# BIDIRECTIONAL ENCODER REPRESENTATIONS FROM TRANSFORMERS

**Pedro Lelis** • Cientista de Dados • pedrolelis@ciandt.com

BERT

**Pedro Lelis** • Cientista de Dados • pedrolelis@ciandt.com

# Word Embedding

"The fat cat sat on the mat."
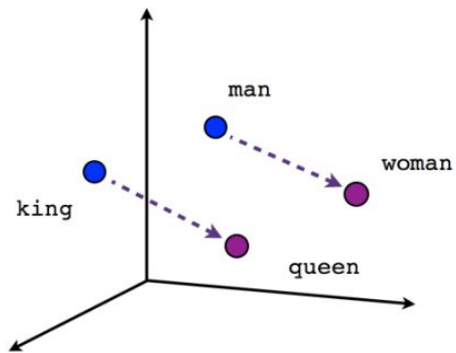
CBOW          Skip-gram
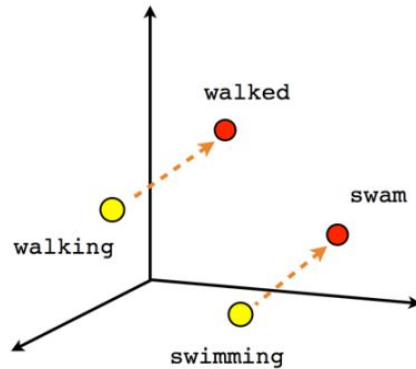
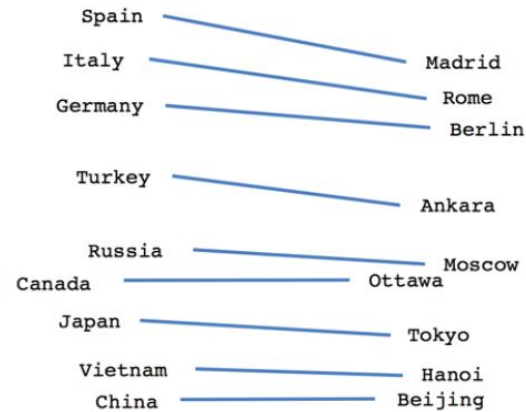# CBOW

$$\max p\,(w\,|\,C)$$
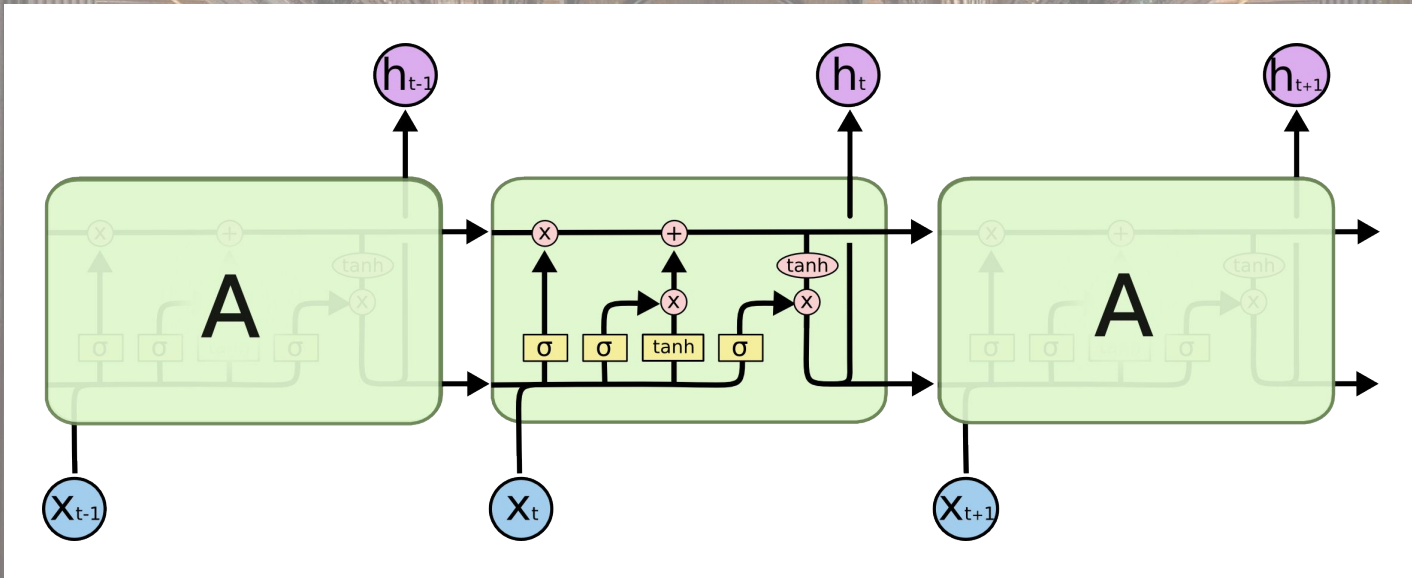
# Skip-gram

$$\max p\,(C\,|\,w)$$

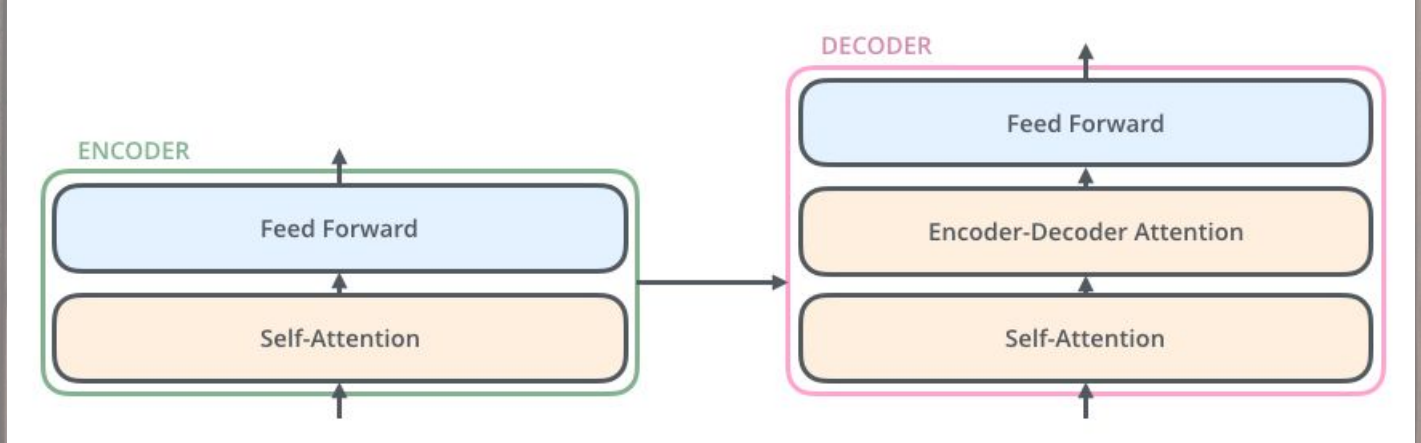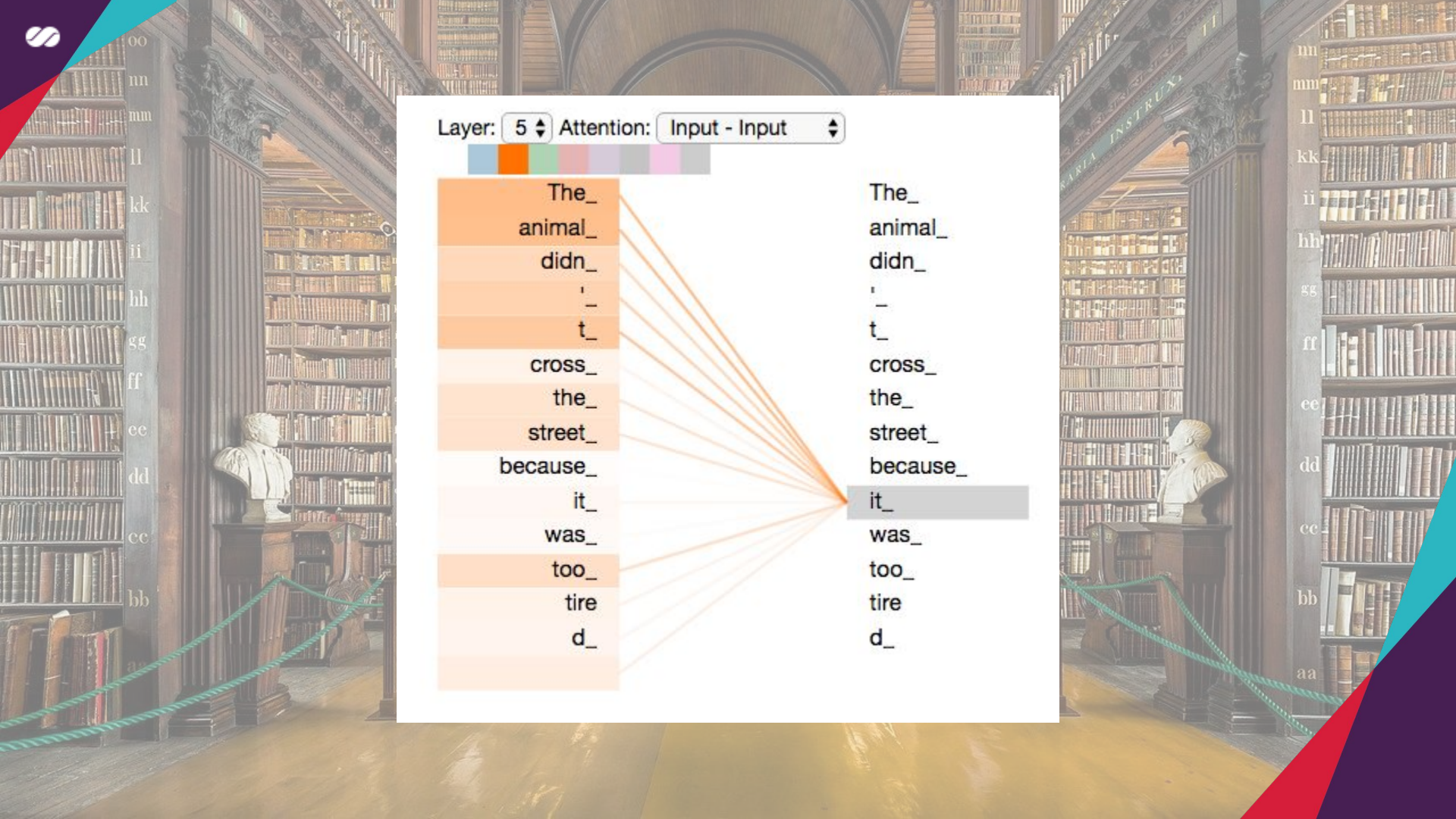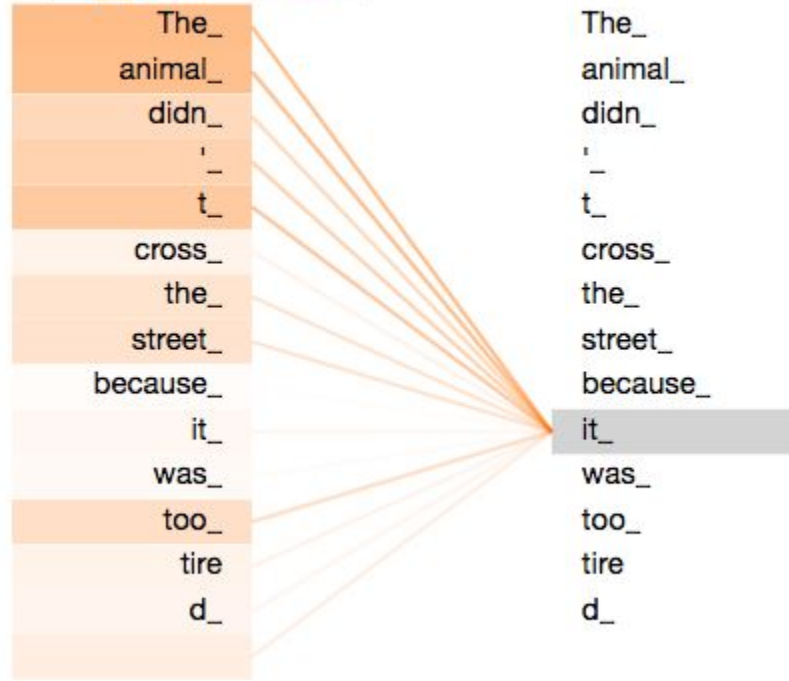Male-Female    Verb tense    Country-Capital

# Long Short-Term Memory

# Transformer

# 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.

**Semi-supervised Learning Step**

**Model:**

BERT

**Dataset:**

WIKIPEDIA
Die freie Enzyklopädie

**Objective:**

Predict the masked word
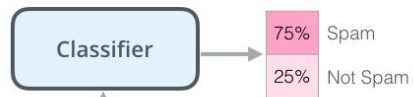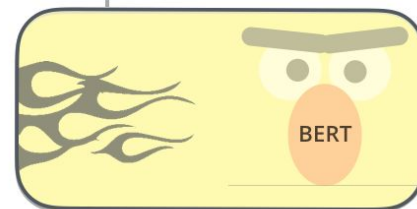(langauge modeling)

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

**Supervised Learning Step**

Classifier → 75% Spam / 25% Not Spam

**Model:**
(pre-trained in step #1)

BERT

**Dataset:**

| Email message | Class |
|---|---|
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

# ELMo

ULMFiT

# BERT

Input = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

Pre-training

Fine-Tuning

(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks: SST-2, CoLA

(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

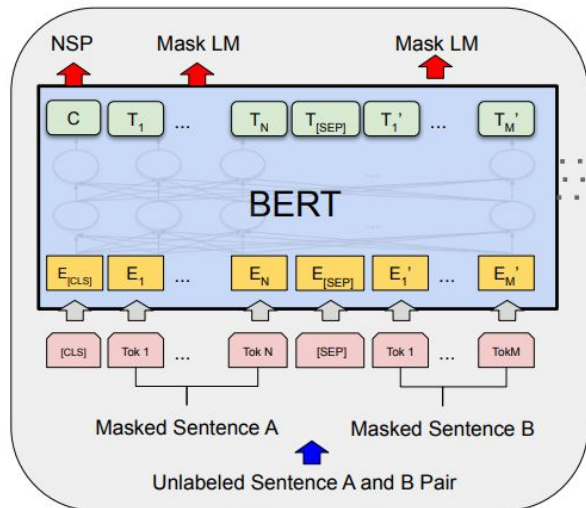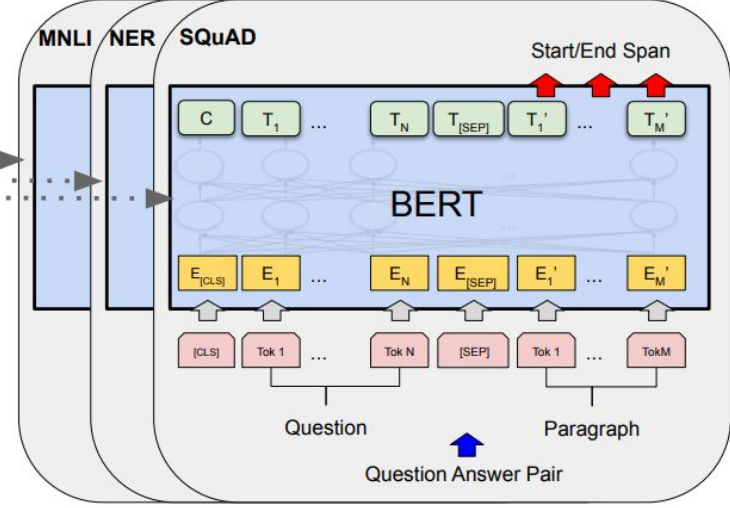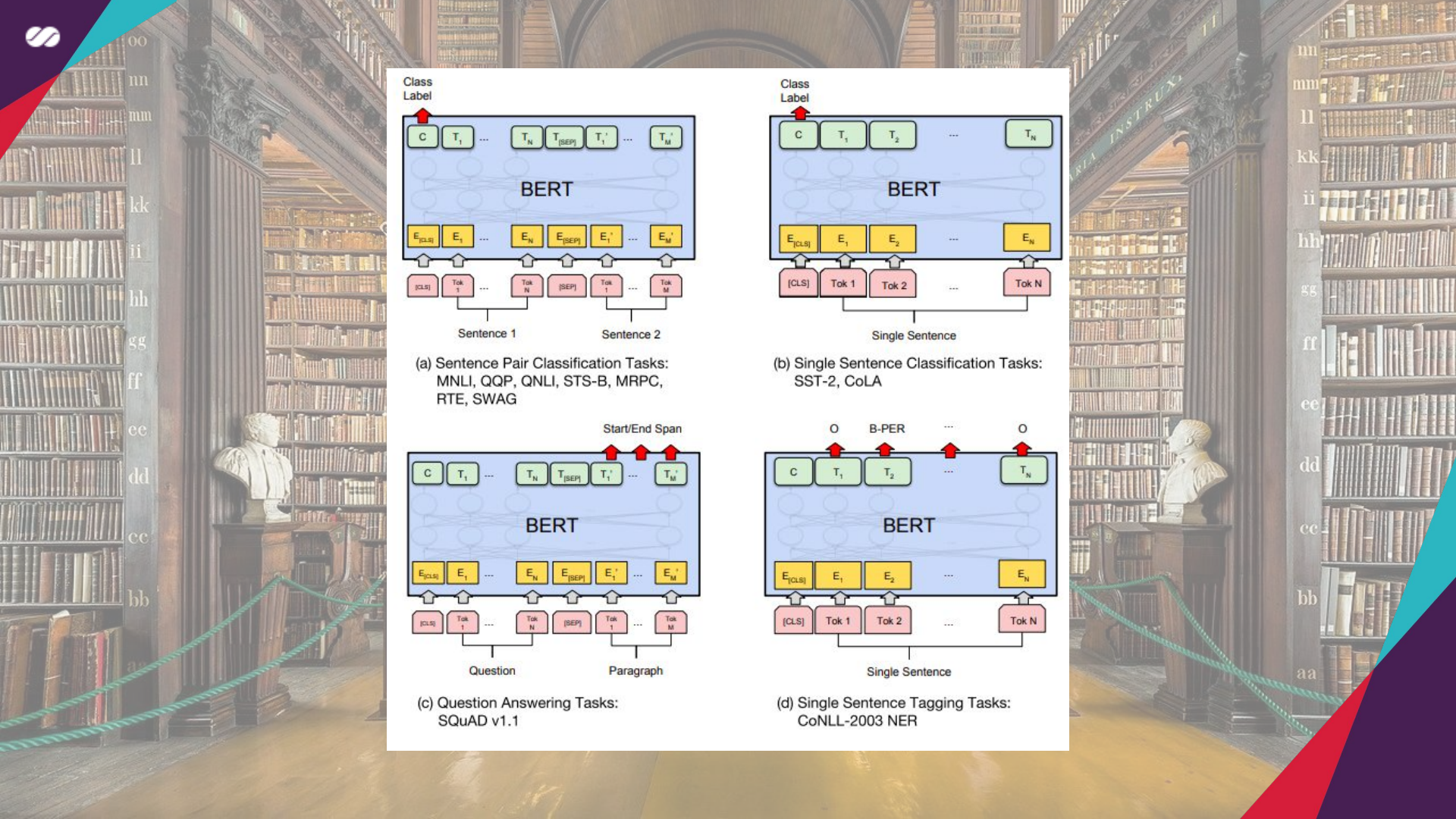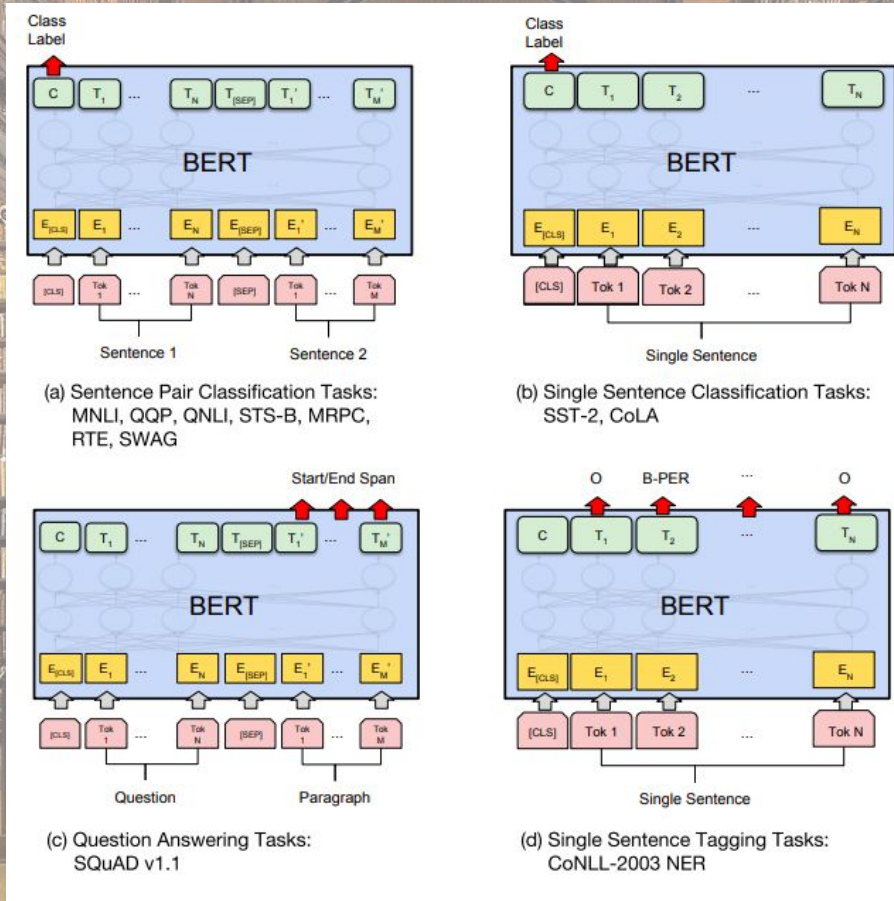| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | Average - |
|---|---|---|---|---|---|---|---|---|---|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| $BERT_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| $BERT_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

Table 1: GLUE Test results, scored by the evaluation server (https://gluebenchmark.com/leaderboard). The number below each task denotes the number of training examples. The "Average" column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.[8] BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

# Pedro Lelis

Moved by curiosity, I'm an entrepreneurial scientist who seeks to improve people's lives through Artificial Intelligence solutions.

Main researches and interesses: Fairness in Machine Learning, Few-shot Learning, Transformer (Attention), Generative Adversarial Network, Reinforcement Learning and Memory-Augmented Neural Network.

**https://forms.gle/eWuLqj3yoQNg15SC8**

CIST
NO
TDCSP