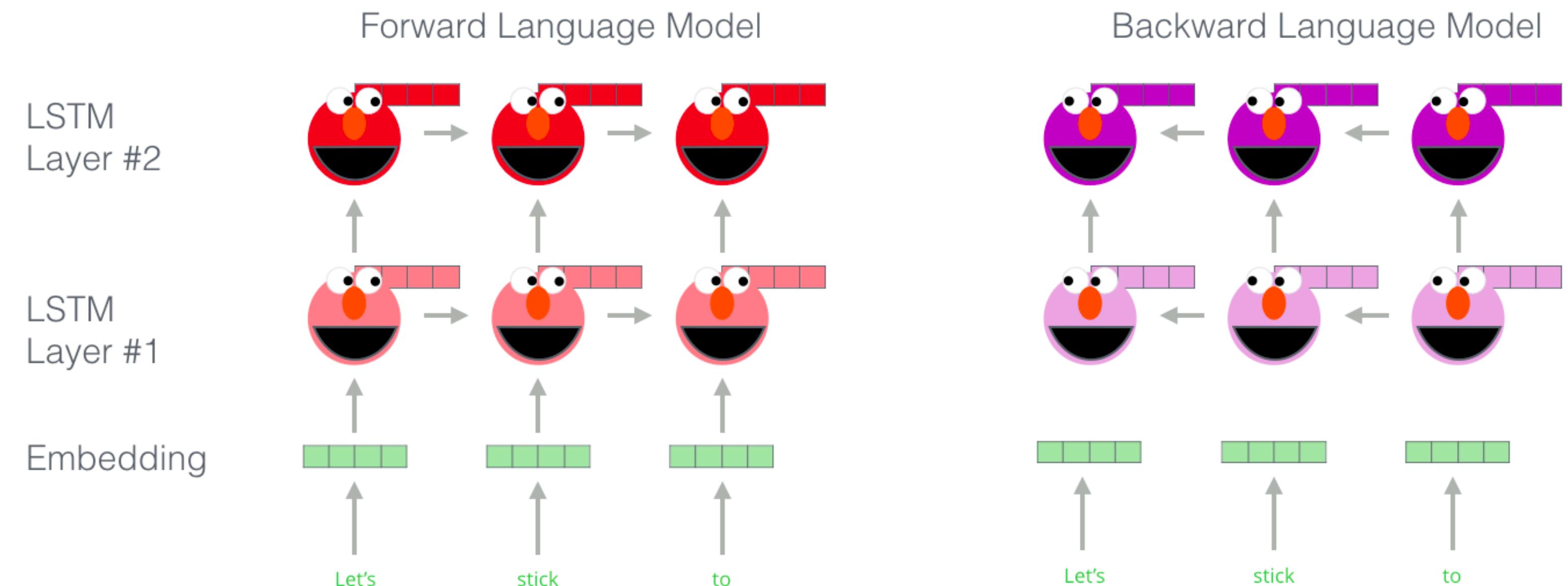
LSTM bidireccional



**Revisa la oración
completa**

Entrenamiento

Embedding of “stick” in “Let’s stick to” - Step #1



Peters et. al., 2018 en el artículo de ELMo

ELMo

Embeddings from Language Models

ELMo

*Embeddings from
Language Models*

LSTM bidireccional

↓
**Revisa la oración
completa**

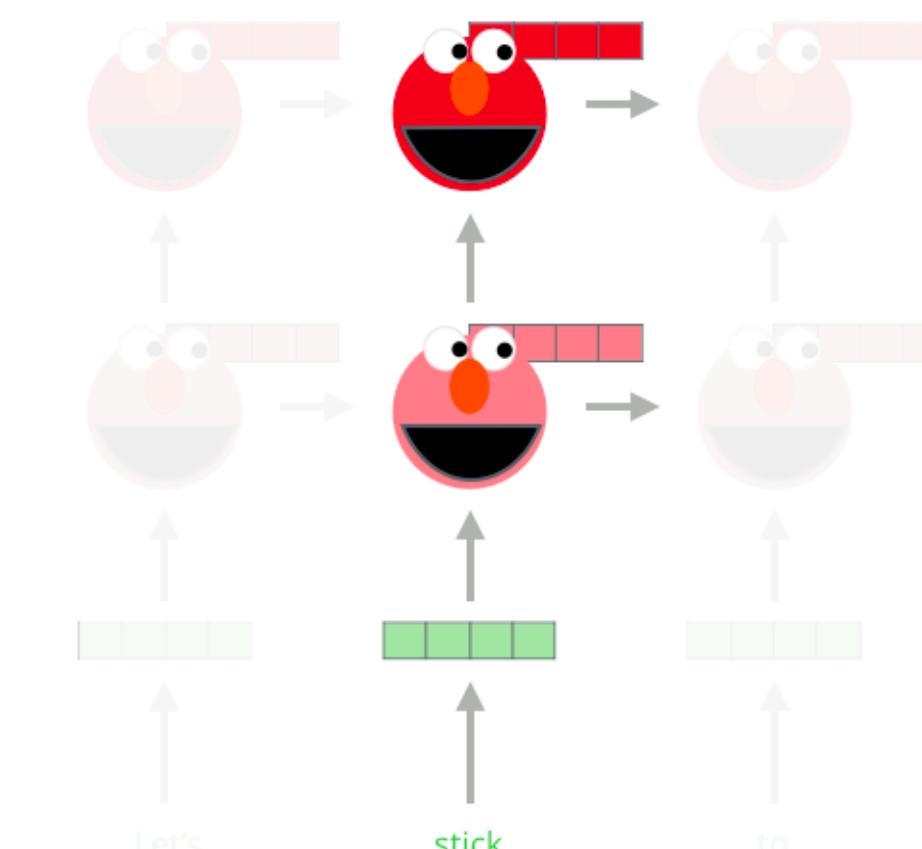
Entrenamiento

Embedding of “stick” in “Let’s stick to” - Step #2

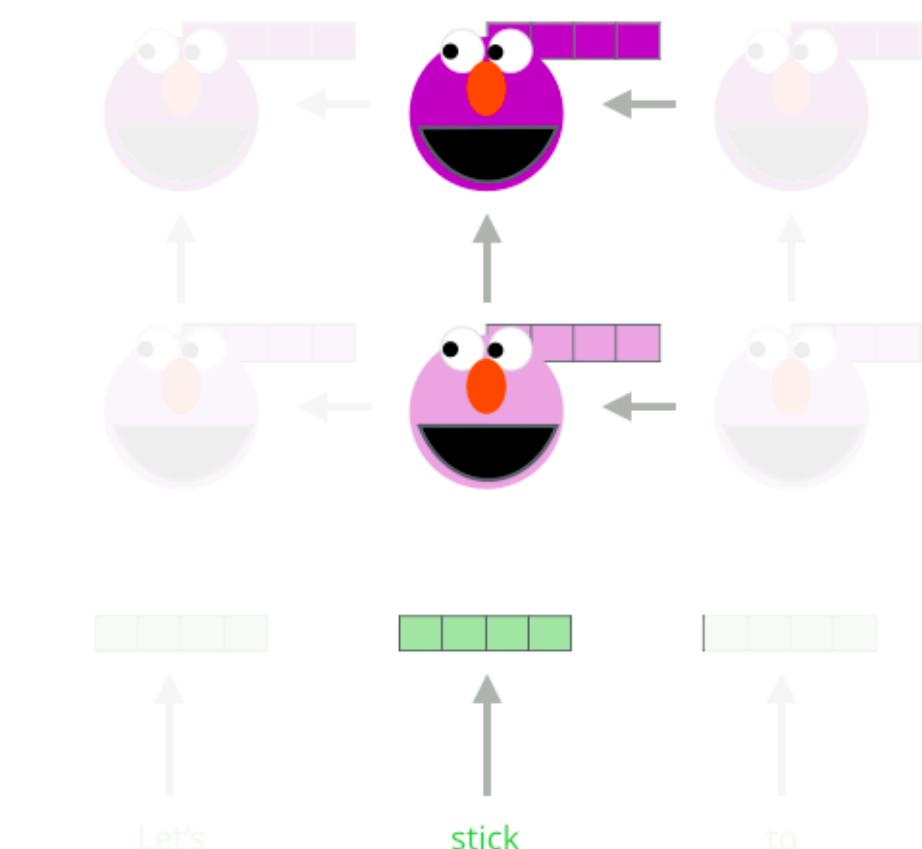
1- Concatenate hidden layers



Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task



3- Sum the (now weighted) vectors



ELMo embedding of “stick” for this task in this context

Concatena y hace suma ponderada de vectores

Peters et. al., 2018 en el artículo de ELMo

BERT

Bidirectional Encoder Representation for Transformers

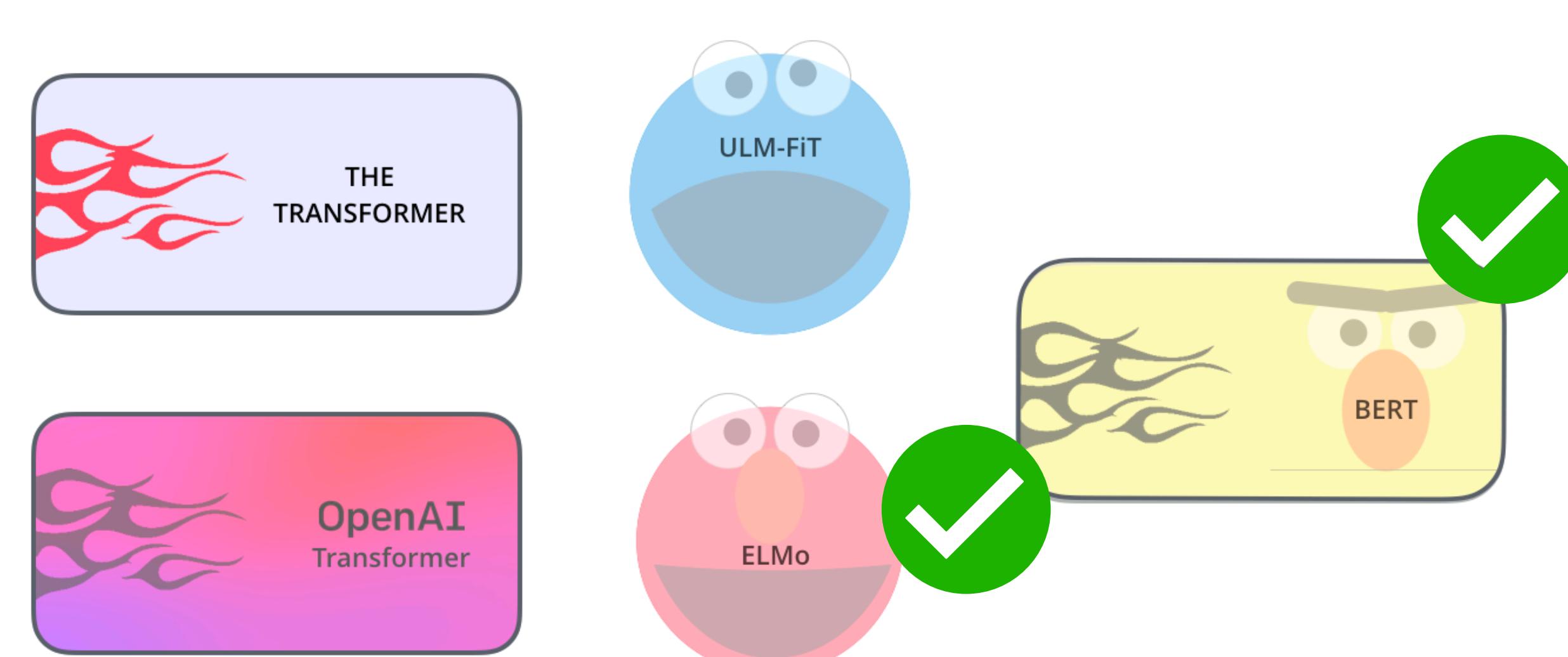
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

BERT family... (RoBERT, RoBERTA, alBERT, Hasta un tal BETO existe)



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

BERT

Bidirectional Encoder Representation for Transformers



*Estructura
Encoder-Decoder*

*Maneja mejor
dependencias a largo
plazo que los LSTM*

*Reemplazaría los LSTM
en algunos campos*



*Modelo de
lenguaje*

*Proceso para ajustar el
modelo a diversas
tareas*

*Introdujo el
aprendizaje por
transferencia*

BERT

Bidirectional Encoder Representation for Transformers

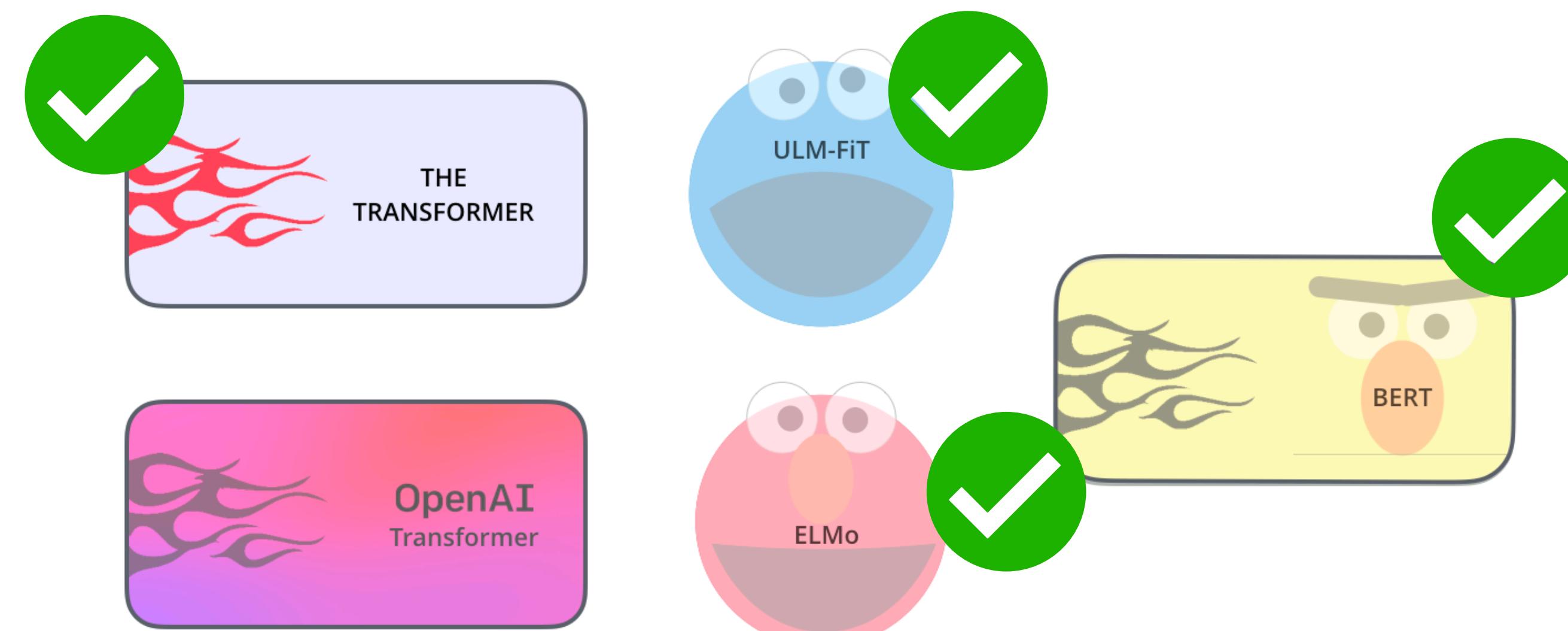
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

BERT family... (RoBERT, RoBERTA, alBERT, Hasta un tal BETO existe)



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

Open AI Transformer

Pre-training a transformer Decoder for Language Modeling

Open AI

Se conforma con el
codificador del
transformer

Función de perdida:

- NSP



Open AI Transformer

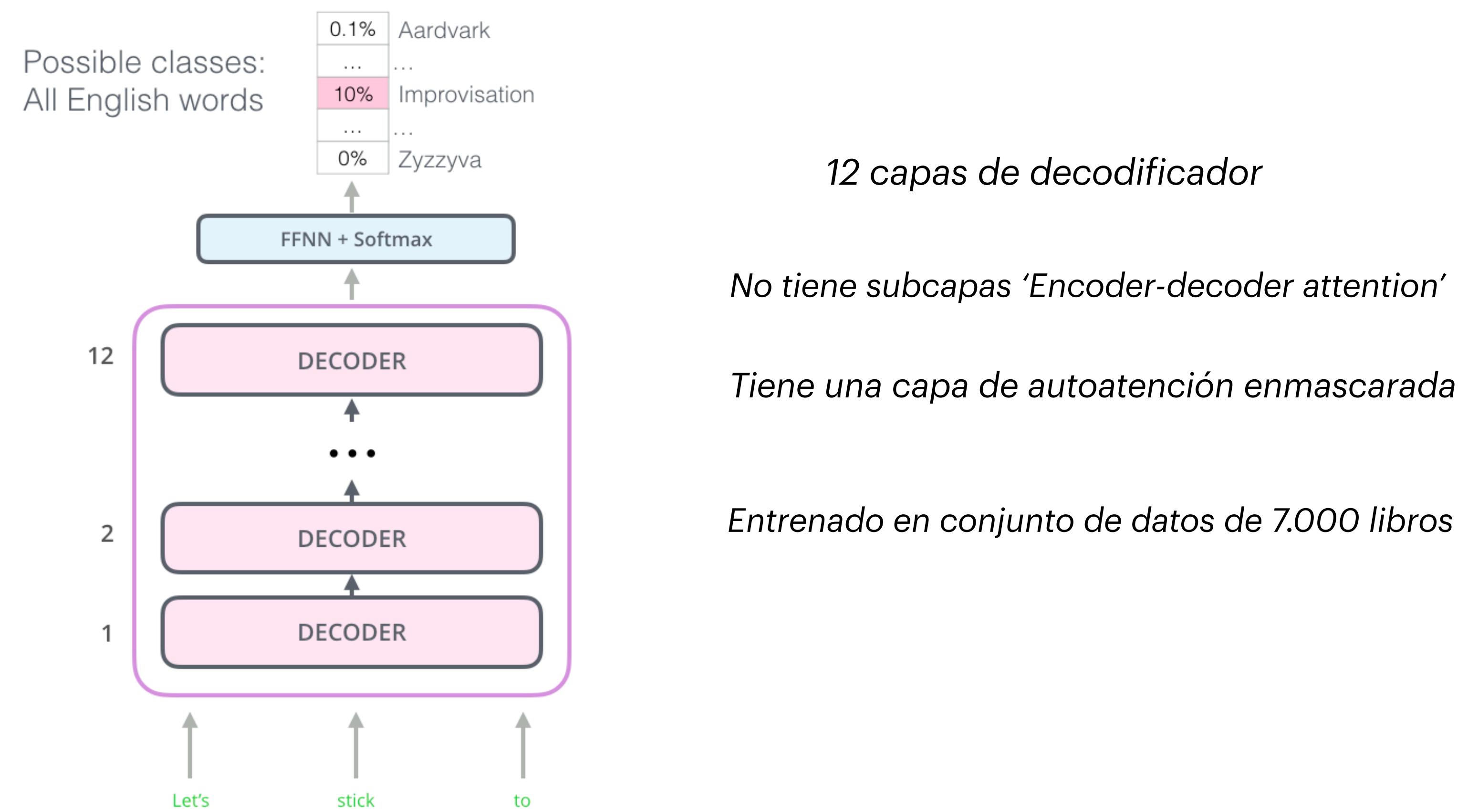
Pre-training a transformer Decoder for Language Modeling

Open AI

Se conforma con el codificador del transformer

Función de perdida:

- NSP



Open AI Transformer

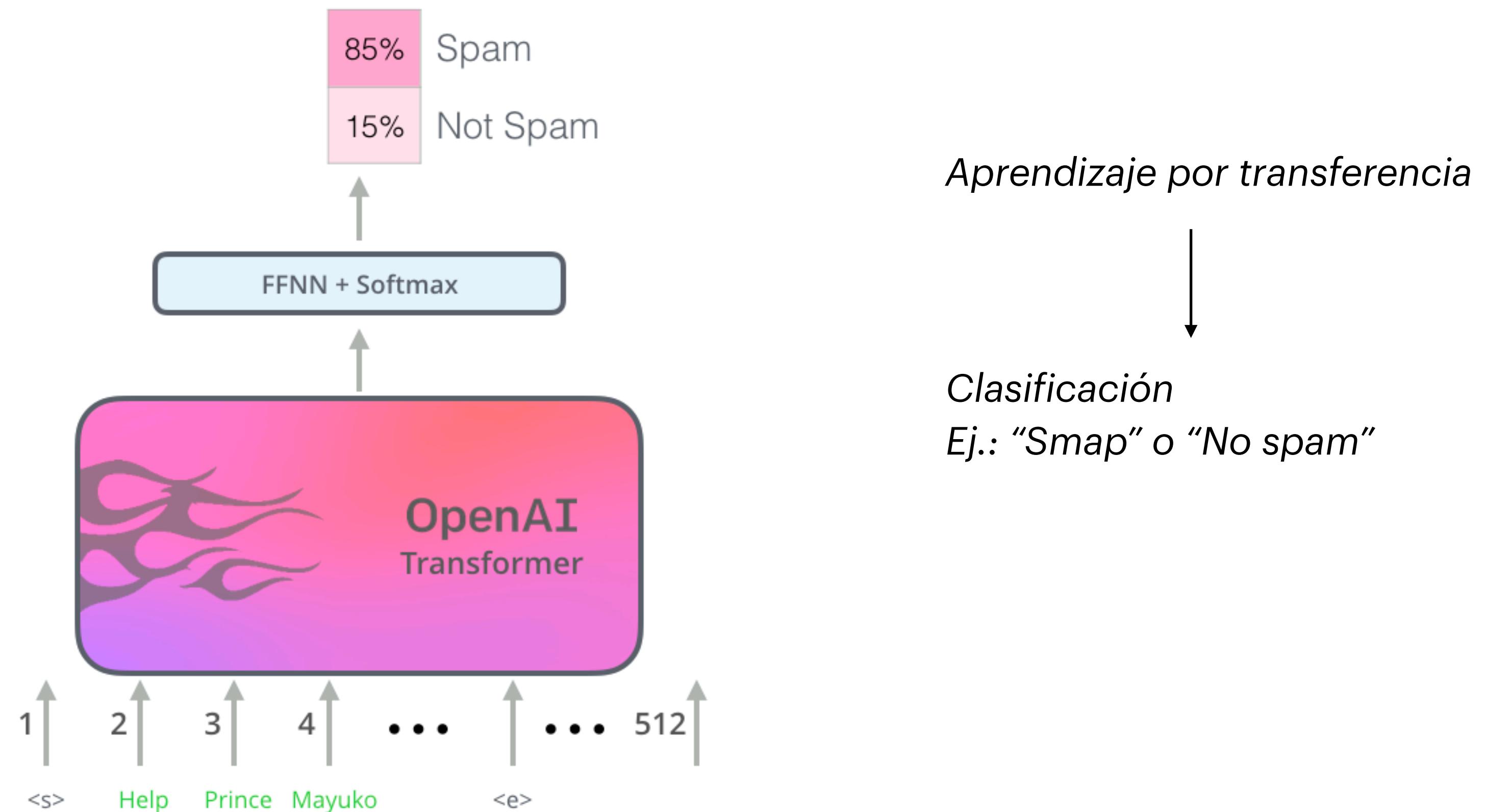
Pre-training a transformer Decoder for Language Modeling

Open AI

Se conforma con el codificador del transformer

Función de perdida:

- NSP



Open AI Transformer

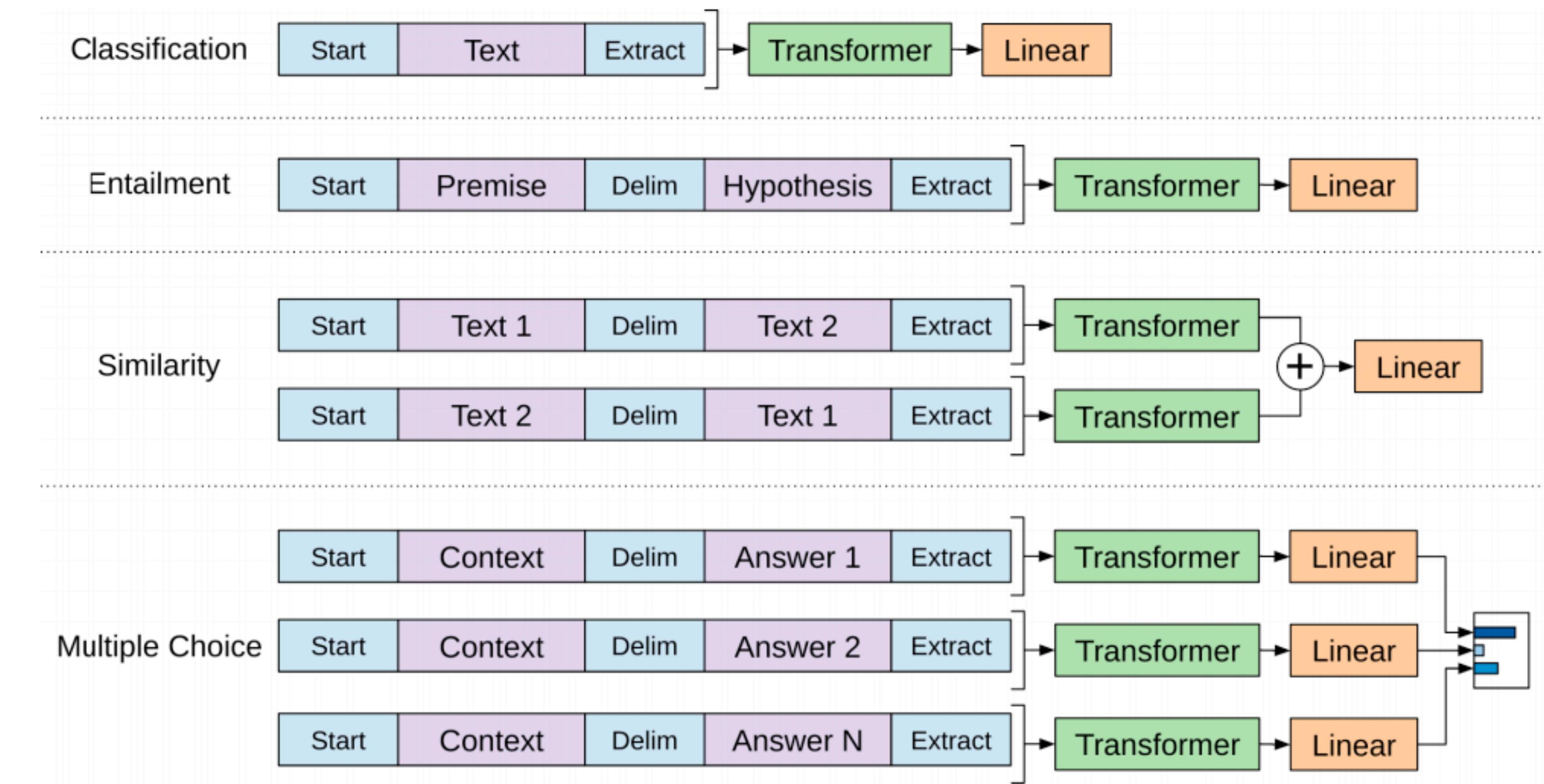
Pre-training a transformer Decoder for Language Modeling

Open AI

Se conforma con el codificador del transformer

Función de perdida:

- NSP



Transformaciones de entrada para manejar diferentes tareas

BERT

Bidirectional Encoder Representation for Transformers

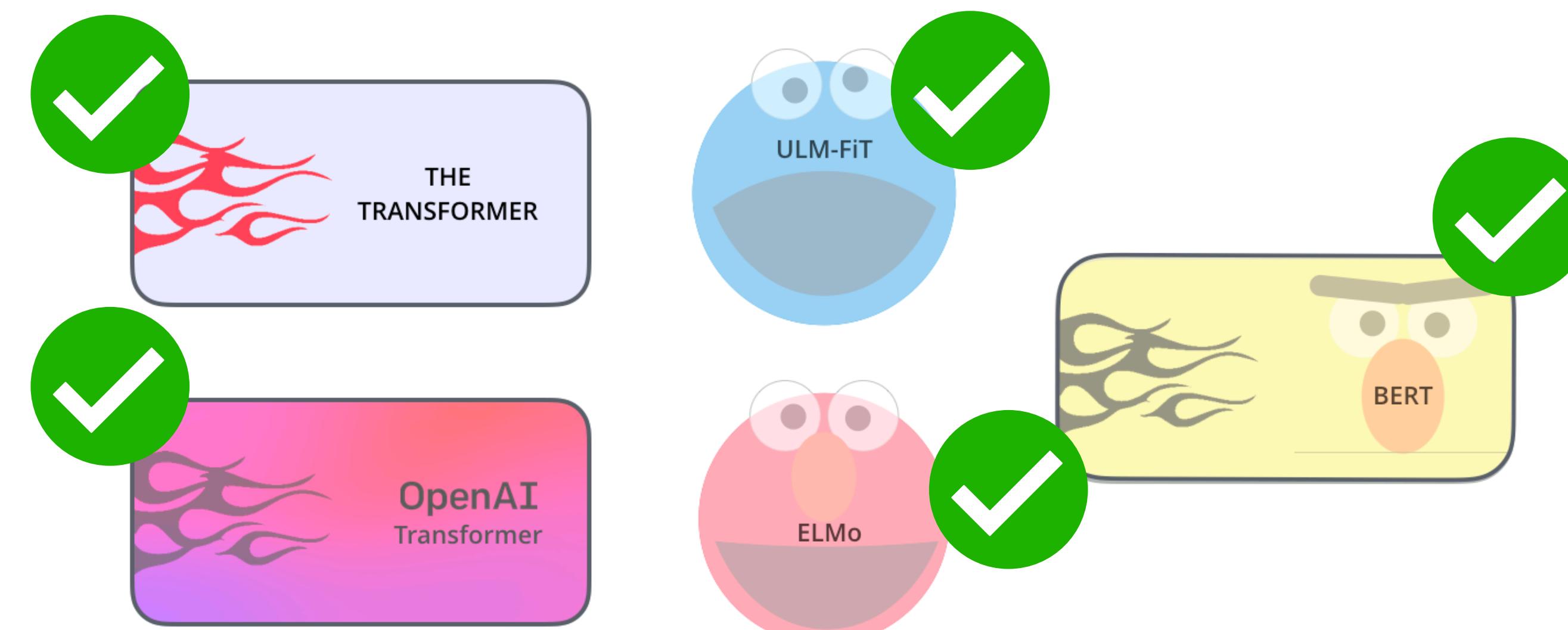
BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

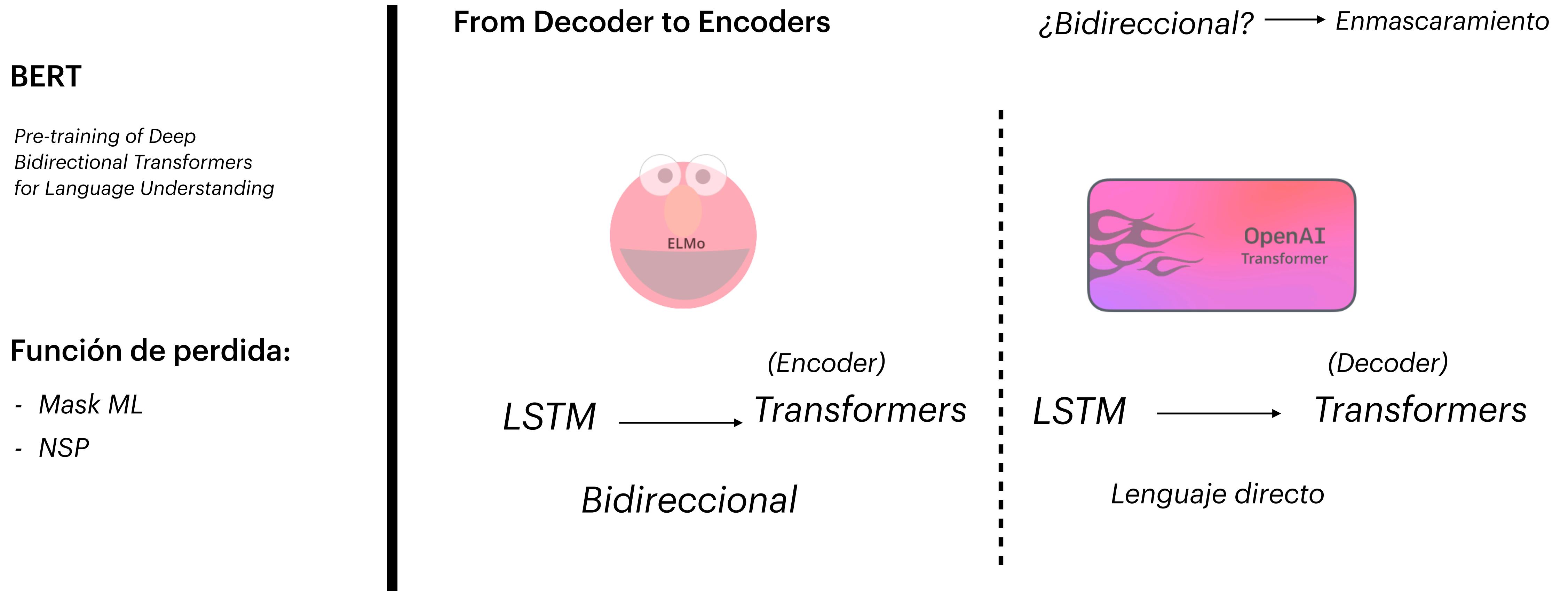
BERT family... (RoBERT, RoBERTA, alBERT, Hasta un tal BETO existe)



(ULMFiT no tiene nada que ver con Cookie Monster. Pero no se me ocurre nada más...)

BERT

Bidirectional Encoder Representation for Transformers



BERT

Bidirectional Encoder Representation for Transformers

BERT

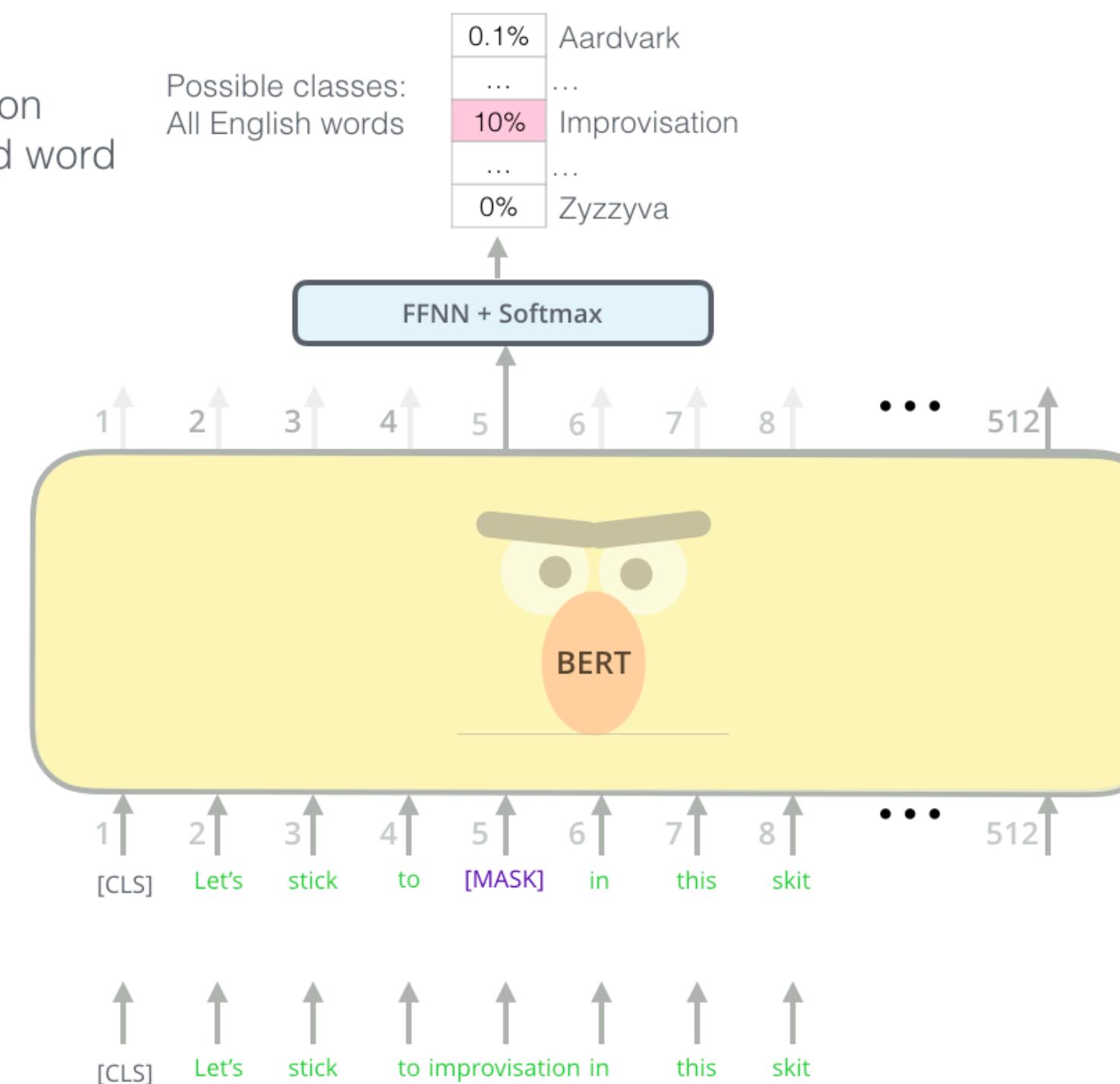
*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Marked Language Model (modelo de lenguaje enmascarado)

Use the output of the
masked word's position
to predict the masked word



*BERT enmascara al 15%
de palabras de entrada*

*Se le pide al modelo
predecir la siguiente
palabra*

(Reemplazo aleatorio)

BERT

Bidirectional Encoder Representation for Transformers

BERT

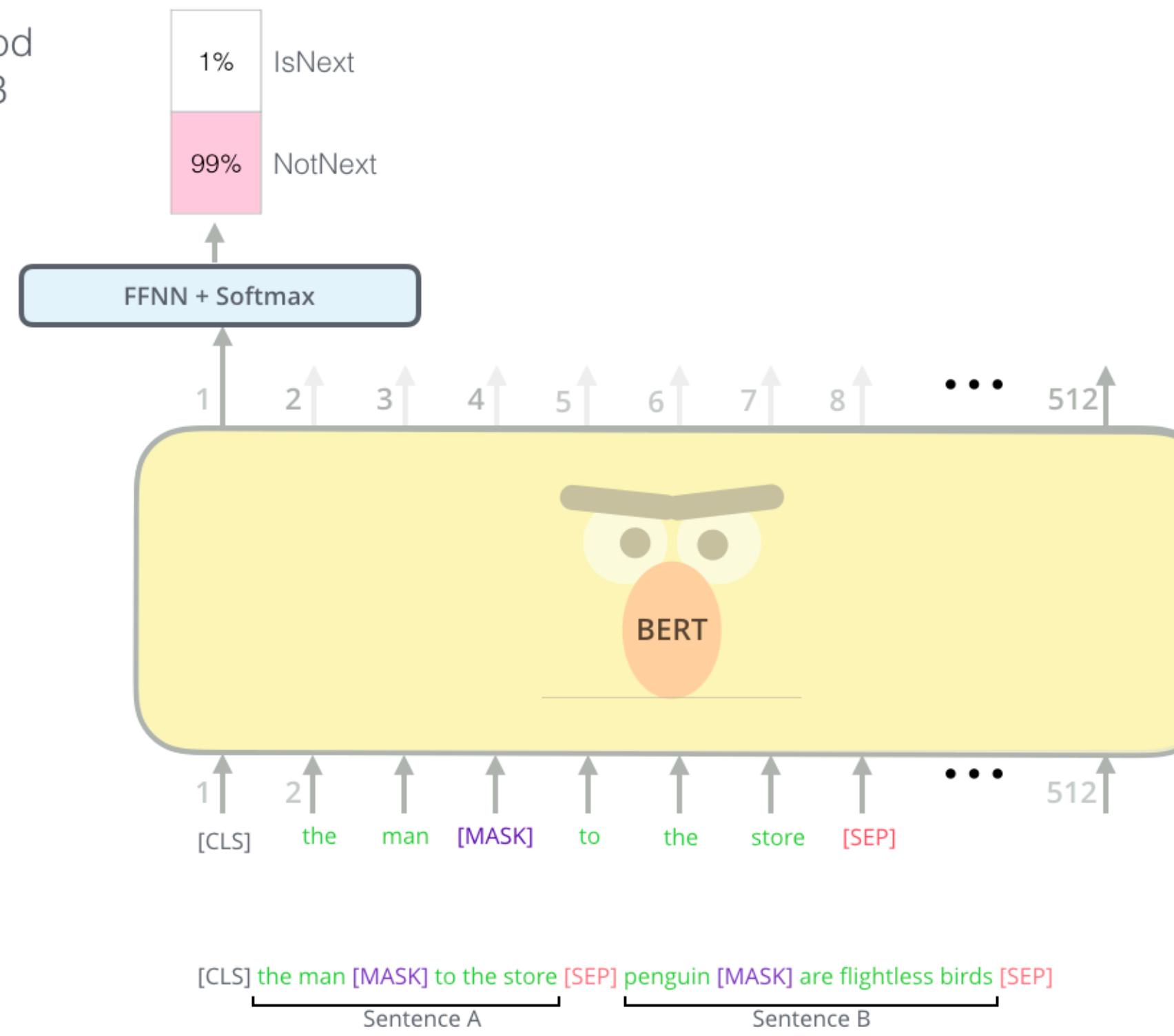
*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

- Mask ML
- NSP

Two-sentence Tasks (Tarea de dos oraciones)

Predict likelihood
that sentence B
belongs after
sentence A



*Clasifica la siguiente
oración. Si continua o
no.*

*Bert usa WordPieces,
no token por palabras.*

BERT

Bidirectional Encoder Representation for Transformers

BERT

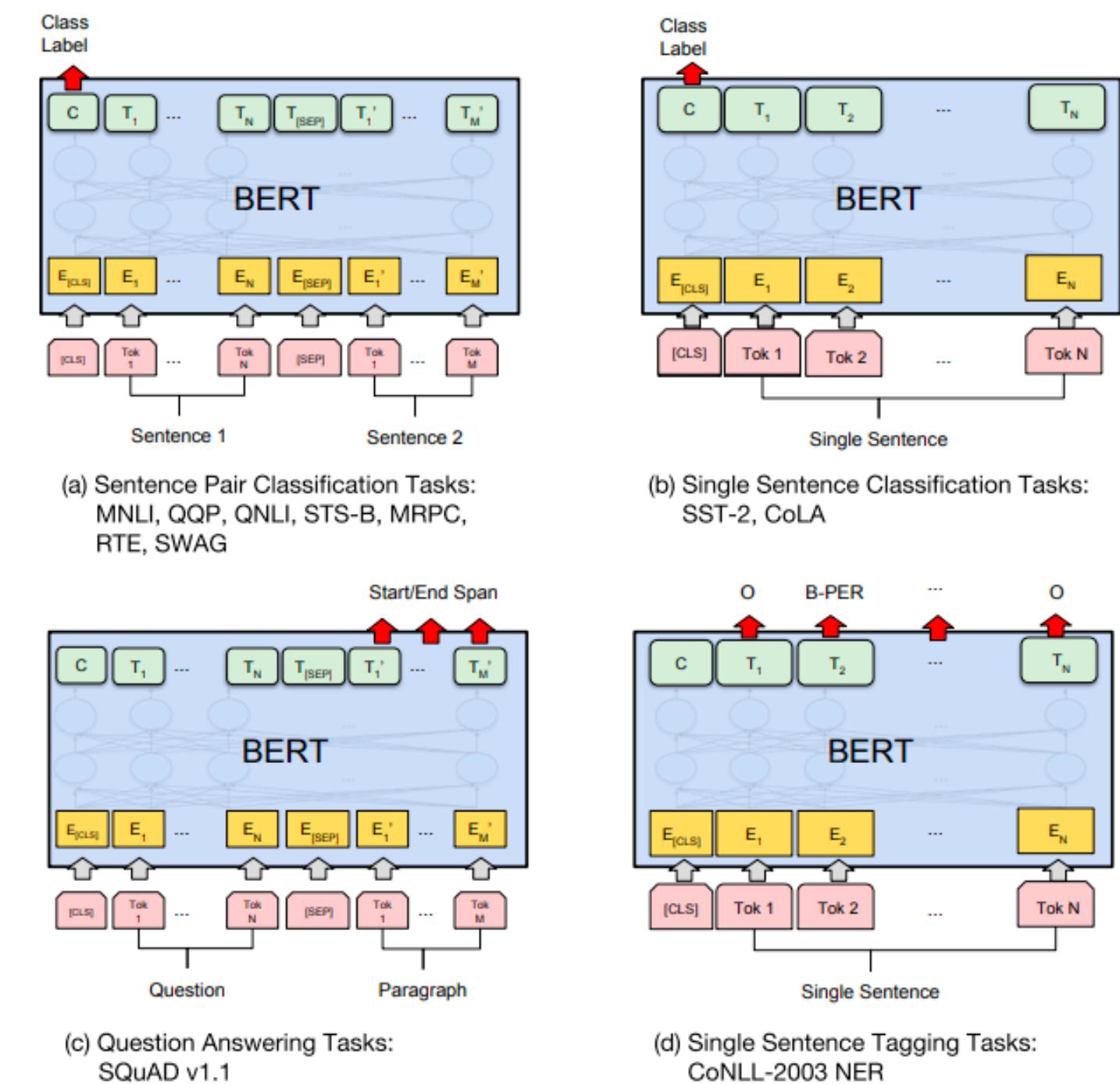
Pre-training of Deep Bidirectional Transformers for Language Understanding



Función de perdida:

- Mask ML
- NSP

Task specific-Models (BERT for different tasks)



Sentence pair classification task

Single sentence classification tasks

Question Answering tasks

Single sentence tagging tasks

BERT

Bidirectional Encoder Representation for Transformers

BERT

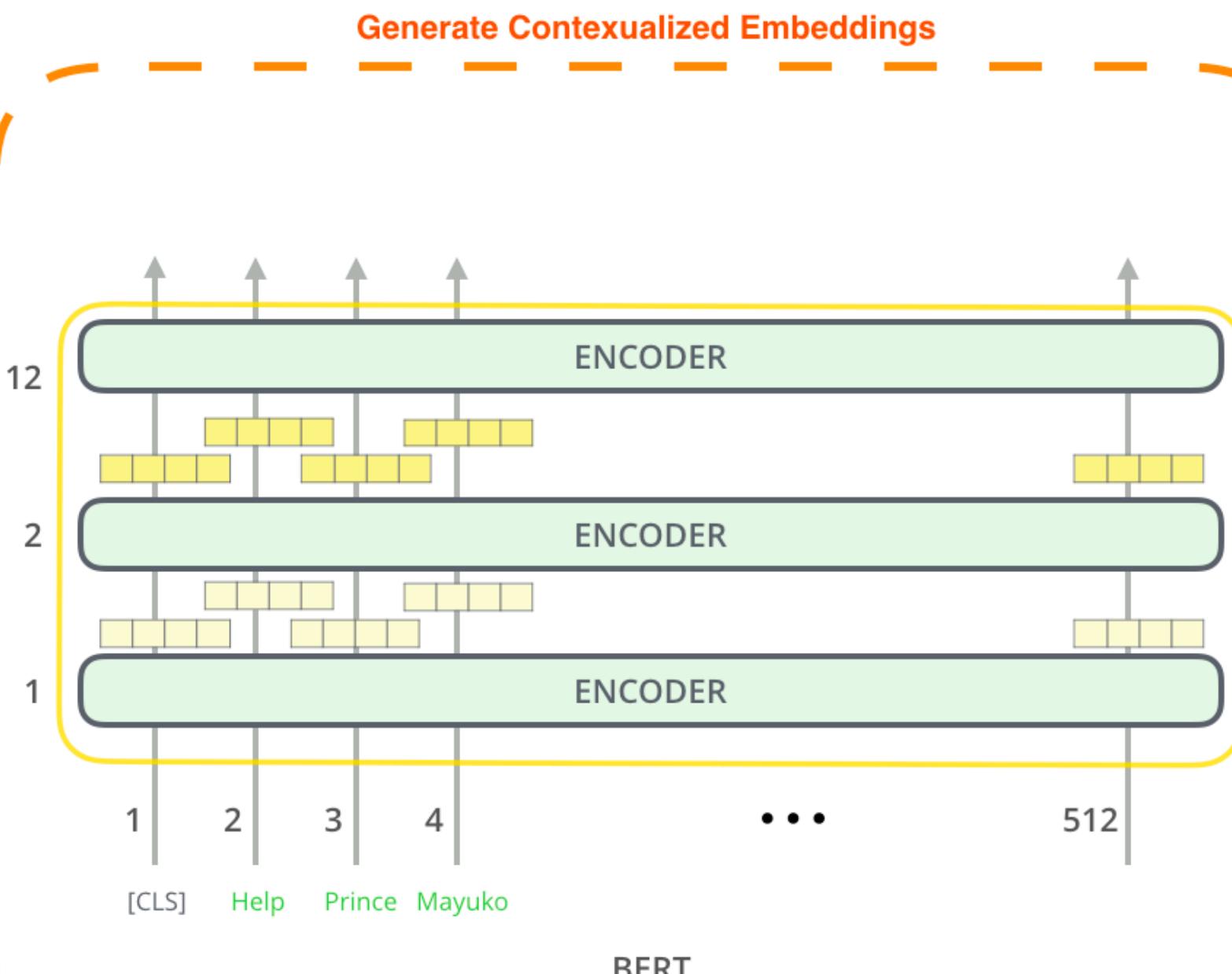
*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*

Función de perdida:

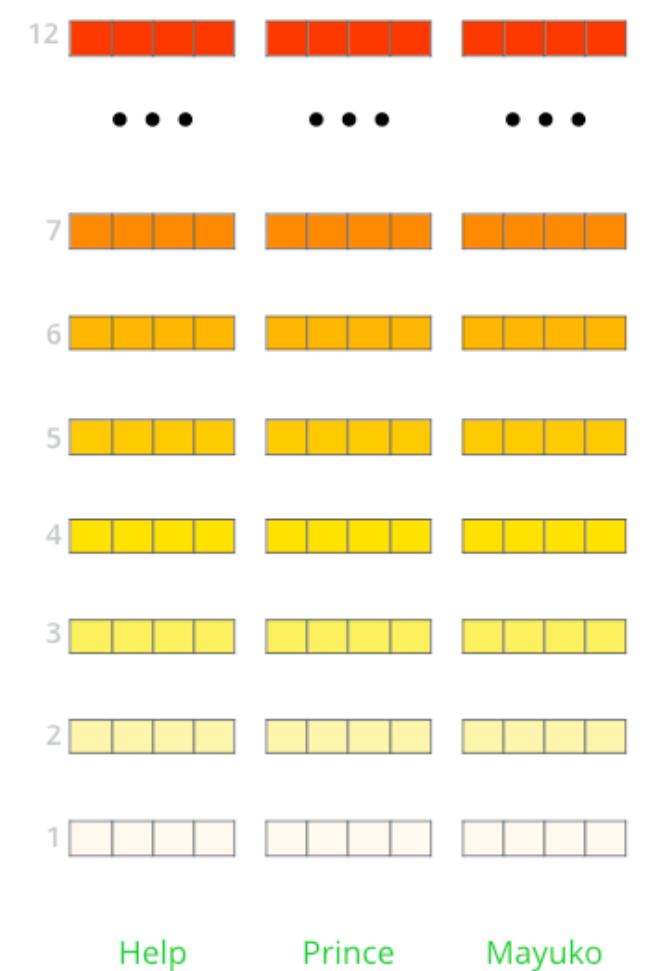
- Mask ML
- NSP

BERT for feature extraction

Like ELMo, with the pre-trained BERT create contextualized word embeddings



The output of each encoder layer along each token's path can be used as a feature representing that token.



But which one should we use?

BERT

Bidirectional Encoder Representation for Transformers

BERT

*Pre-training of Deep
Bidirectional Transformers
for Language Understanding*



Función de perdida:

- Mask ML
- NSP

BERT for feature extraction

Like ELMo, with the pre-trained BERT create contextualized word embeddings

What is the best contextualized embedding for “**Help**” in that context?
For named-entity recognition task CoNLL-2003 NER

