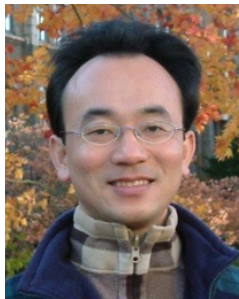# 人工智慧
# (Artificial Intelligence)

# 深度學習自然語言處理
# (Deep Learning for Natural Language Processing)

1092AI10
MBA, IM, NTPU (M5010) (Spring 2021)
Wed 2, 3, 4 (9:10-12:00) (B8F40)

**Min-Yuh Day**
**戴敏育**
**Associate Professor**
副教授
**Institute of Information Management**, **National Taipei University**
國立臺北大學 資訊管理研究所

https://web.ntpu.edu.tw/~myday

2021-05-26

# 課程大綱 (Syllabus)

週次 (Week)    日期 (Date)    內容 (Subject/Topics)

1  2021/02/24  人工智慧概論
(Introduction to Artificial Intelligence)

2  2021/03/03  人工智慧和智慧代理人
(Artificial Intelligence and Intelligent Agents)

3  2021/03/10  問題解決
(Problem Solving)

4  2021/03/17  知識推理和知識表達
(Knowledge, Reasoning and Knowledge Representation)

5  2021/03/24  不確定知識和推理
(Uncertain Knowledge and Reasoning)

6  2021/03/31  人工智慧個案研究 I
(Case Study on Artificial Intelligence I)

# 課程大綱 **(Syllabus)**

週次 (Week)　　日期 (Date)　　內容 (Subject/Topics)

7　2021/04/07 放假一天 (Day off)

8　2021/04/14 機器學習與監督式學習
　　　　　　　(Machine Learning and Supervised Learning)

9　2021/04/21 期中報告
　　　　　　　 (Midterm Project Report)

10　2021/04/28 學習理論與綜合學習
　　　　　　　　(The Theory of Learning and Ensemble Learning)

11　2021/05/05 深度學習
　　　　　　　　(Deep Learning)

12　2021/05/12 人工智慧個案研究 II
　　　　　　　　(Case Study on Artificial Intelligence II)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

13  2021/05/19  強化學習
(Reinforcement Learning)

14  2021/05/26  深度學習自然語言處理
(Deep Learning for Natural Language Processing)

15  2021/06/02  機器人技術
(Robotics)

16  2021/06/09  人工智慧哲學與倫理，人工智慧的未來
 (Philosophy and Ethics of AI, The Future of AI)

17  2021/06/16  期末報告 I

(Final Project Report I)

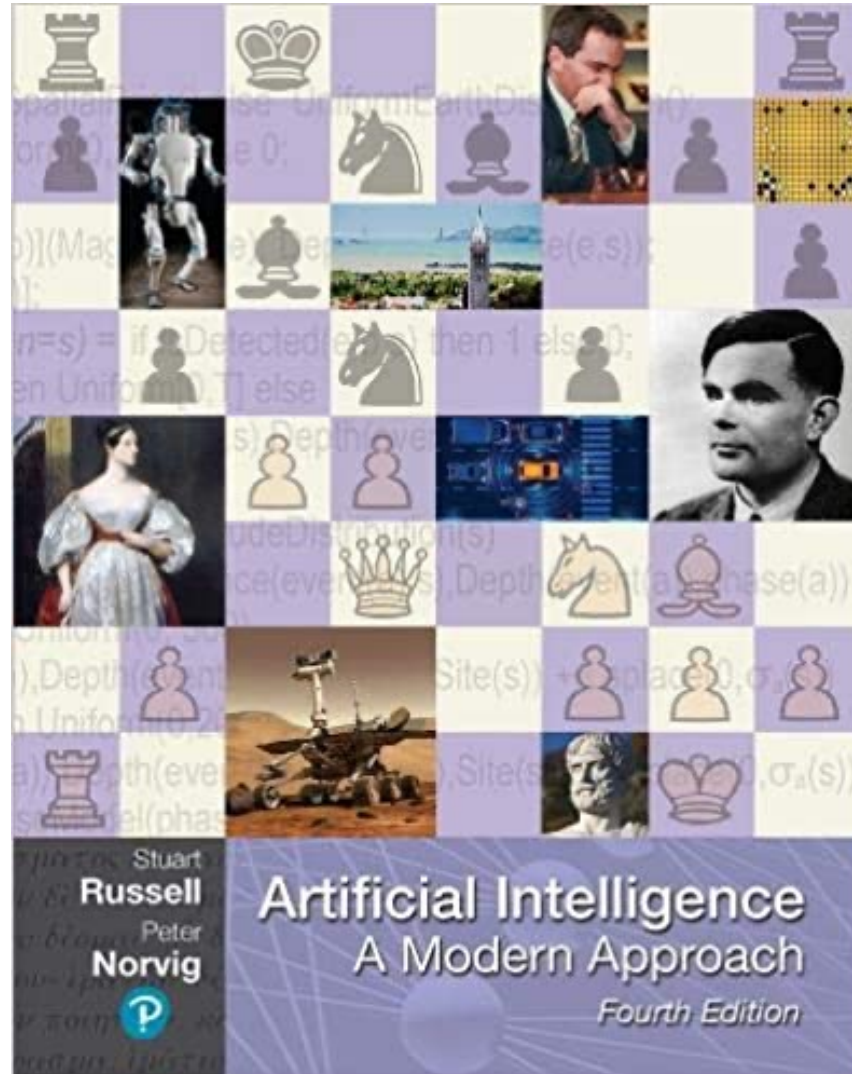18  2021/06/23  期末報告 II
(Final Project Report II)

# Deep Learning for Natural Language Processing

# Outline

- Word Embeddings

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- Pretraining and Transfer Learning

- State of the art (SOTA)

# Stuart Russell and Peter Norvig (2020),
# Artificial Intelligence: A Modern Approach,
## 4th Edition, Pearson



Source: Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson

https://www.amazon.com/Artificial-Intelligence-A-Modern-Approach/dp/0134610997/

# Artificial Intelligence:
# A Modern Approach

1. Artificial Intelligence

2. Problem Solving

3. Knowledge and Reasoning

4. Uncertain Knowledge and Reasoning

5. Machine Learning

6. Communicating, Perceiving, and Acting

7. Philosophy and Ethics of AI

# Artificial Intelligence: Communicating, perceiving, and acting

# Artificial Intelligence:
## 6. Communicating, Perceiving, and Acting

- Natural Language Processing

- Deep Learning for Natural Language Processing

- Computer Vision

- Robotics

# Artificial Intelligence:
## Natural Language Processing

- Language Models

- Grammar

- Parsing

- Augmented Grammars

- Complications of Real Natural Language

- Natural Language Tasks

# Artificial Intelligence:
# Deep Learning for
# Natural Language Processing

- Word Embeddings

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- Pretraining and Transfer Learning

- State of the art (SOTA)

# Reinforcement Learning (DL)

Agent

Environment

# Reinforcement Learning (DL)



Source: Richard S. Sutton & Andrew G. Barto (2018), Reinforcement Learning: An Introduction, 2nd Edition, A Bradford Book.

14
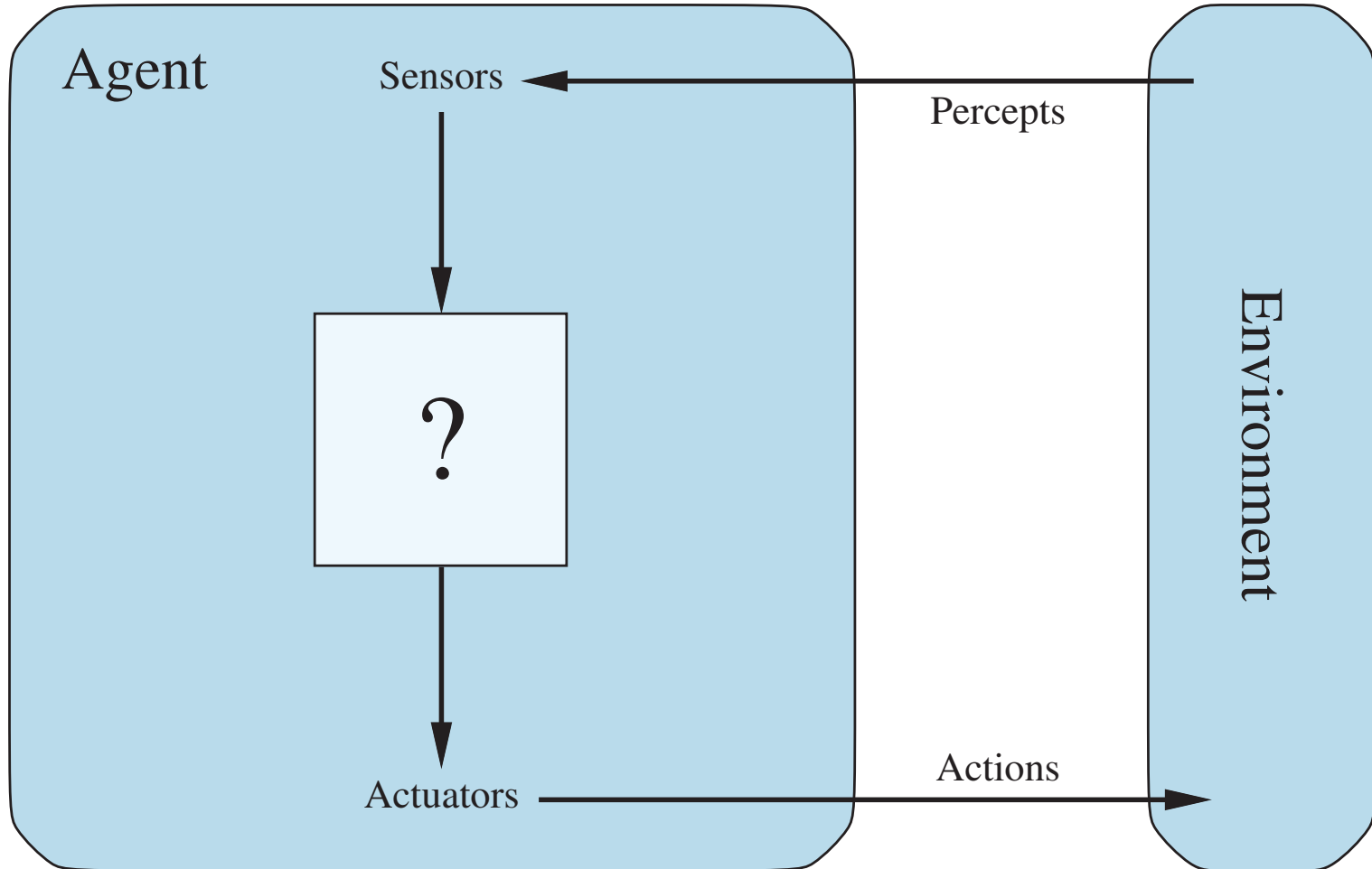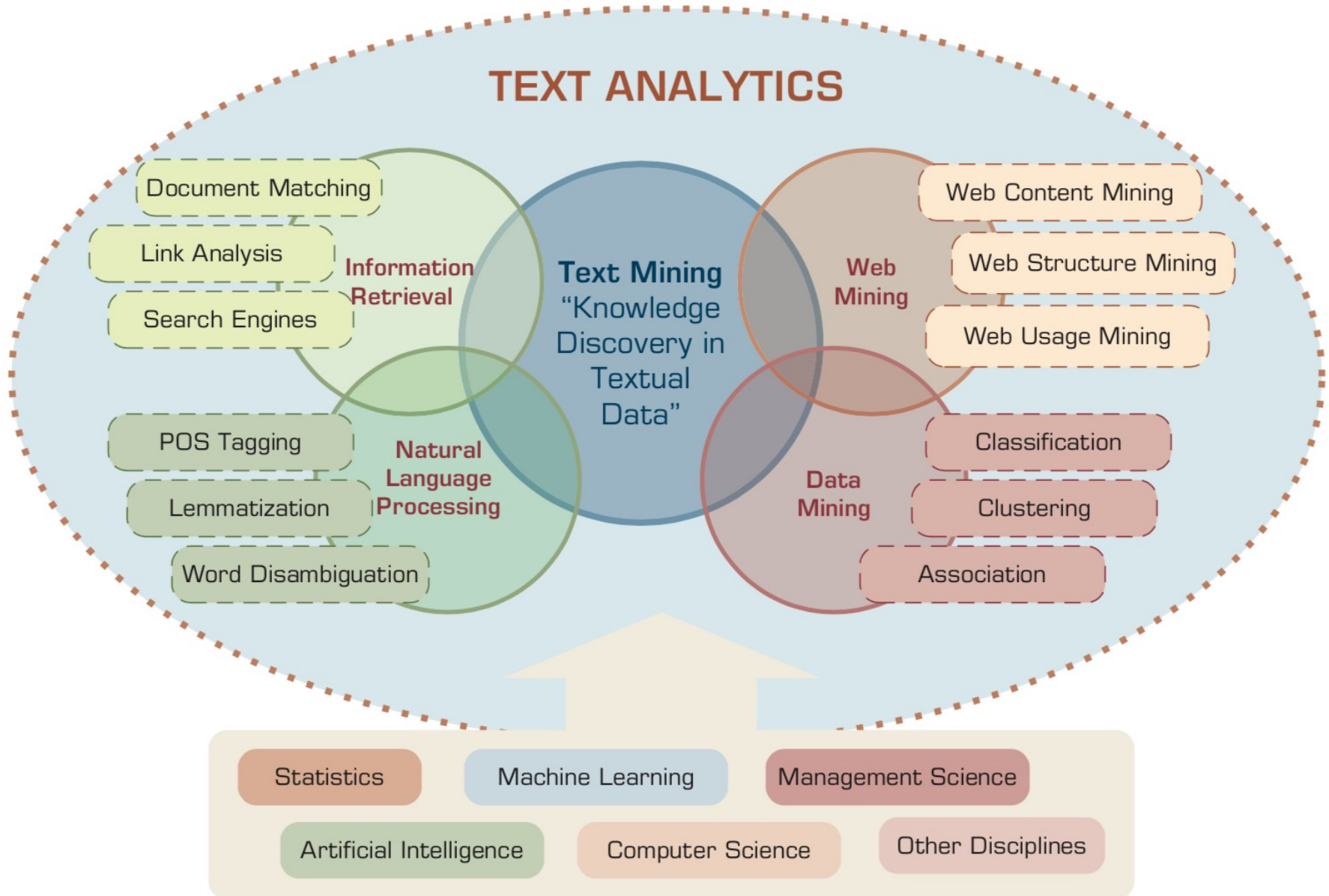
# Reinforcement Learning (DL)

# Agents interact with environments through sensors and actuators

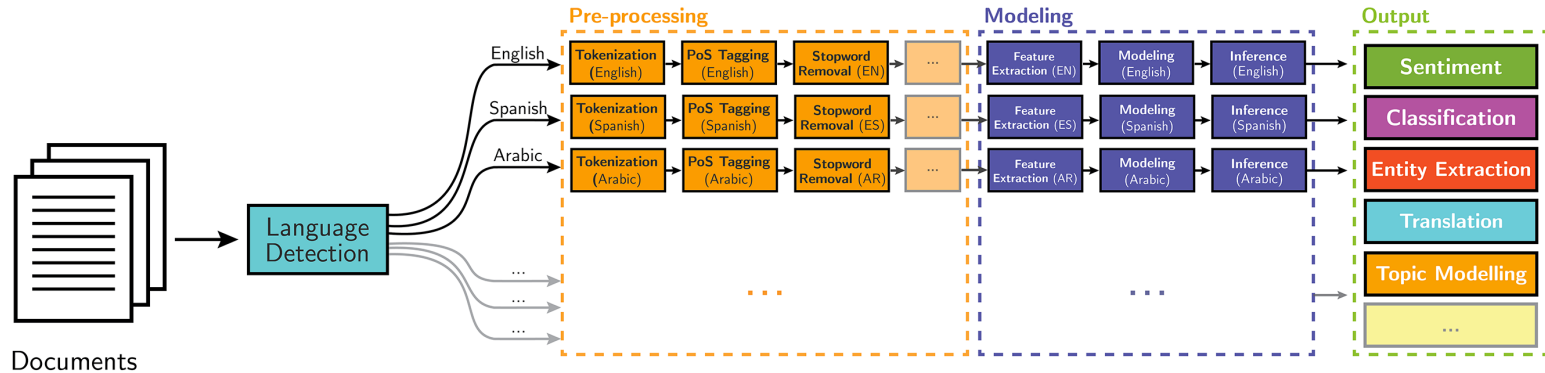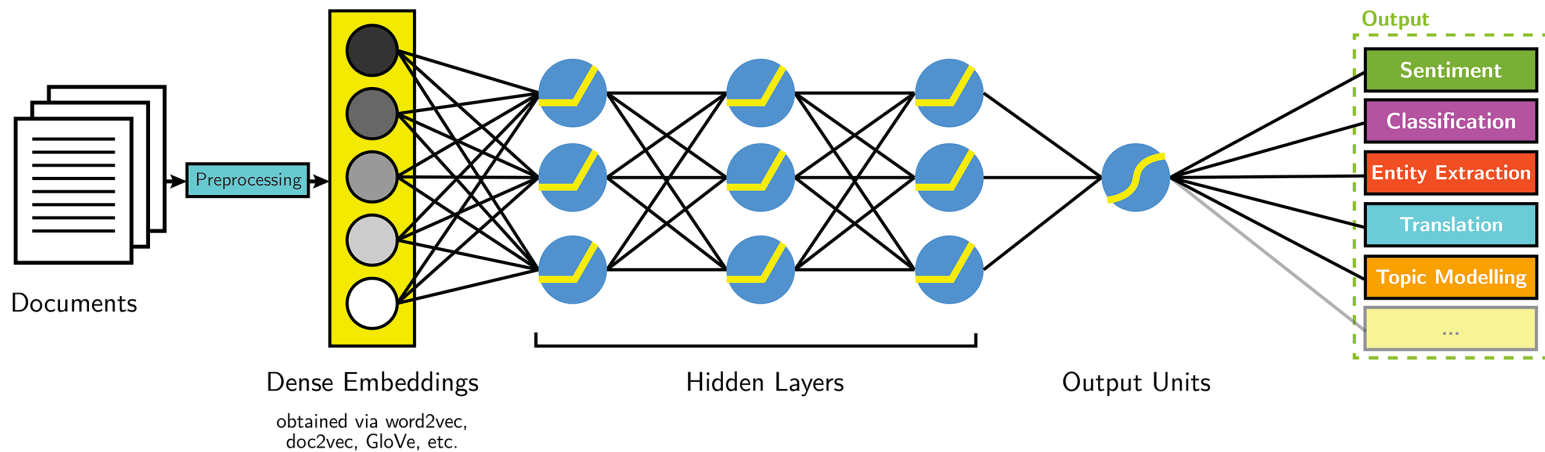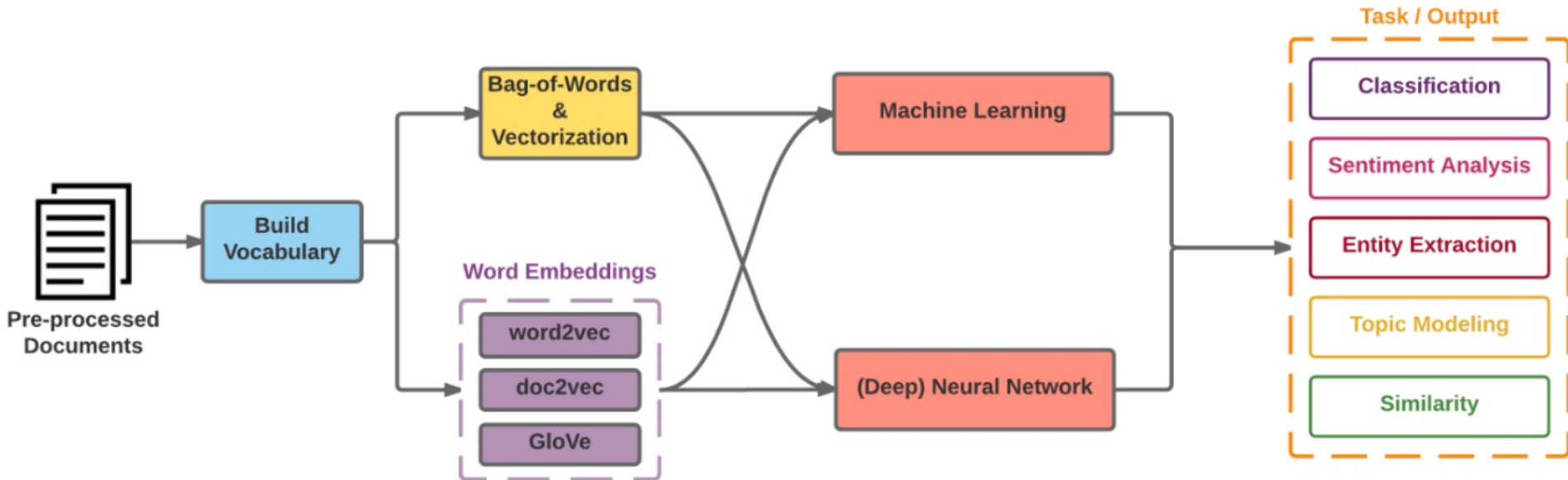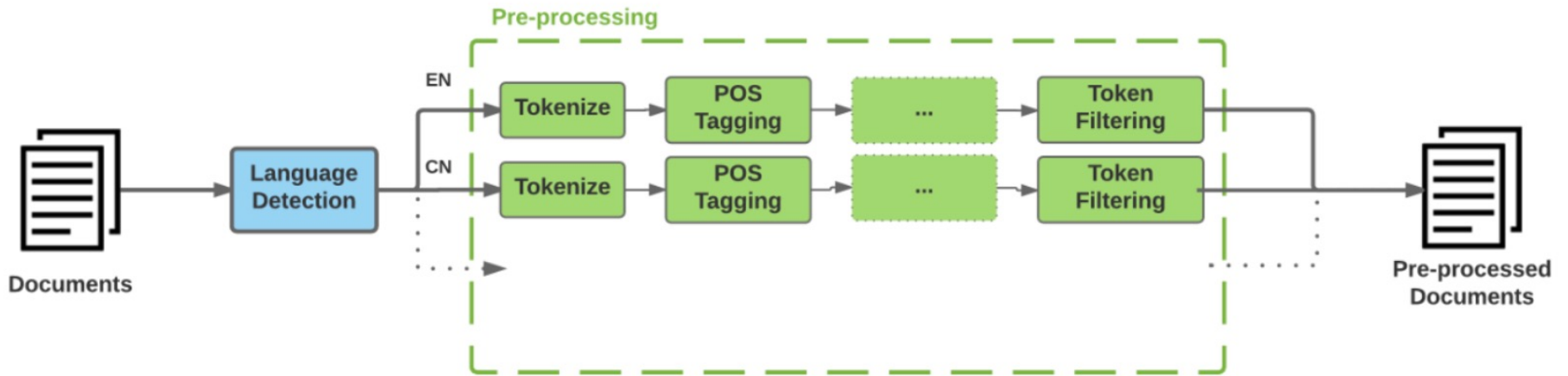# Deep Learning for Natural Language Processing

# AI for Text Analytics

## TEXT ANALYTICS

**Information Retrieval**
- Document Matching
- Link Analysis
- Search Engines

**Natural Language Processing**
- POS Tagging
- Lemmatization
- Word Disambiguation

**Text Mining** "Knowledge Discovery in Textual Data"

**Web Mining**
- Web Content Mining
- Web Structure Mining
- Web Usage Mining

**Data Mining**
- Classification
- Clustering
- Association

- Statistics
- Machine Learning
- Management Science
- Artificial Intelligence
- Computer Science
- Other Disciplines

18

# NLP



## Classical NLP

**Pre-processing** | **Modeling** | **Output**

Documents → Language Detection → English / Spanish / Arabic → ...

Pre-processing: Tokenization, PoS Tagging, Stopword Removal (English/Spanish/Arabic)

Modeling: Feature Extraction, Modeling, Inference (English/Spanish/Arabic)

Output: Sentiment, Classification, Entity Extraction, Translation, Topic Modelling, ...

## Deep Learning-based NLP

Documents → Preprocessing → Dense Embeddings → Hidden Layers → Output Units → Output

**Dense Embeddings** obtained via word2vec, doc2vec, GloVe, etc.

**Hidden Layers**

**Output Units**

Output: Sentiment, Classification, Entity Extraction, Translation, Topic Modelling, ...

AYLIEN

# Modern NLP Pipeline

# Modern NLP Pipeline

# Deep Learning NLP



**Documents** → **Preprocessing** → **Dense Word Embeddings** → **Deep Neural Network** → **Task / Output**

*Pre-generated Lookup*
*OR*
*Generated in 1st level*
*of NeuralNet*

**Task / Output**
- Classification
- Sentiment Analysis
- Entity Extraction
- Topic Modeling
- Document Similarity

# Natural Language Processing (NLP) and Text Mining

| Raw text |
| :---: |

| Sentence Segmentation |
| :---: |

| Tokenization |
| :---: |

| Part-of-Speech (POS) |
| :---: |

| Stop word removal |
| :---: |

| **Stemming** / **Lemmatization** |
| :---: |

| Dependency Parser |
| :---: |

| String Metrics & Matching |
| :---: |

word's stem
am → am
having → hav

word's lemma
am → be
having → have

# Outline

- <span style="color:red">Word Embeddings</span>

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- Pretraining and Transfer Learning

- State of the art (SOTA)

# One-hot encoding

'The mouse ran up the clock' =

```
The     1      [ [0, 1, 0, 0, 0, 0, 0],
mouse   2        [0, 0, 1, 0, 0, 0, 0],
ran     3        [0, 0, 0, 1, 0, 0, 0],
up      4        [0, 0, 0, 0, 1, 0, 0],
the     1        [0, 1, 0, 0, 0, 0, 0],
clock   5        [0, 0, 0, 0, 0, 1, 0] ]
```

```
[0, 1, 2, 3, 4, 5, 6]
```

Source: https://developers.google.com/machine-learning/guides/text-classification/step-3

# Word embedding

## GloVe (trained on 6 billion words of text)
### 100-dimensional word vectors are projected down onto two dimensions



france
greece
germany

nephew
niece
uncle
aunt

car
bicycle
truck

banana
pizza
apple

# Word Embedding model

## answer the question "A is to B as C is to [what]?"

| A | B | C | $D = C + (B - A)$ | Relationship |
|---|---|---|---|---|
| Athens | Greece | Oslo | Norway | *Capital* |
| Astana | Kazakhstan | Harare | Zimbabwe | *Capital* |
| Angola | kwanza | Iran | rial | *Currency* |
| copper | Cu | gold | Au | *Atomic Symbol* |
| Microsoft | Windows | Google | Android | *Operating System* |
| New York | New York Times | Baltimore | Baltimore Sun | *Newspaper* |
| Berlusconi | Silvio | Obama | Barack | *First name* |
| Switzerland | Swiss | Cambodia | Cambodian | *Nationality* |
| Einstein | scientist | Picasso | painter | *Occupation* |
| brother | sister | grandson | granddaughter | *Family Relation* |
| Chicago | Illinois | Stockton | California | *State* |
| possibly | impossibly | ethical | unethical | *Negative* |
| mouse | mice | dollar | dollars | *Plural* |
| easy | easiest | lucky | luckiest | *Superlative* |
| walking | walked | swimming | swam | *Past tense* |

# Word embeddings



Male-Female

Verb Tense

Country-Capital

# Word embeddings

Source: https://developers.google.com/machine-learning/guides/text-classification/step-3

# Feedforward part-of-speech (POS) tagging model

Class = PastTenseVerb

| Output Layer |
| --- |

| Hidden Layer 2 |
| --- |

| Hidden Layer 1 |
| --- |

| Embedding lookup | Embedding lookup | Embedding lookup | Embedding lookup | Embedding lookup |
| --- | --- | --- | --- | --- |

Yesterday     they     cut     the     rope

# Universal Sentence Encoder (USE)

- The Universal Sentence Encoder encodes **text** into high-dimensional **vectors** that can be used for text classification, semantic similarity, clustering and other natural language tasks.

- The universal-sentence-encoder model is trained with a **deep averaging network (DAN)** encoder.

# Universal Sentence Encoder (USE)
# Semantic Similarity

# Universal Sentence Encoder (USE) Classification

# Universal Sentence Encoder (USE)

```python
import tensorflow_hub as hub

embed = hub.Module("https://tfhub.dev/google/"
    "universal-sentence-encoder/1")

embedding = embed([
    "The quick brown fox jumps over the lazy dog."])
```

Source: Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Céspedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, Ray Kurzweil. Universal Sentence Encoder. arXiv:1803.11175, 2018.

# Multilingual Universal Sentence Encoder (MUSE)

```python
import tensorflow_hub as hub

module = hub.Module("https://tfhub.dev/google/"
    "universal-sentence-encoder-multilingual/1")

multilingual_embeddings = module([
  "Hola Mundo!", "Bonjour le monde!", "Ciao mondo!"
  "Hello World!", "Hallo Welt!", "Hallo Wereld!",
  "你好世界!", "Привет, мир!", "مرحبا بالعالم!"])
```

Source: Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez Abrego , Steve Yuan, Chris Tar, Yun-hsuan Sung, Ray Kurzweil. Multilingual Universal Sentence Encoder for Semantic Retrieval. July 2019

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT



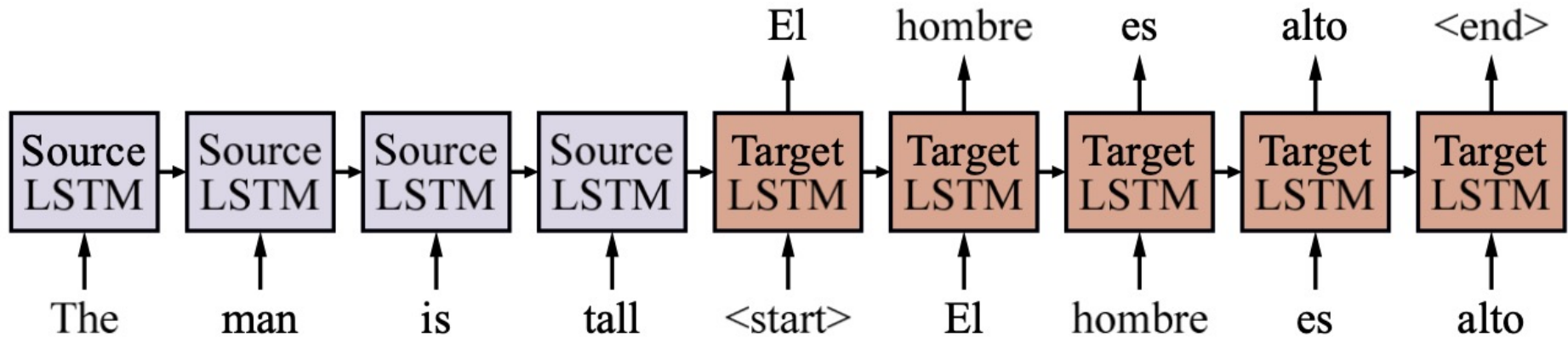Semantic Textual Similarity

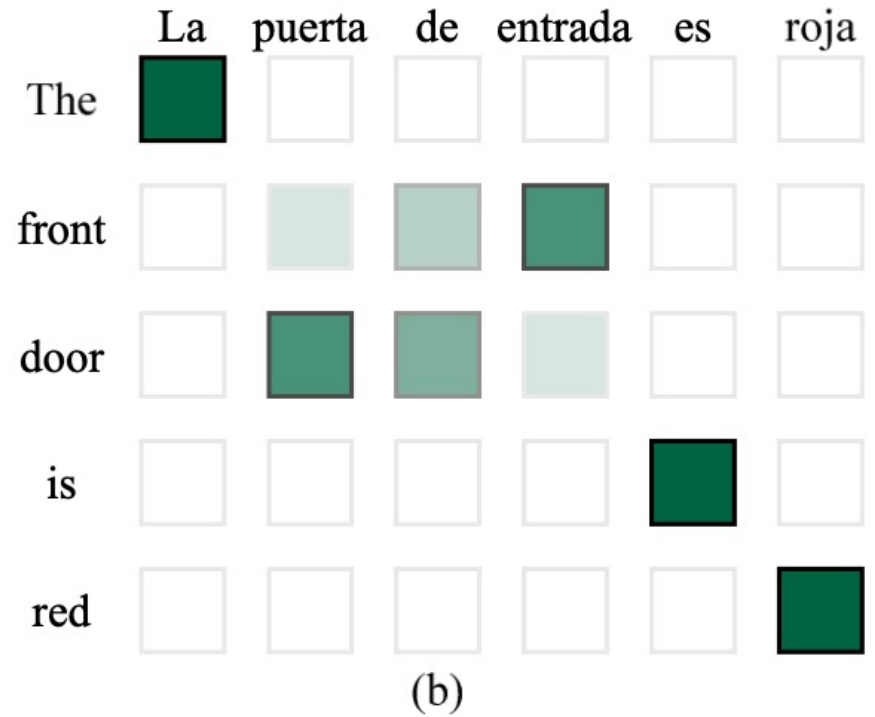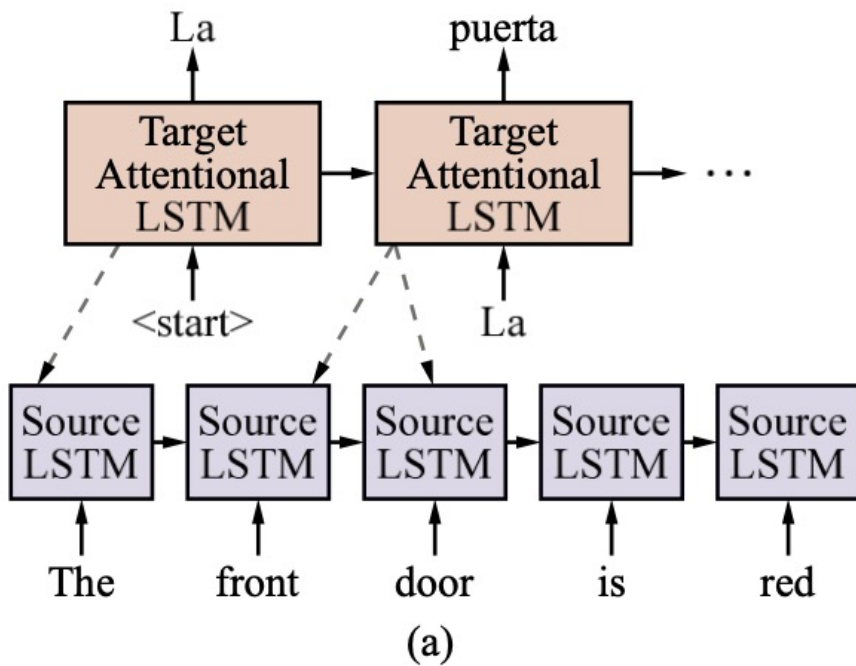https://tinyurl.com/aintpupython101

# Outline

- Word Embeddings
- <span style="color:red">Recurrent Neural Networks for NLP</span>
- Sequence-to-Sequence Models
- The Transformer Architecture
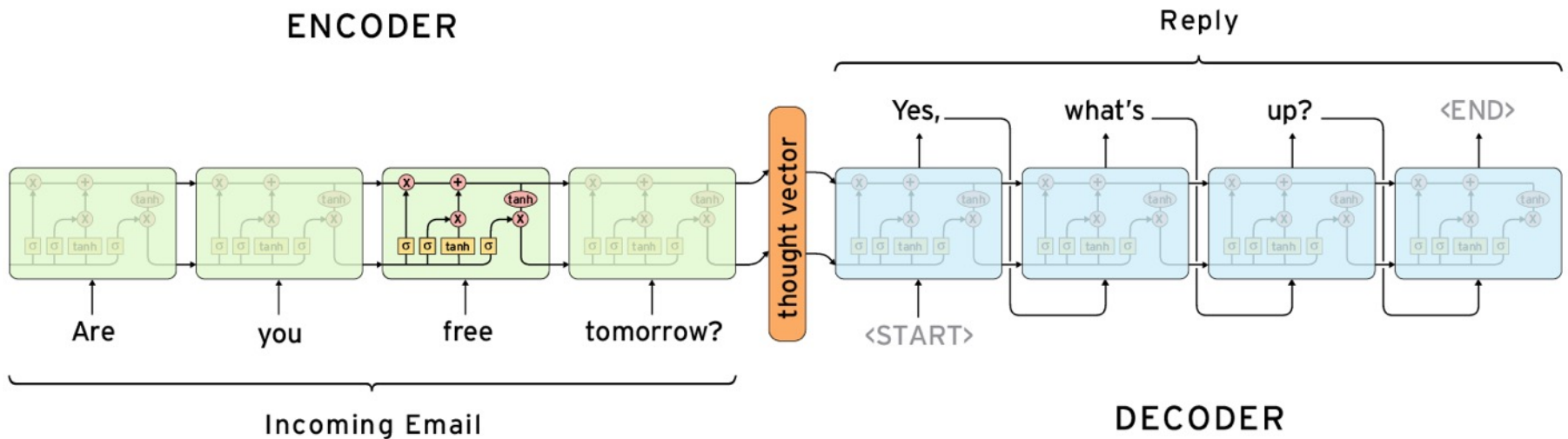- Pretraining and Transfer Learning
- State of the art (SOTA)

# RNN



(a)

(b)

# Bidirectional RNN network for POS tagging

# LSTM Recurrent Neural Network

# Outline

- Word Embeddings
- Recurrent Neural Networks for NLP
- <span style="color:red">Sequence-to-Sequence Models</span>
- The Transformer Architecture
- Pretraining and Transfer Learning
- State of the art (SOTA)

# Sequence-to-Sequence model

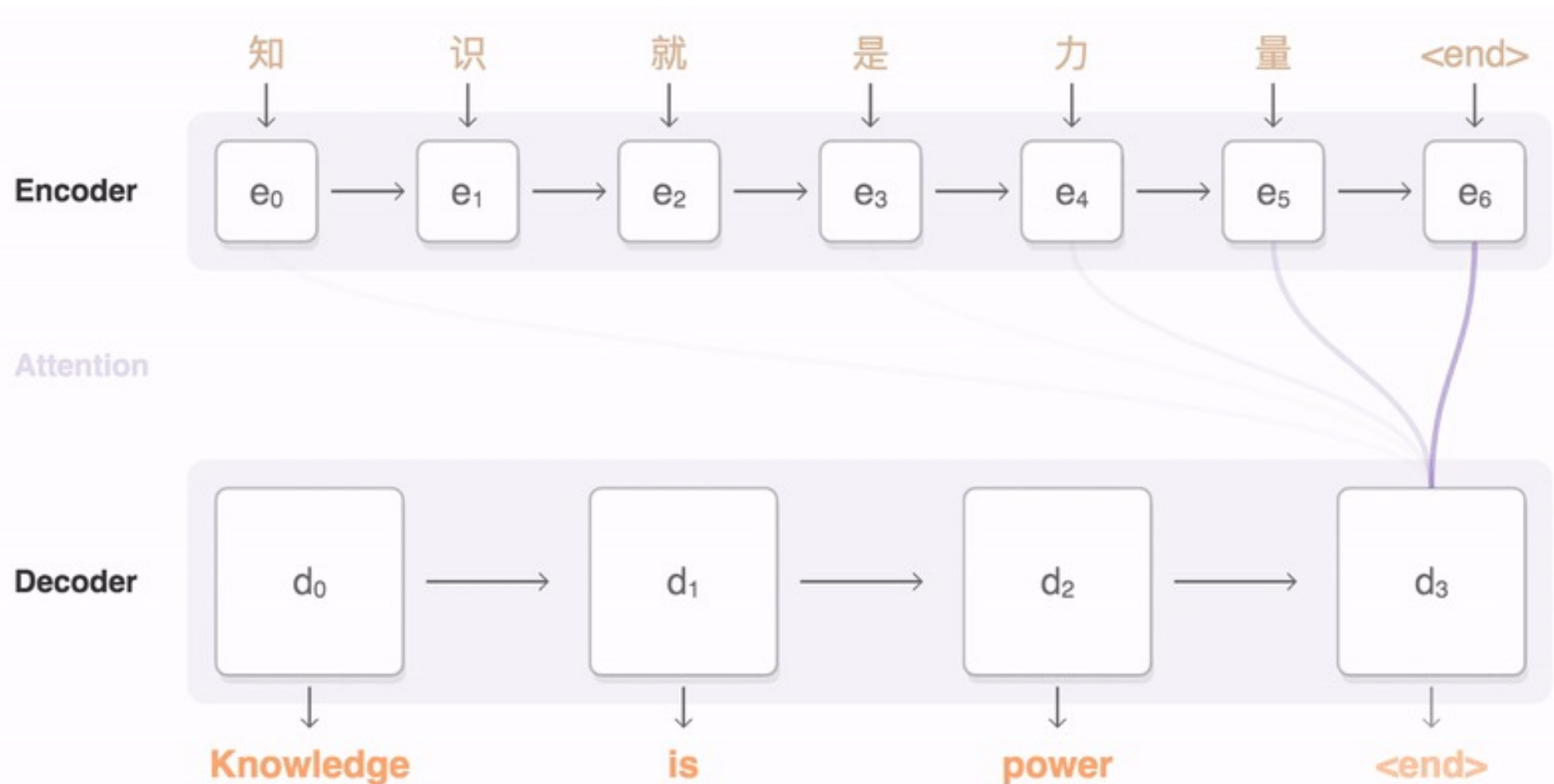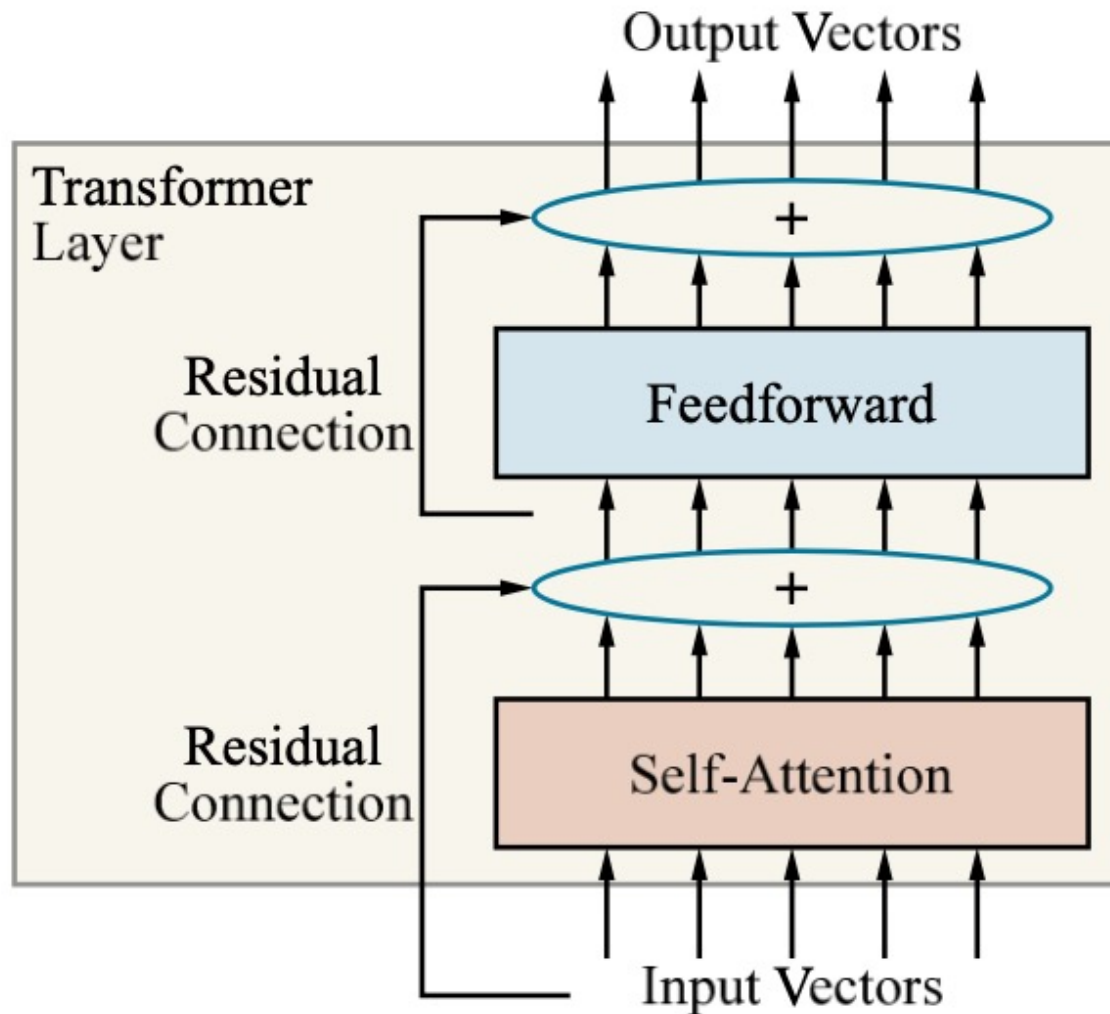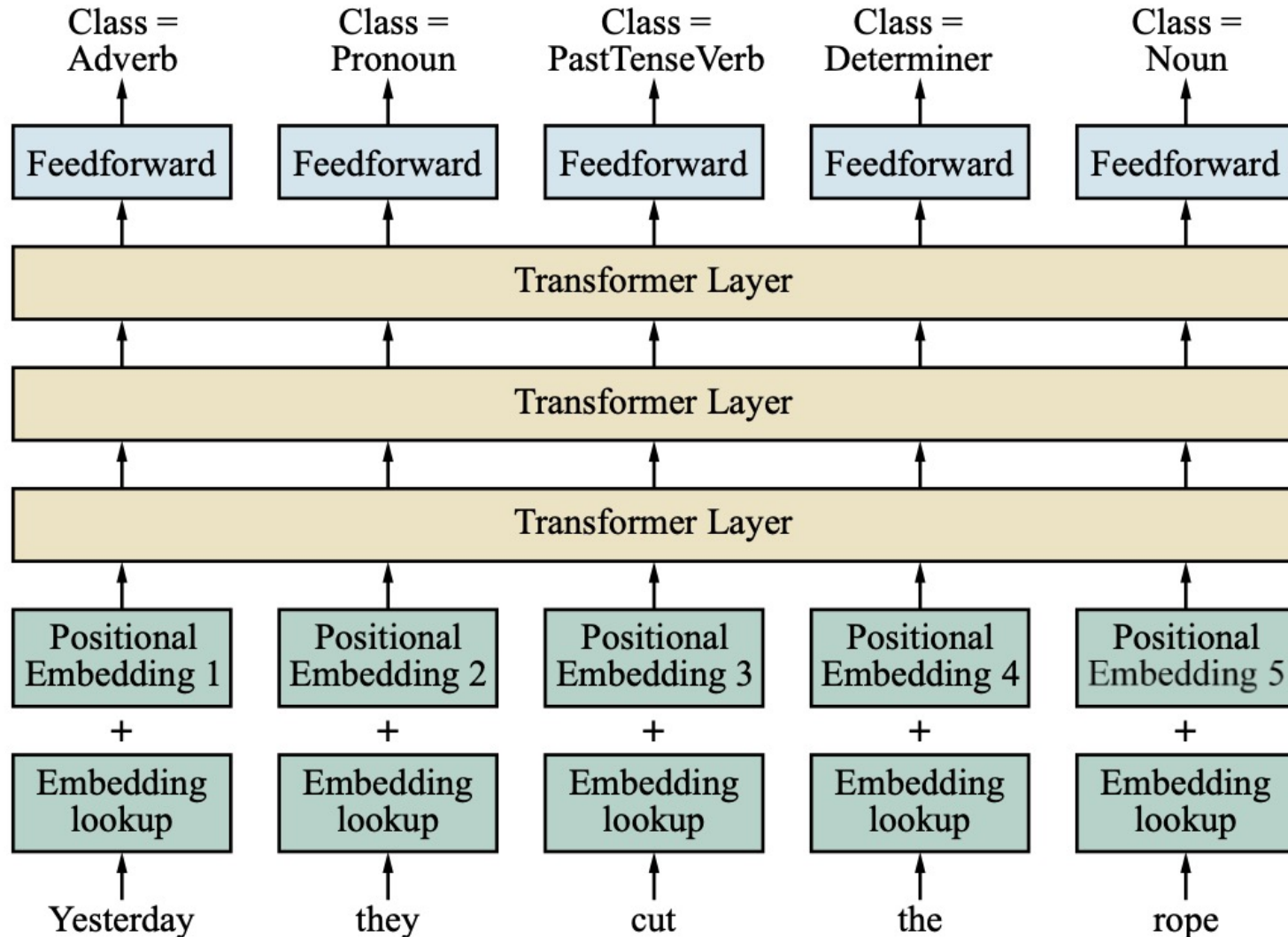# Attentional Sequence-to-Sequence model
# for English-to-Spanish translation



(a)

(b)

# The Sequence to Sequence model (seq2seq)

45

# Sequence to Sequence (Seq2Seq)

# Outline

- Word Embeddings
- Recurrent Neural Networks for NLP
- Sequence-to-Sequence Models
- The Transformer Architecture
- Pretraining and Transfer Learning
- State of the art (SOTA)

# Single-layer Transformer
## consists of self-attention,
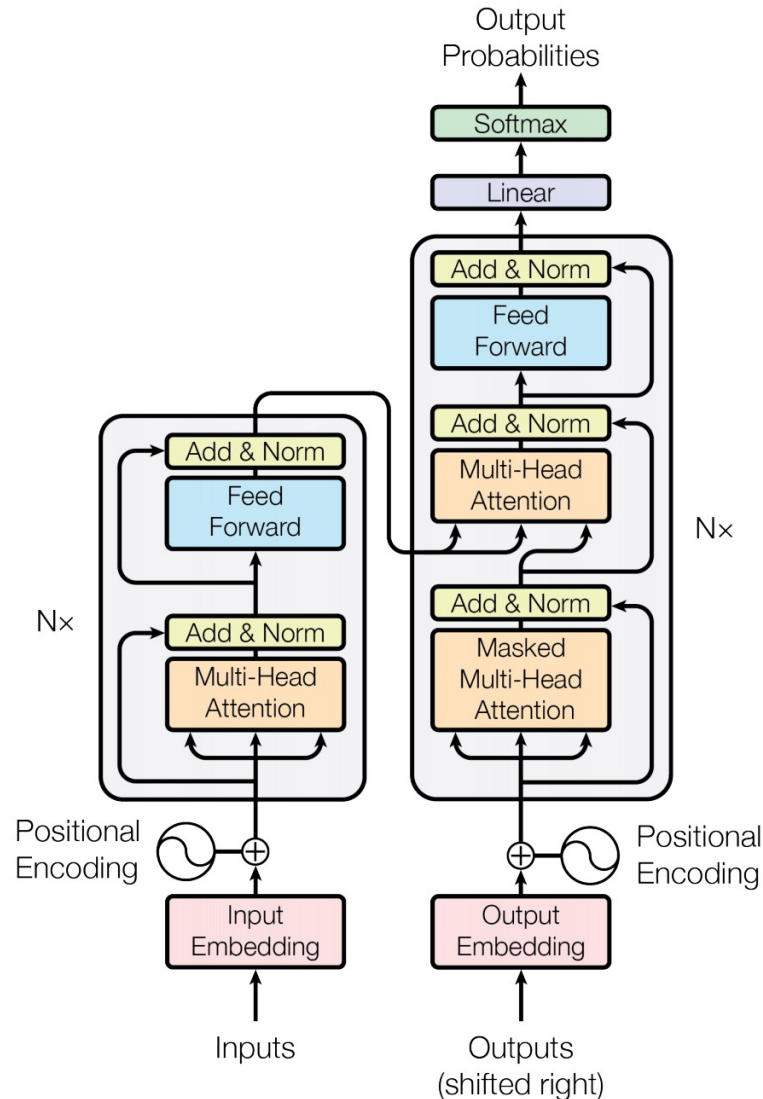## a feedforward network, and residual connection

# Transformer Architecture for POS Tagging
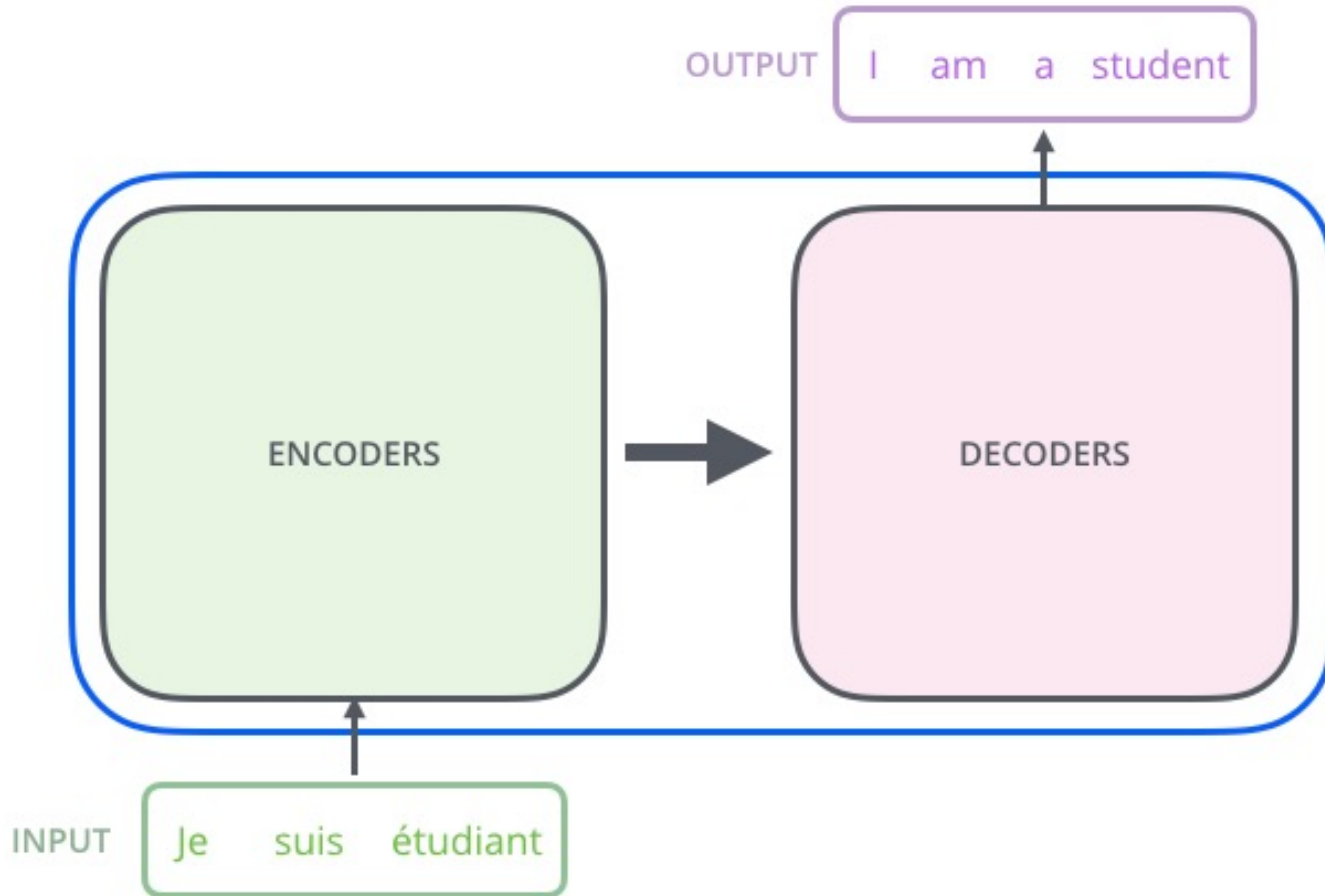
# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)

# Transformer



INPUT

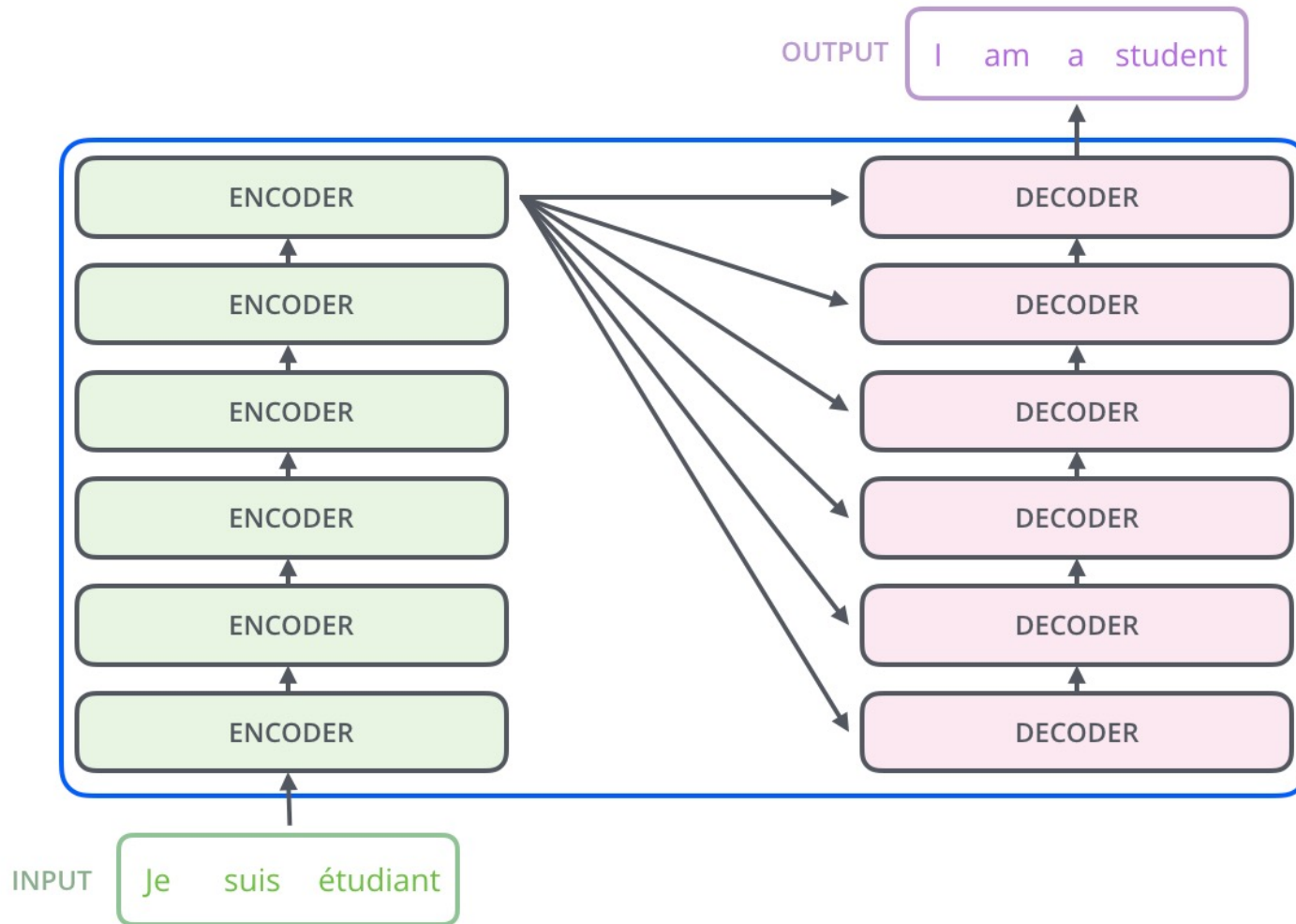Je suis étudiant

THE TRANSFORMER

OUTPUT

I am a student
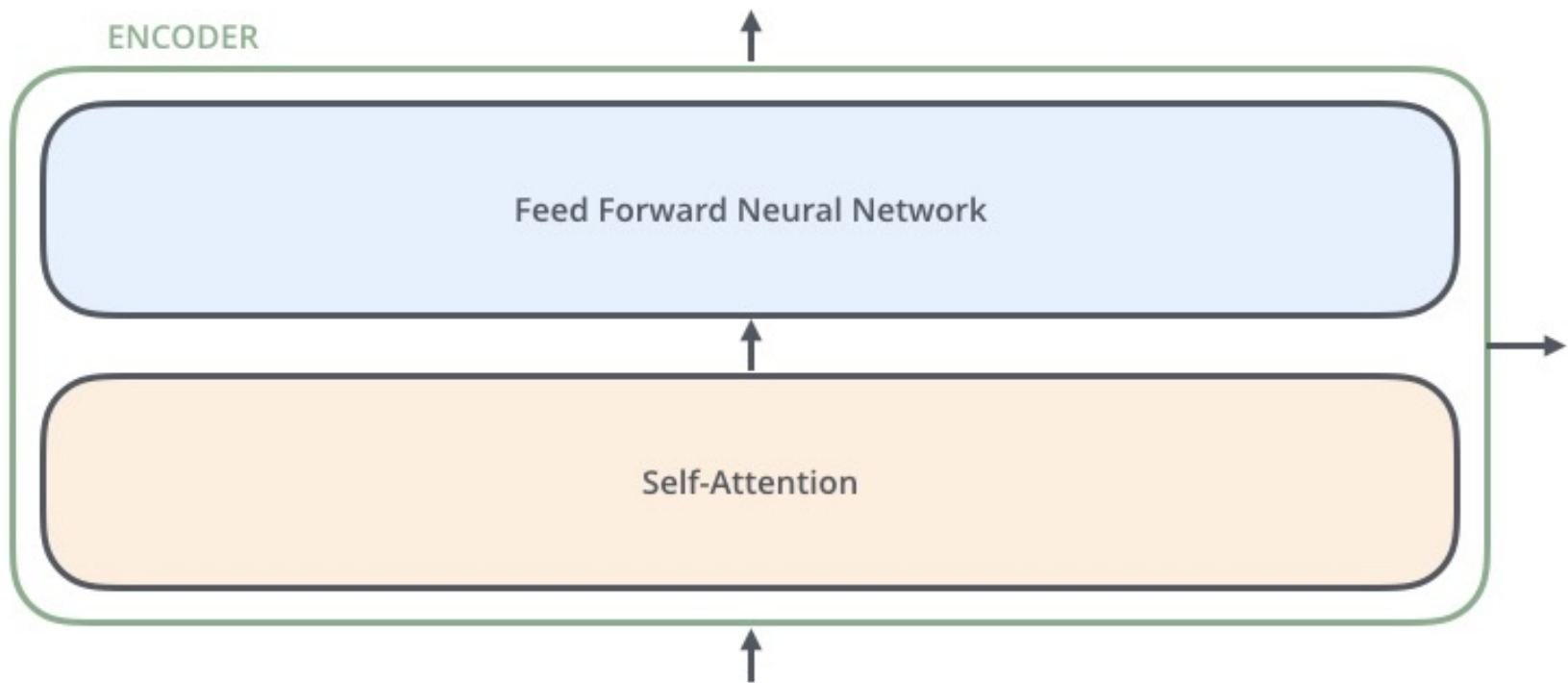
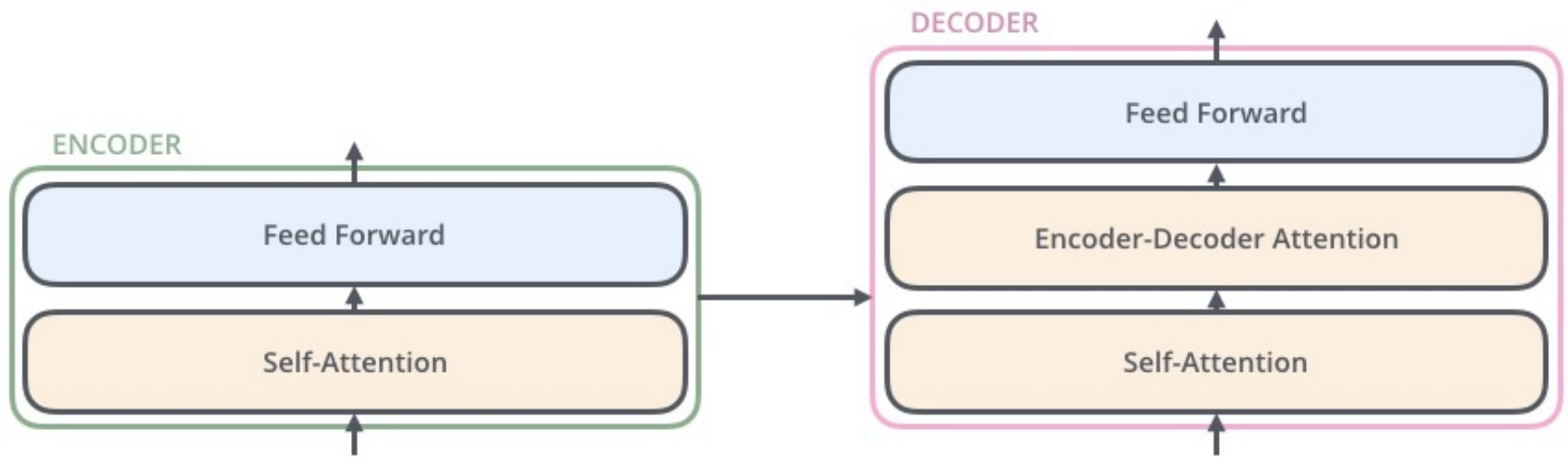# Transformer Encoder Decoder

# Transformer
# Encoder Decoder Stack

# Transformer
# Encoder Self-Attention



54
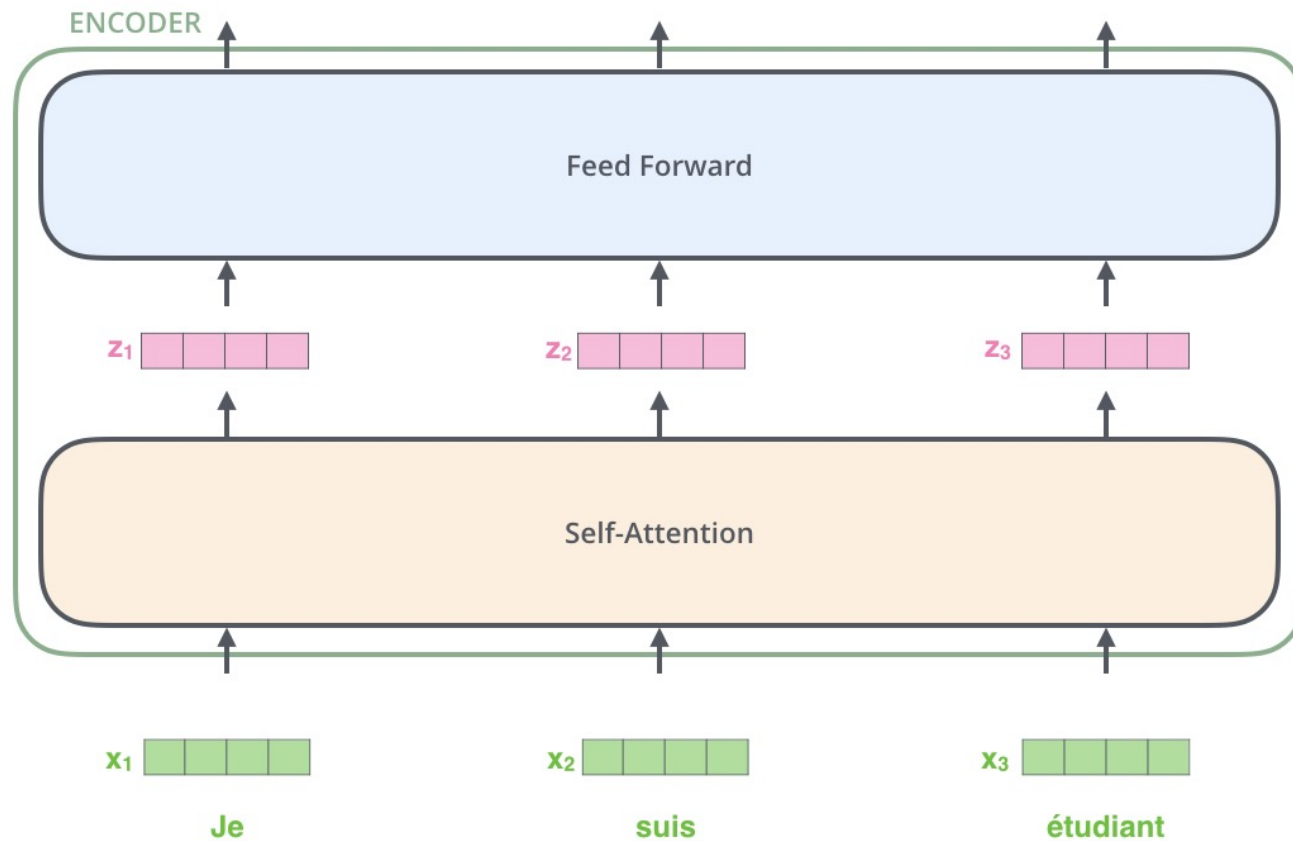
# Transformer Decoder

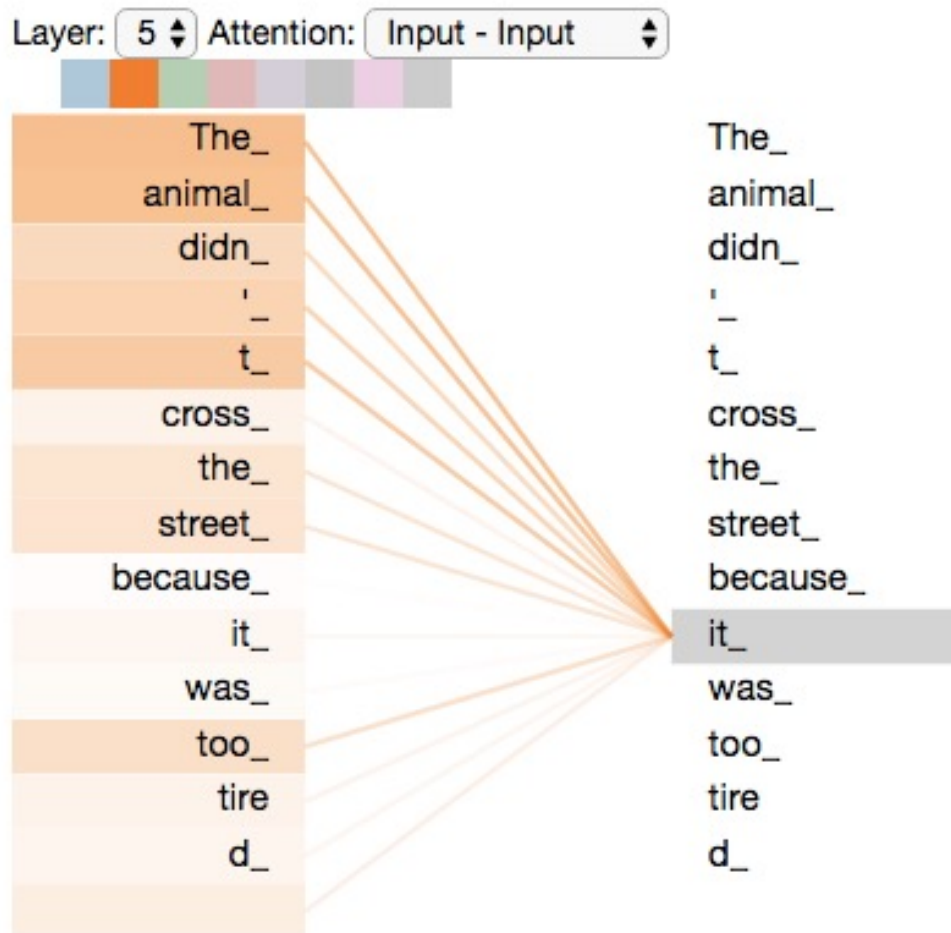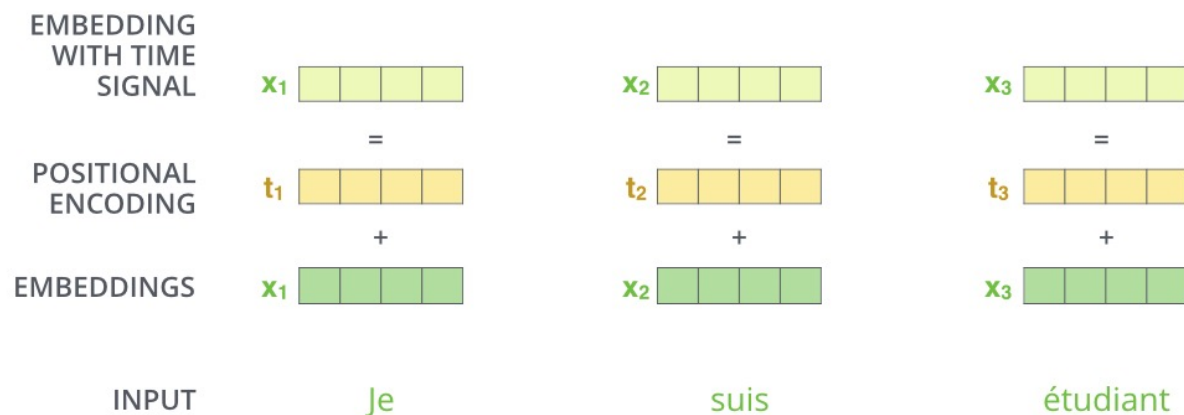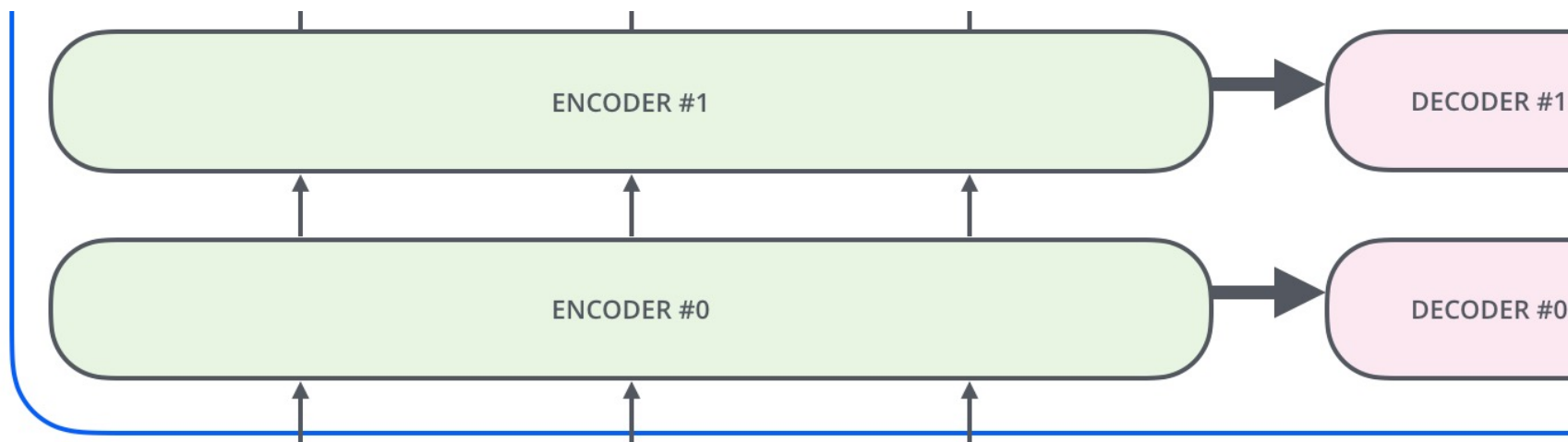# Transformer Encoder with Tensors Word Embeddings

# Transformer
# Self-Attention Visualization

# Transformer
# Positional Encoding Vectors

# Transformer
# Self-Attention Softmax Output



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| 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 | $z_1$ | $z_2$ |

# Outline

- Word Embeddings

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- <span style="color:red">Pretraining and Transfer Learning</span>

- State of the art (SOTA)

# BERT:
# Pre-training of Deep Bidirectional Transformers for Language Understanding
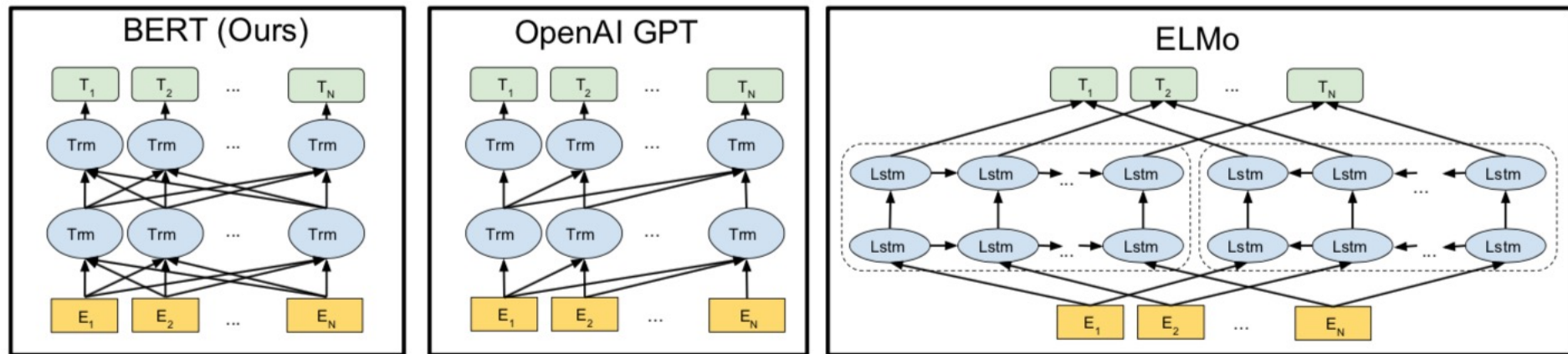
**Jacob Devlin**    **Ming-Wei Chang**    **Kenton Lee**    **Kristina Toutanova**

Google AI Language

{jacobdevlin,mingweichang,kentonl,kristout}@google.com

# BERT

## Bidirectional Encoder Representations from Transformers



## Pre-training model architectures

**BERT** uses a bidirectional Transformer.
**OpenAI GPT** uses a left-to-right Transformer.
**ELMo** uses the concatenation of independently trained left-to-right and right- to-left LSTM to generate features for downstream tasks.
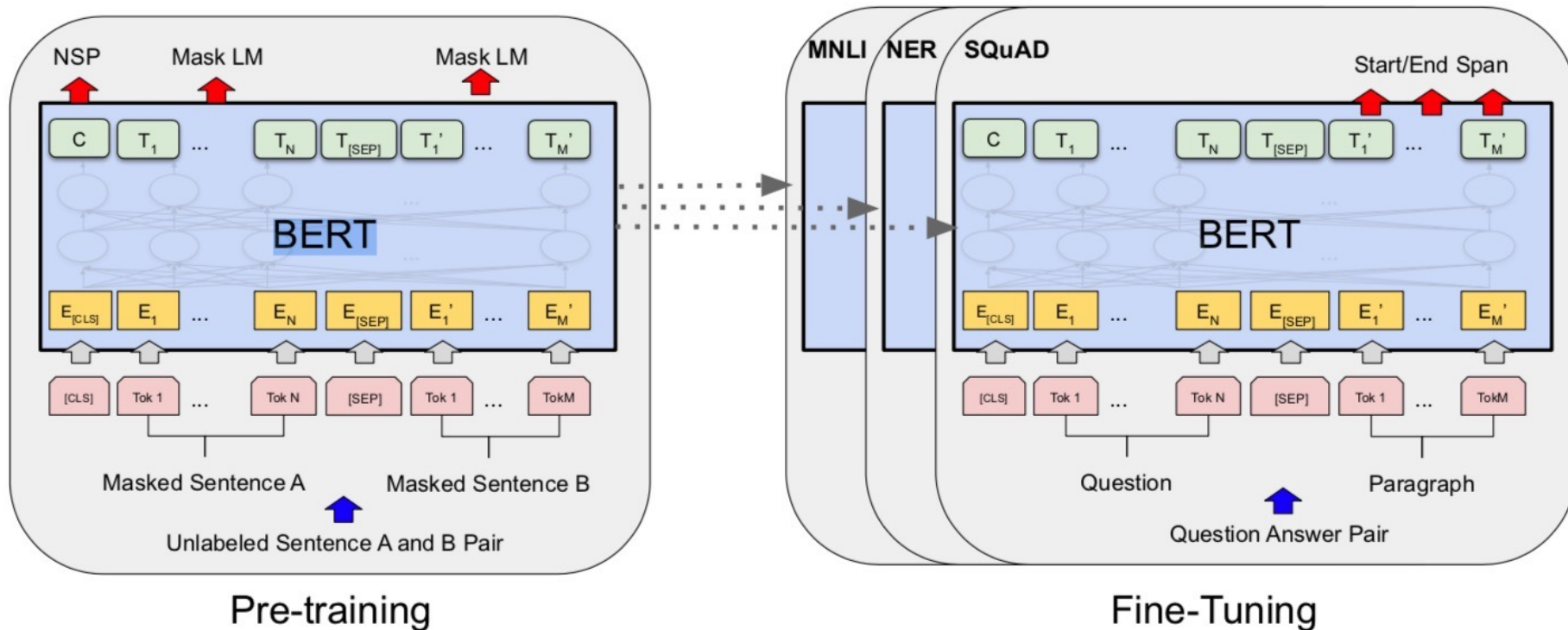Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT
## (Bidirectional Encoder Representations from Transformers)

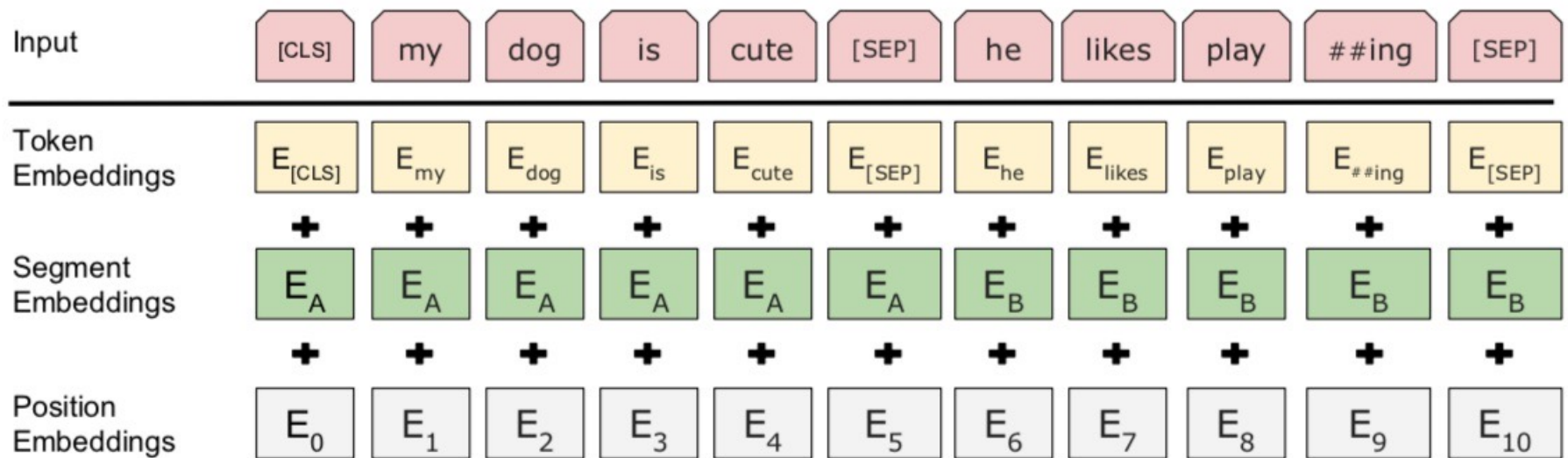## Overall pre-training and fine-tuning procedures for BERT



Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
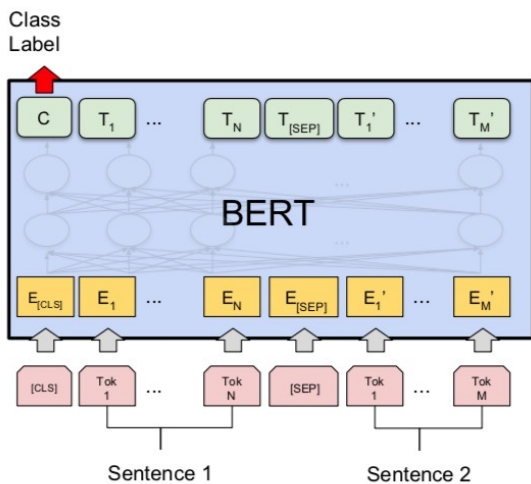
BERT (Bidirectional Encoder Representations from Transformers)
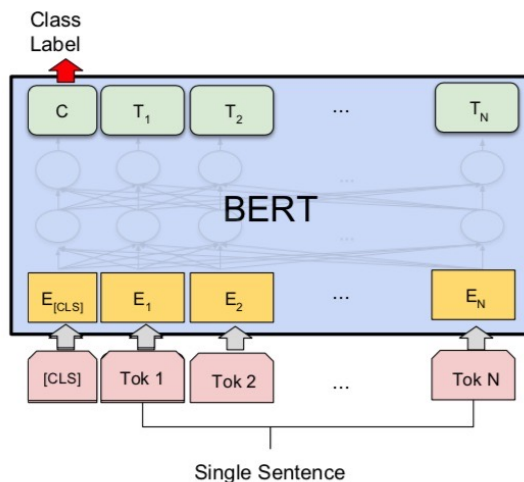
## BERT input representation



The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.
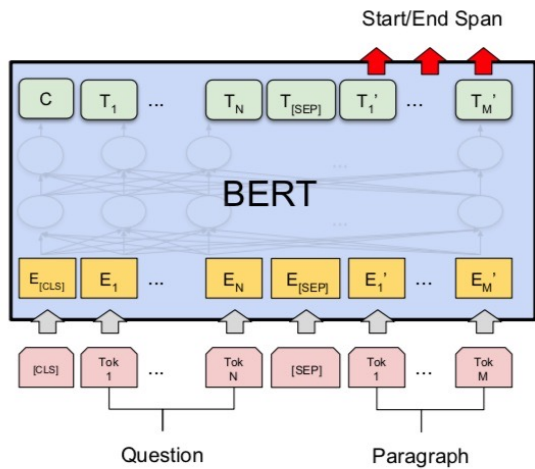
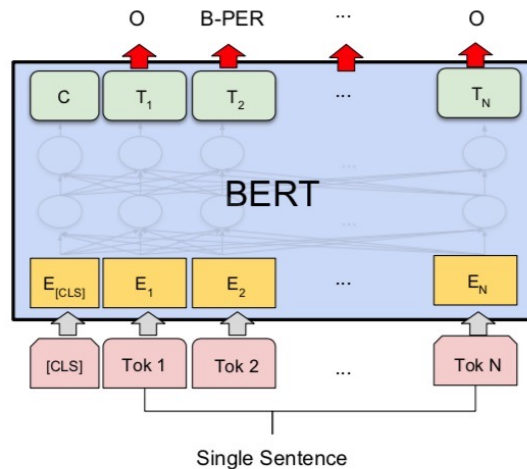# Fine-tuning BERT on Different Tasks



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

(b) Single Sentence Classification Tasks:
SST-2, CoLA
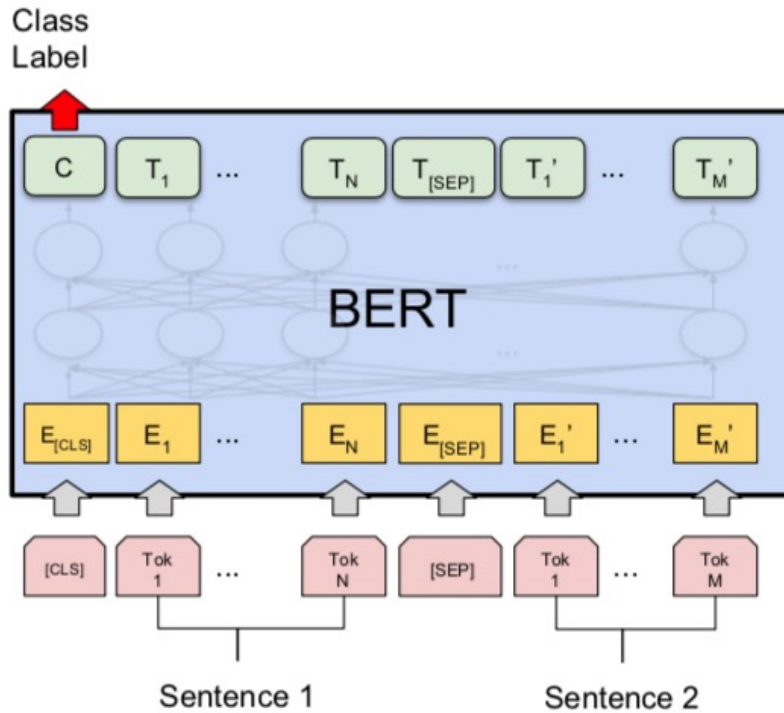
(c) Question Answering Tasks:
SQuAD v1.1

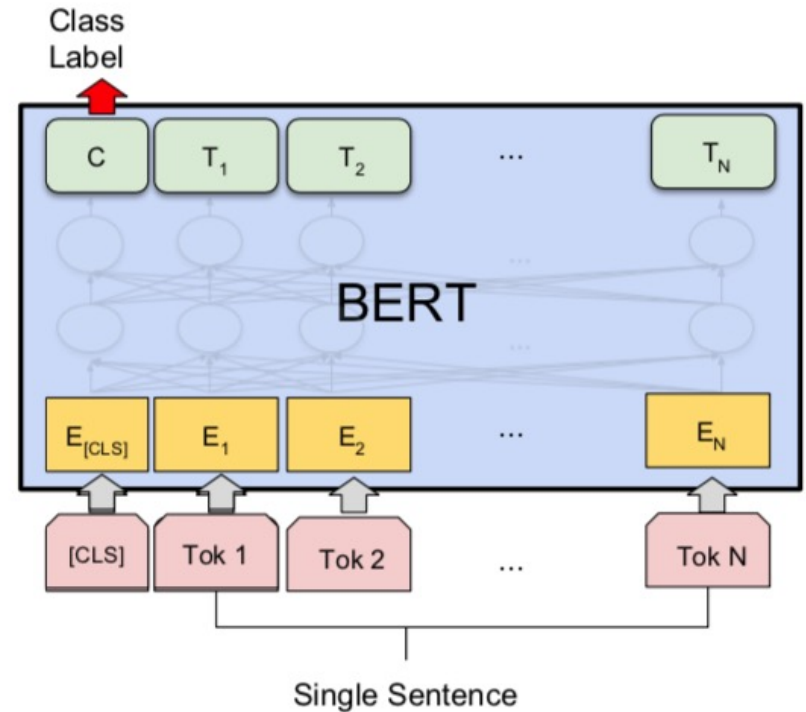(d) Single Sentence Tagging Tasks:
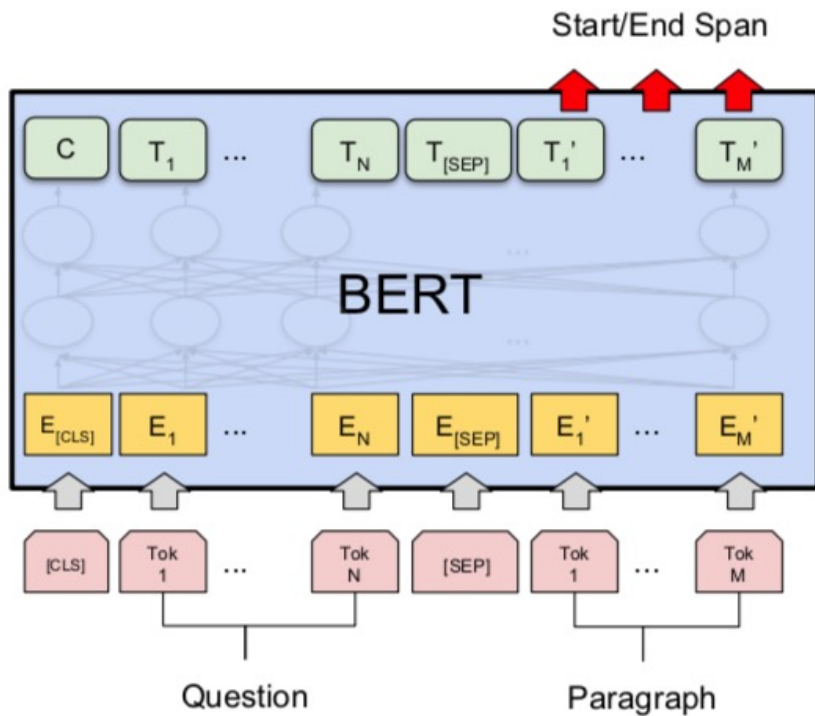CoNLL-2003 NER

65

# BERT Sequence-level tasks



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

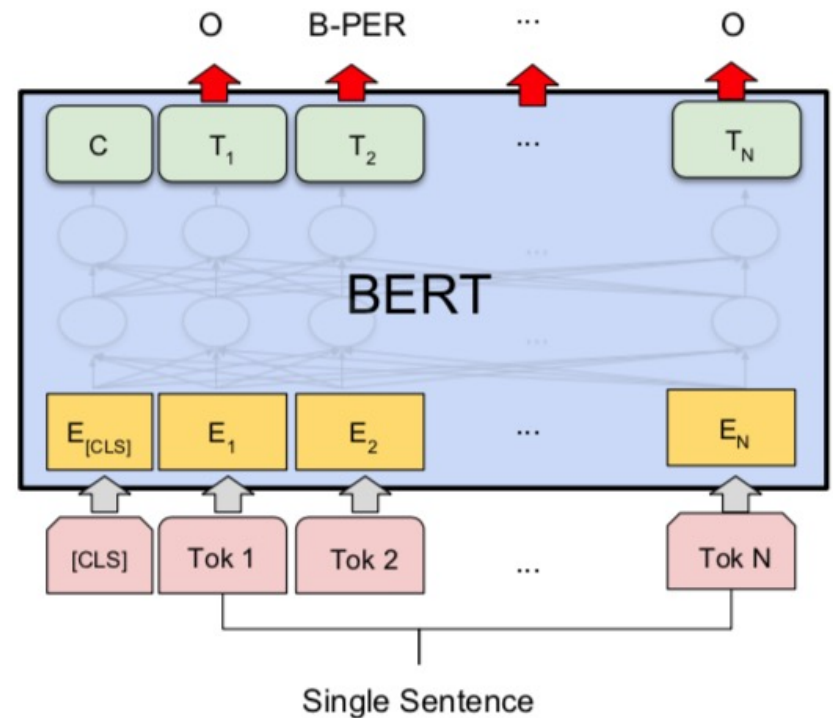(b) Single Sentence Classification Tasks:
SST-2, CoLA

# BERT Token-level tasks



(c) Question Answering Tasks: SQuAD v1.1

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

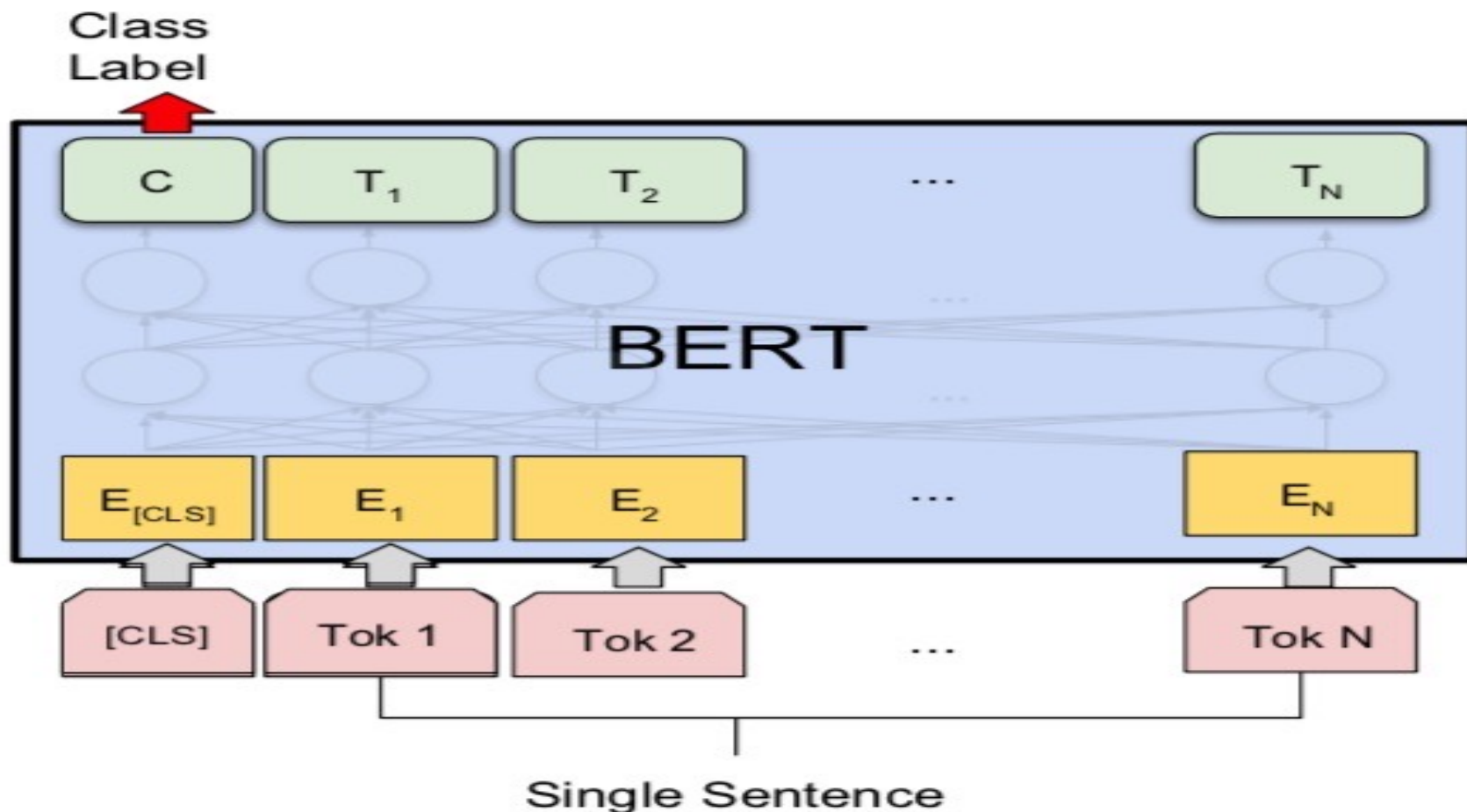# Fine-tuning BERT on Question Answering (QA)



Start/End Span

BERT

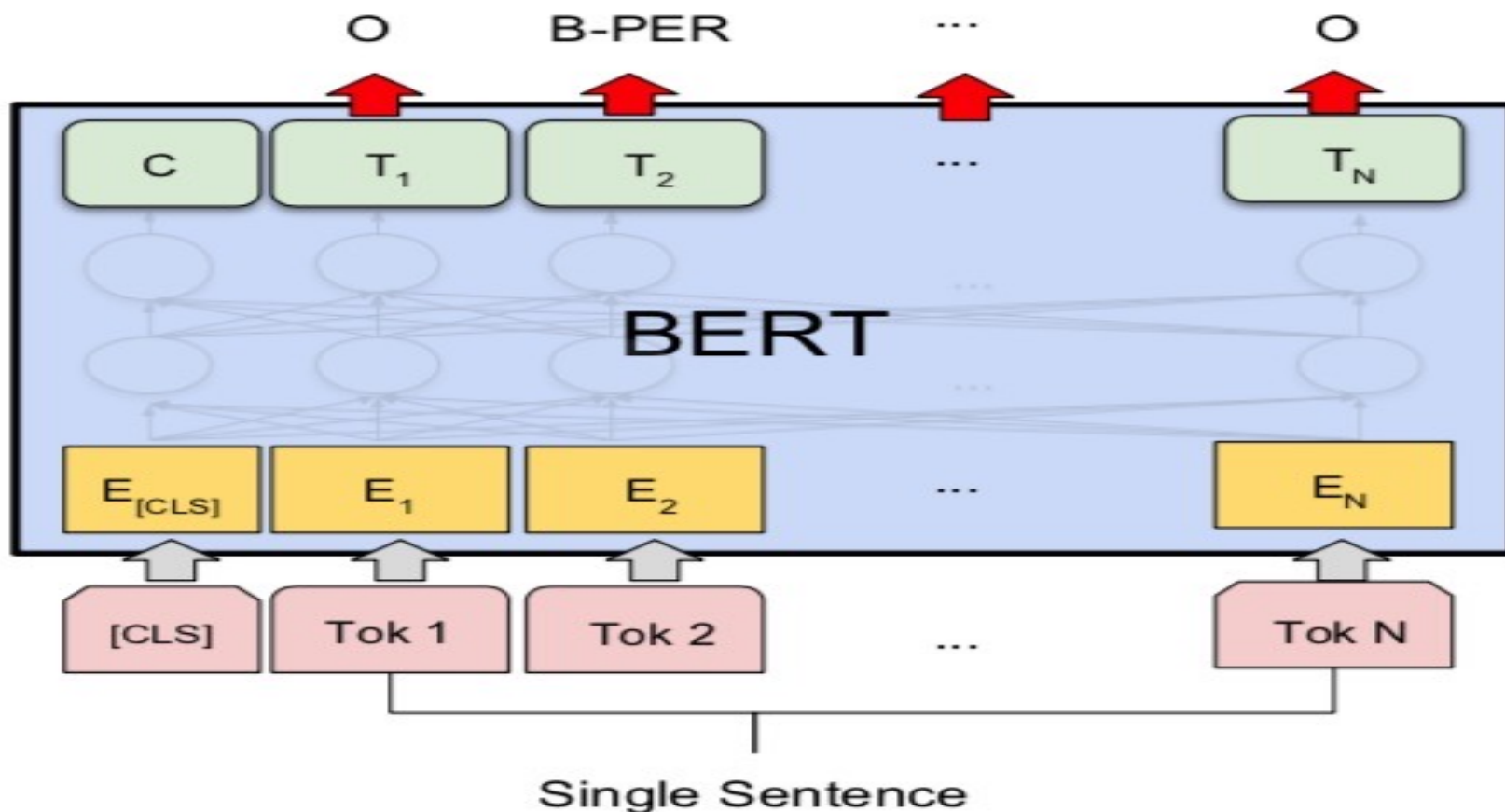(c) Question Answering Tasks: SQuAD v1.1

# Fine-tuning BERT on Dialogue
# Intent Detection (ID; Classification)



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

# Fine-tuning BERT on Dialogue Slot Filling (SF)



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

# General Language Understanding Evaluation (GLUE) benchmark

## GLUE Test results

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| 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 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| $BERT_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| $BERT_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

**MNLI**: Multi-Genre Natural Language Inference
**QQP**: Quora Question Pairs
**QNLI**: Question Natural Language Inference
**SST-2**: The Stanford Sentiment Treebank
**CoLA**: The Corpus of Linguistic Acceptability
**STS-B**:The Semantic Textual Similarity Benchmark
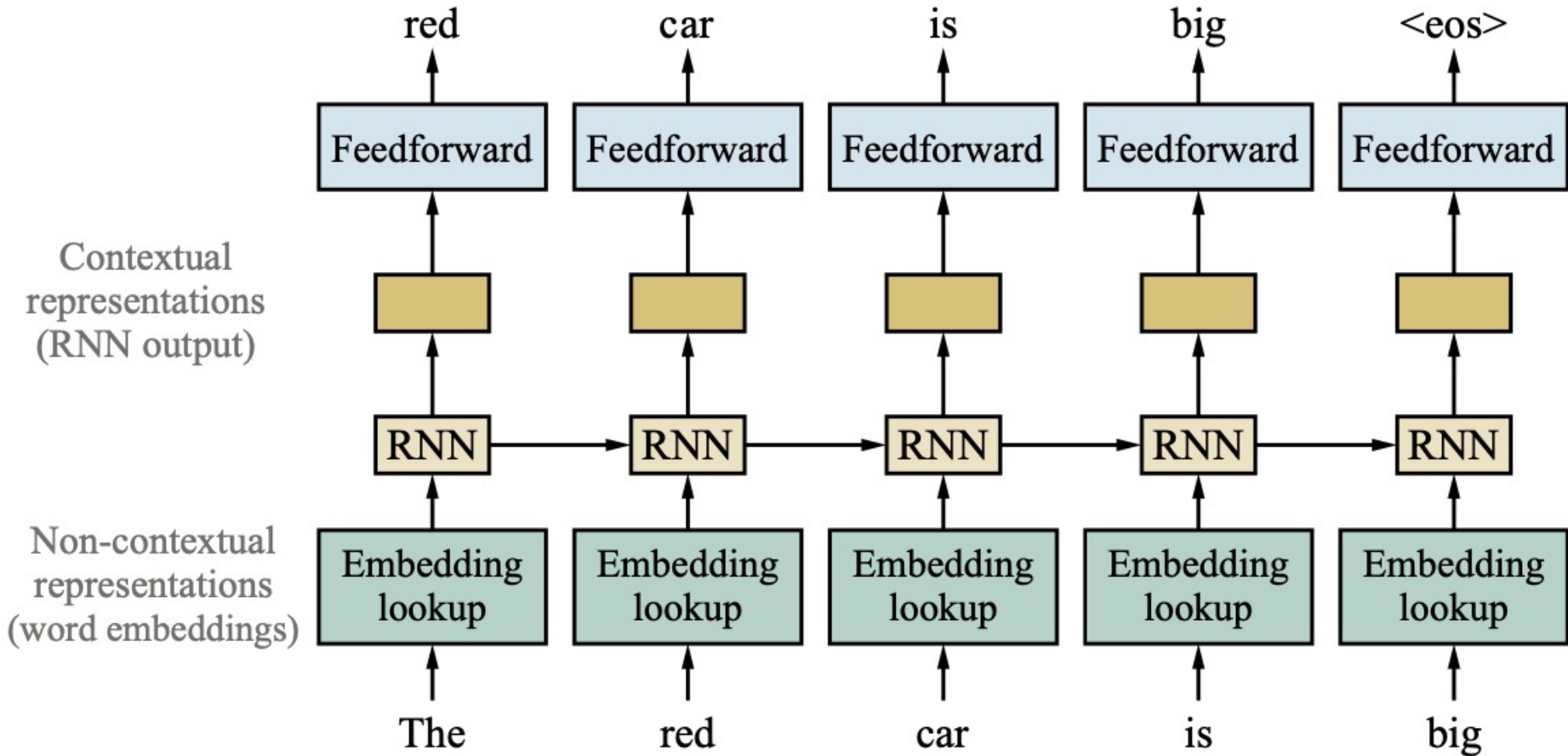**MRPC**: Microsoft Research Paraphrase Corpus
**RTE**: Recognizing Textual Entailment

# Training Contextual Representations
## using a left-to-right Language Model

# Masked Language Modeling:
## Pretrain a Bidirectional Model

# Illustrated BERT

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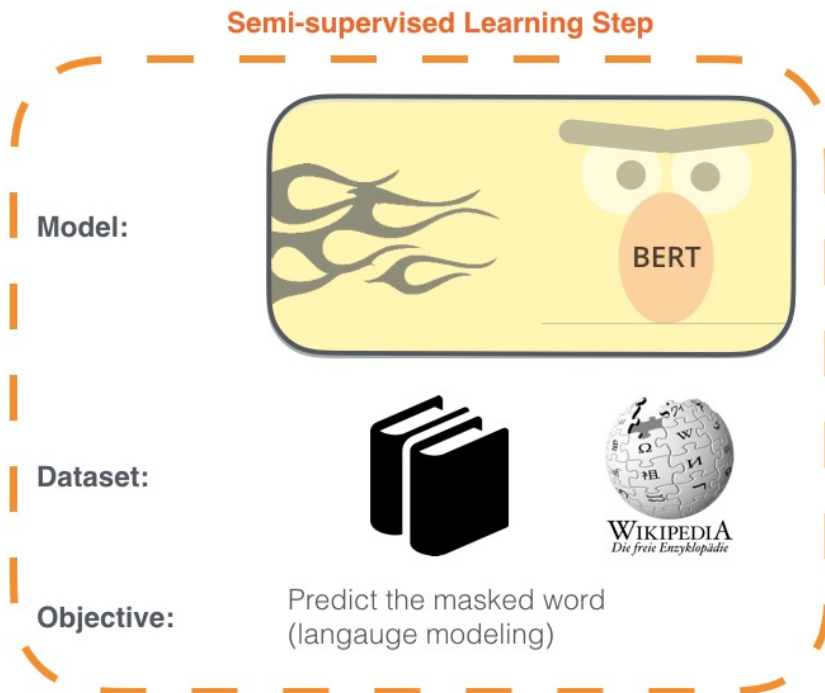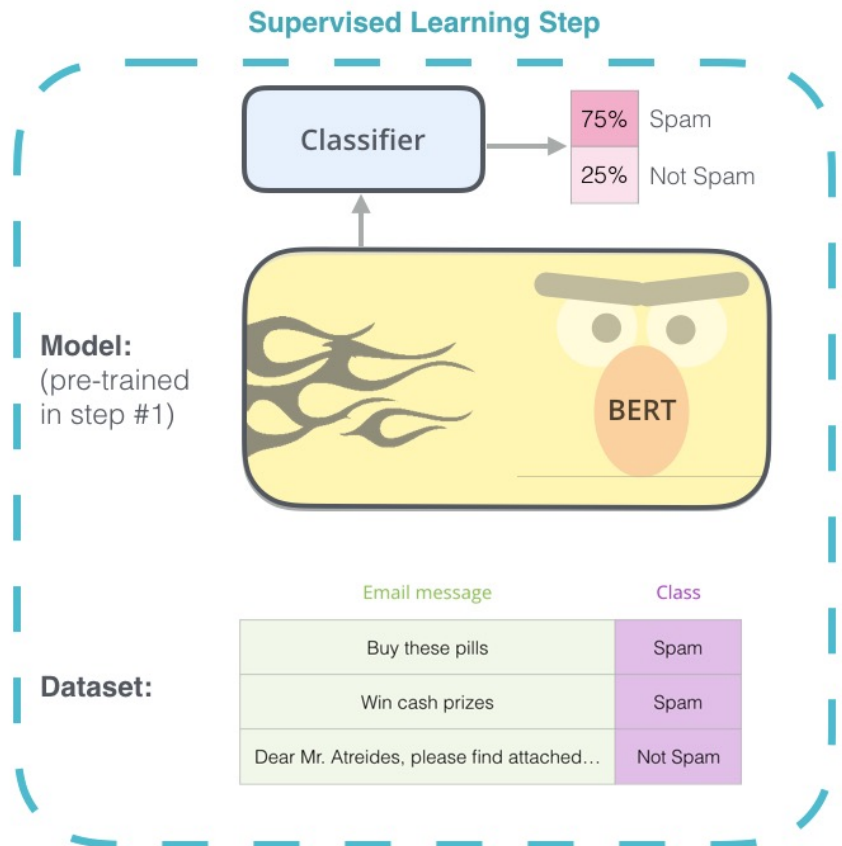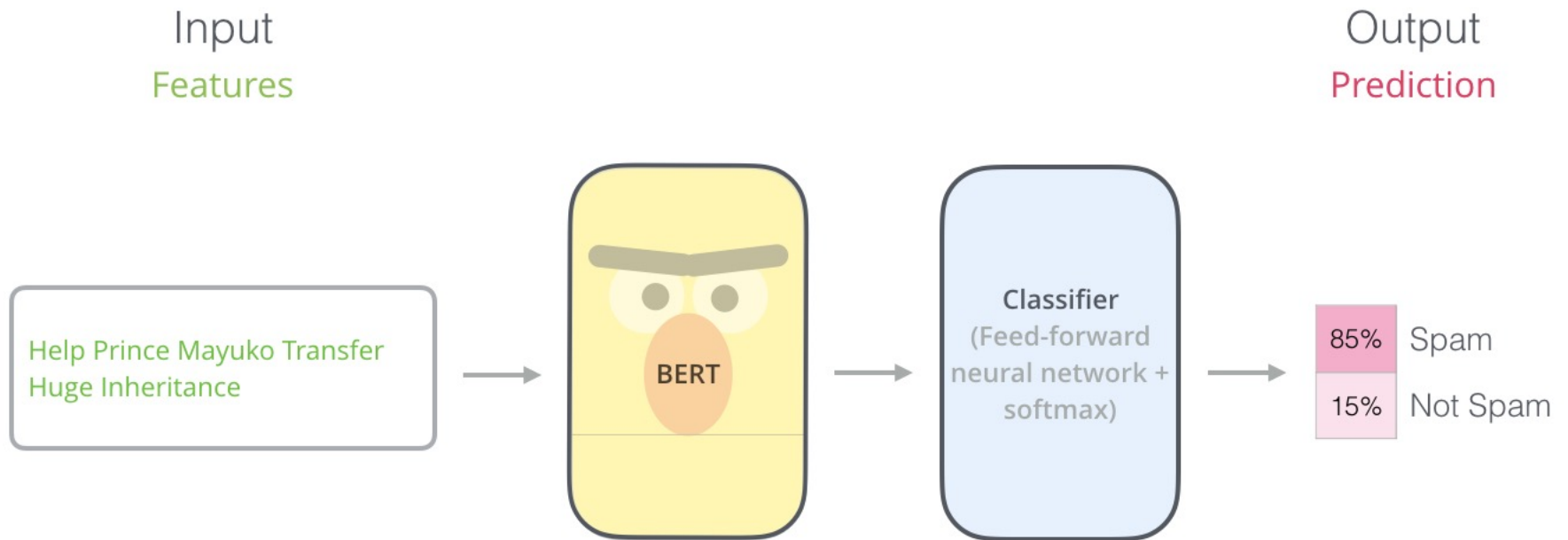 has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

**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 |

# BERT Classification Input Output



Input
Features

Output
Prediction

Help Prince Mayuko Transfer Huge Inheritance

BERT

Classifier
(Feed-forward neural network + softmax)

85% Spam
15% Not Spam

# BERT Encoder Input



BERT

# BERT Classifier



85% Spam

15% Not Spam

**Classifier**
(Feed-forward neural network + softmax)

1    2    3    4    •••    512

BERT

1    2    3    4    •••    512

[CLS]   Help   Prince   Mayuko

# Sentiment Analysis: Single Sentence Classification



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

# A Visual Guide to
# Using BERT for the First Time
## (Jay Alammar, 2019)

# Sentiment Classification: SST2
# Sentences from movie reviews

| sentence | label |
|---|---|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

# Movie Review Sentiment Classifier

# Movie Review Sentiment Classifier

# Movie Review Sentiment Classifier Model Training

# Step # 1 Use distilBERT to Generate Sentence Embeddings

# Step #2:Test/Train Split for Model #2, Logistic Regression

# Step #3 Train the logistic regression model using the training set



Step #3: Train the logistic regression model using the training set

Sentence Embeddings    label

|  | 0 | 1 | ... | 767 |  |
|---|---|---|---|---|---|
| 0 | −0.215 | −0.1402 | ... | 0.201 | 1 |
| 1 | | | | | |
| 2 | | | | | |
| 3 | | | | | |
| 4 | | | | | |
| 5 | | | | | |
| 6 | | | | | |
| ... | | | | | |
| 1,499 | | | | | |

Model Training

Logistic Regression

scikit learn

# Tokenization

[CLS] a visually stunning rum ##ination on love [SEP]
a visually stunning rumination on love



**Tokenization**
`DistilBertTokenizer`

| [CLS] | a | visually | stunning | rum | ##ination | on | love | [SEP] |
|---|---|---|---|---|---|---|---|---|

2) Add [CLS] and [SEP] tokens

| a | visually | stunning | rum | ##ination | on | love |
|---|---|---|---|---|---|---|

1) Break words into tokens
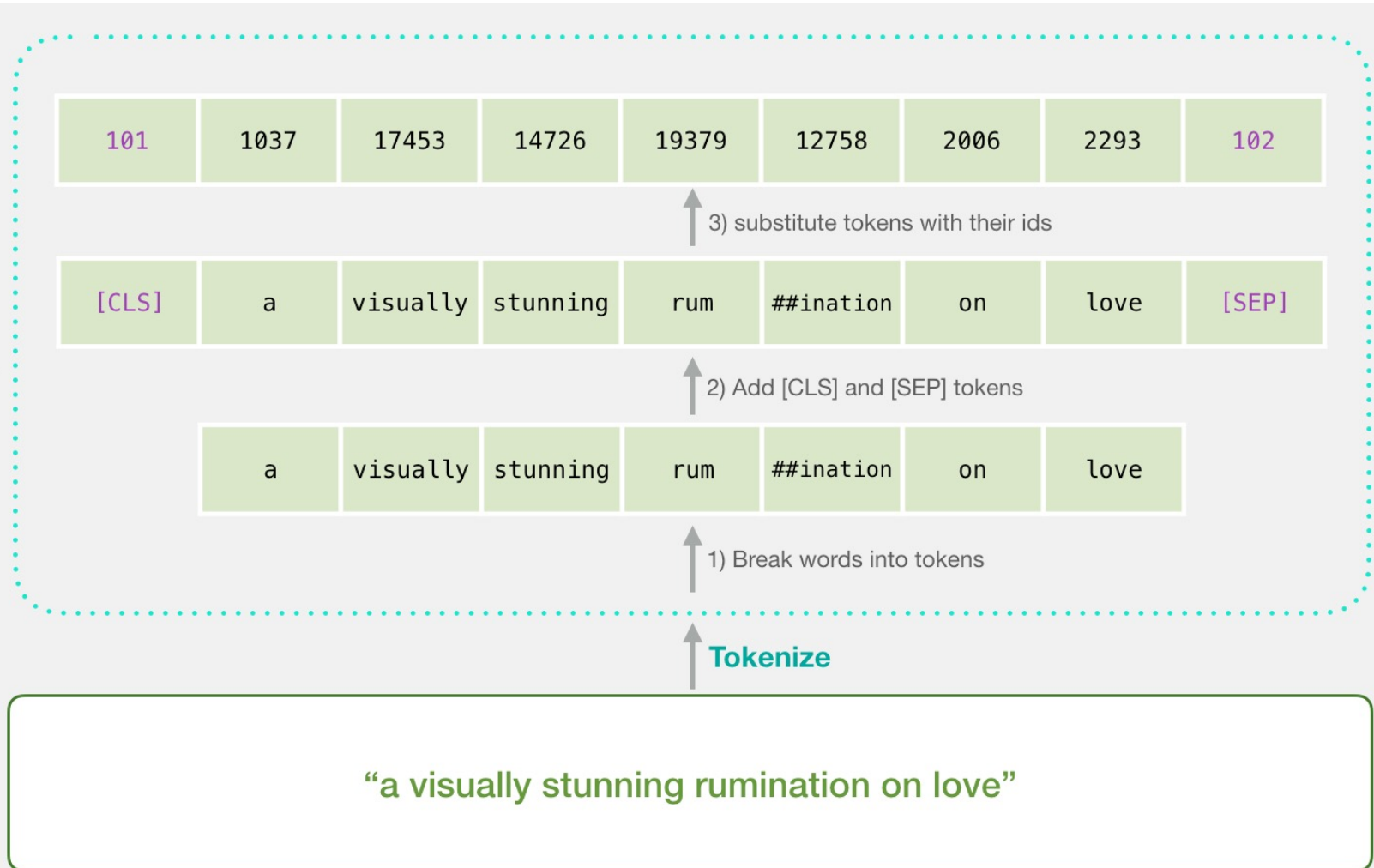
**Tokenize**

"a visually stunning rumination on love"
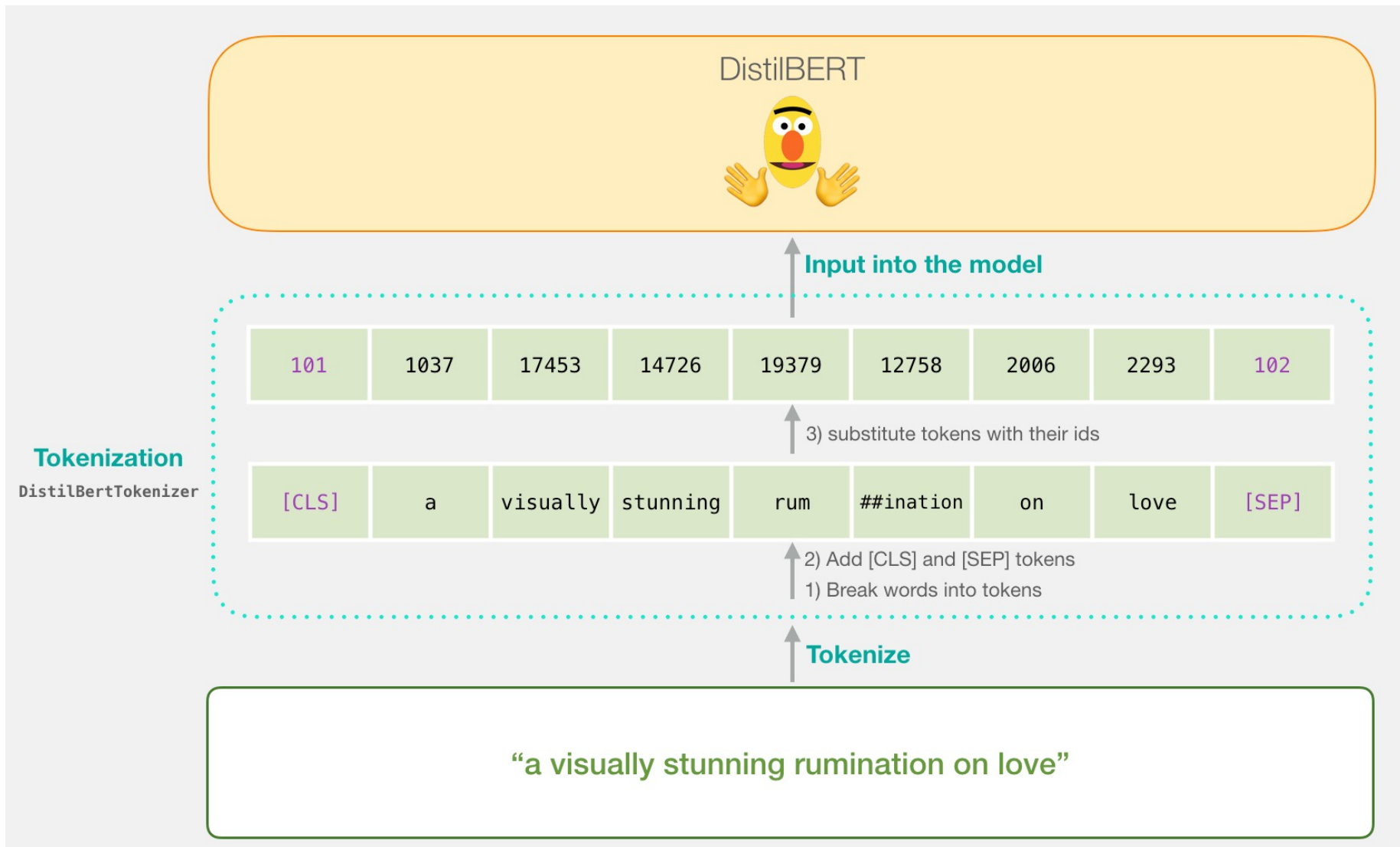
# Tokenization

```
tokenizer.encode("a visually stunning rumination on love",
                 add_special_tokens=True)
```

**Tokenization**

**DistilBertTokenizer**

| 101 | 1037 | 17453 | 14726 | 19379 | 12758 | 2006 | 2293 | 102 |
|---|---|---|---|---|---|---|---|---|

↑ 3) substitute tokens with their ids

| [CLS] | a | visually | stunning | rum | ##ination | on | love | [SEP] |
|---|---|---|---|---|---|---|---|---|

↑ 2) Add [CLS] and [SEP] tokens

| a | visually | stunning | rum | ##ination | on | love |
|---|---|---|---|---|---|---|

↑ 1) Break words into tokens

↑ **Tokenize**

"a visually stunning rumination on love"
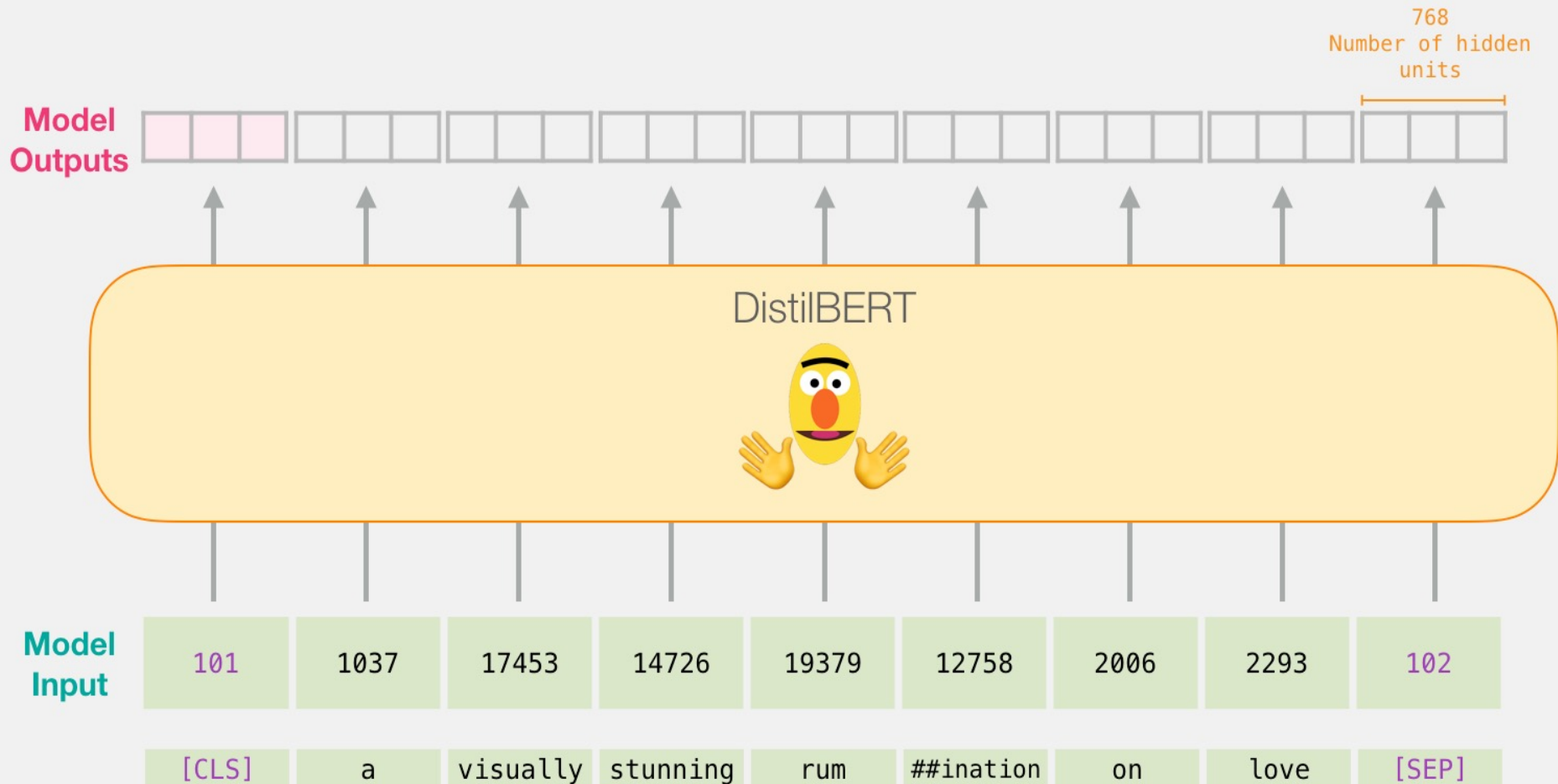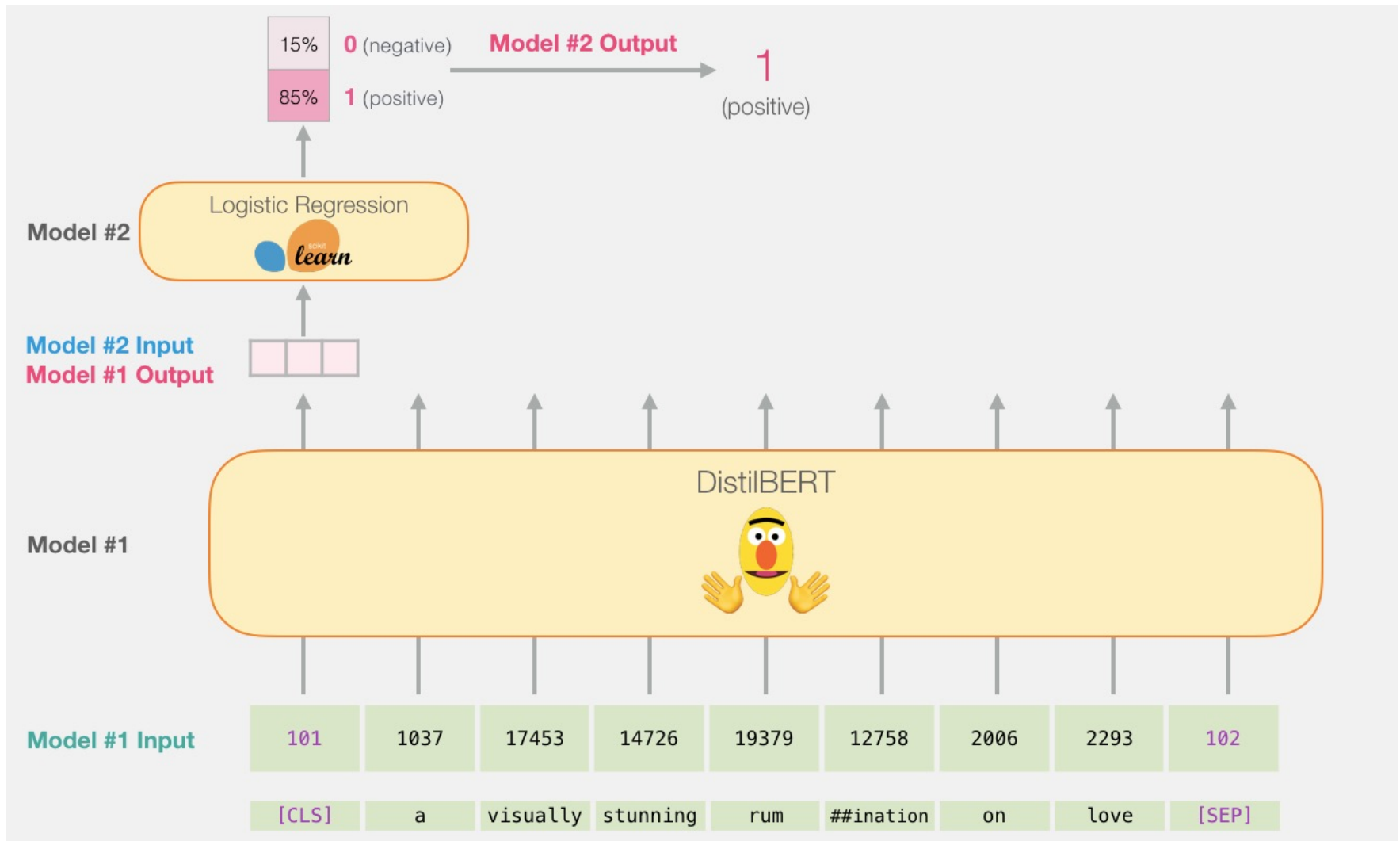
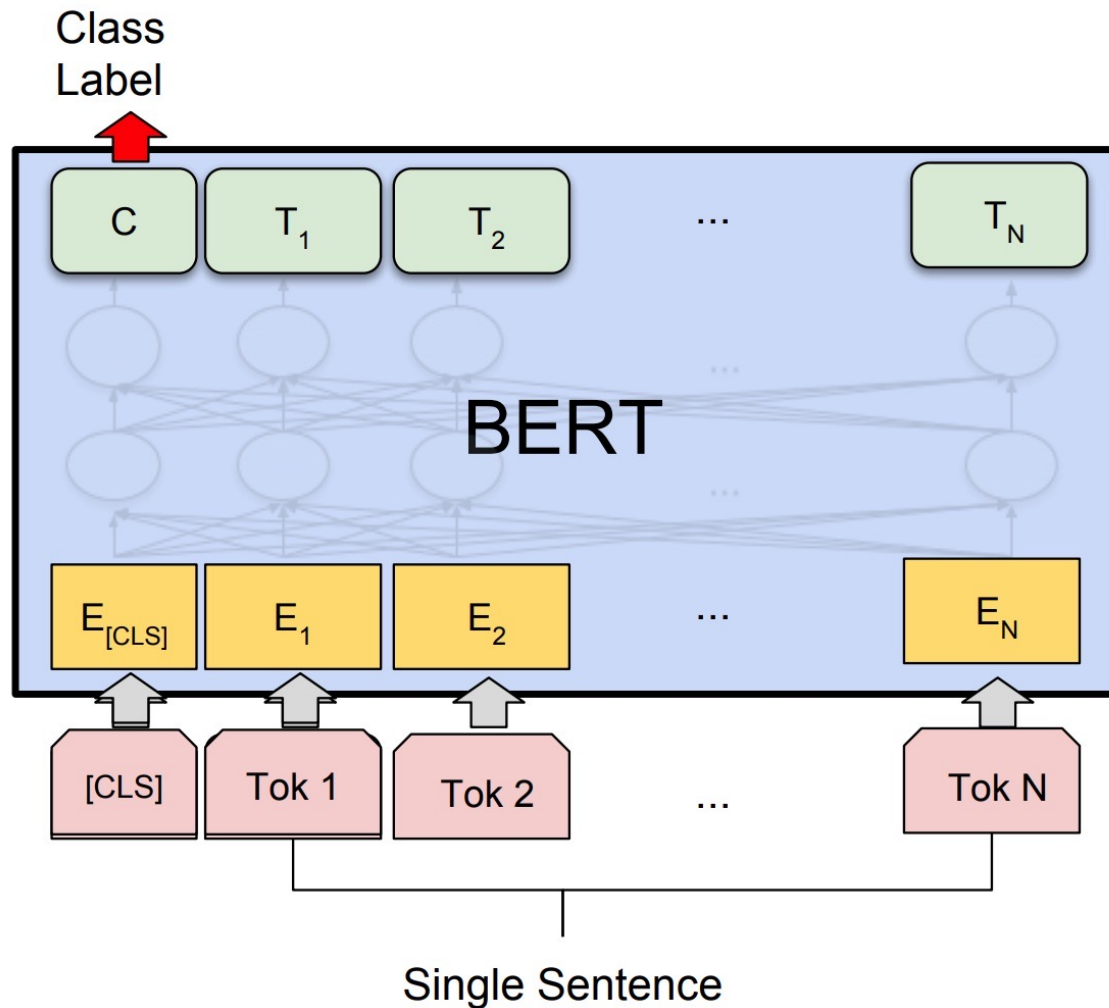# Tokenization for BERT Model

# Flowing Through DistilBERT
# (768 features)

# Model #1 Output Class vector as Model #2 Input



Source: Jay Alammar (2019), A Visual Guide to Using BERT for the First Time,
http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/

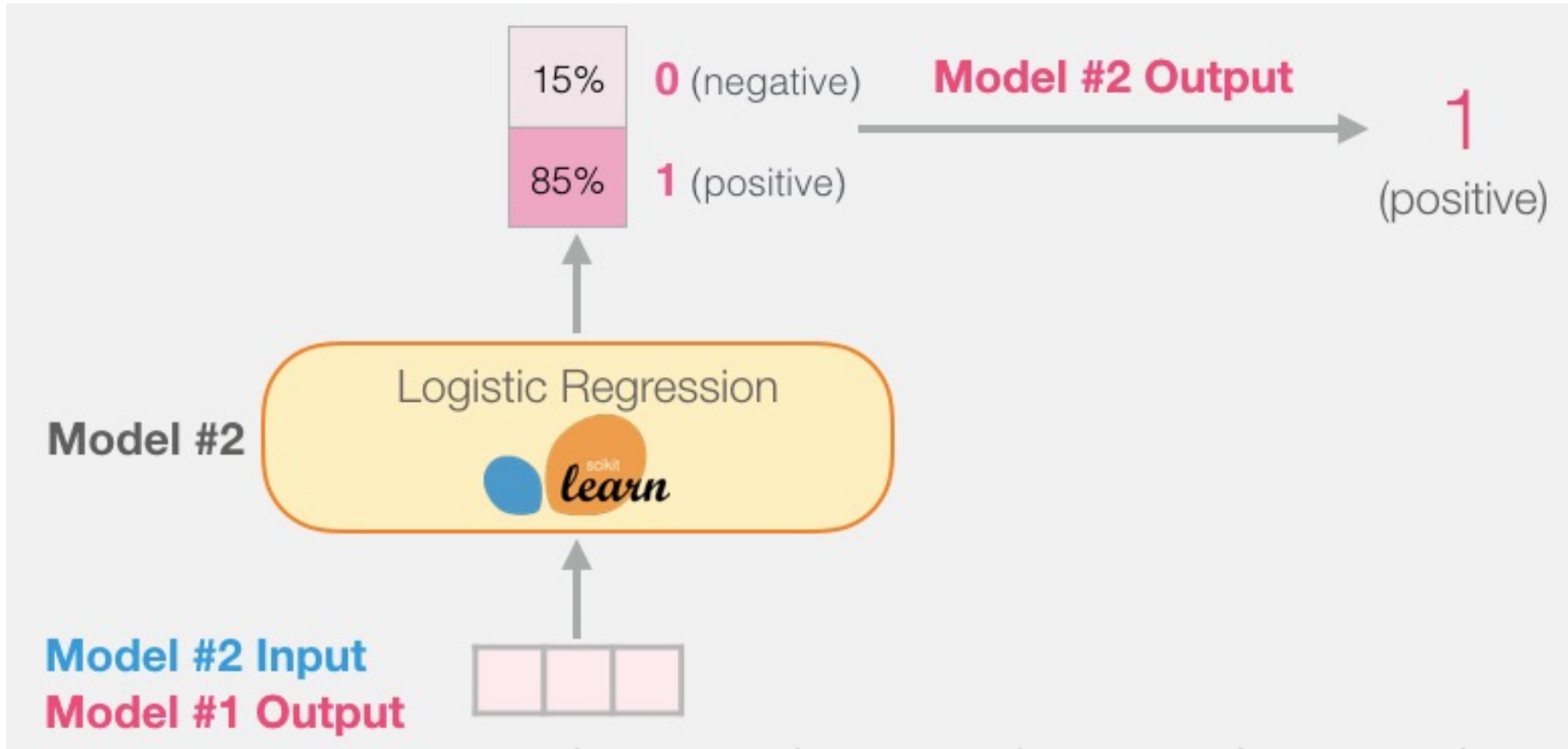# Fine-tuning BERT on Single Sentence Classification Tasks

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
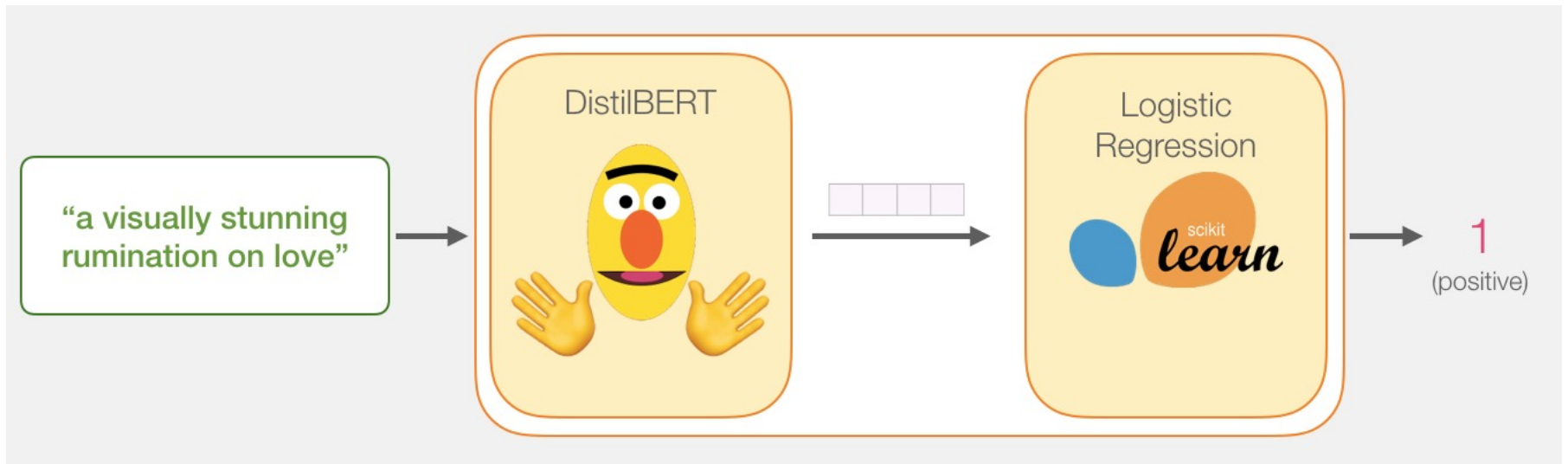"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Model #1 Output Class vector as Model #2 Input

93

# Logistic Regression Model to classify **Class** vector

```python
df = pd.read_csv('https://github.com/clairett/pytorch-
sentiment-classification/raw/master/data/SST2/train.tsv',
delimiter='\t', header=None)

df.head()
```

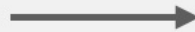|   | 0 | 1 |
|---|---|---|
| 0 | a stirring , funny and finally transporting re... | 1 |
| 1 | apparently reassembled from the cutting room f... | 0 |
| 2 | they presume their audience wo n't sit still f... | 0 |
| 3 | this is a visually stunning rumination on love... | 1 |
| 4 | jonathan parker 's bartleby should have been t... | 1 |

# Tokenization

```
tokenized = df[0].apply((lambda x: tokenizer.encode(x,
add_special_tokens=True)))
```



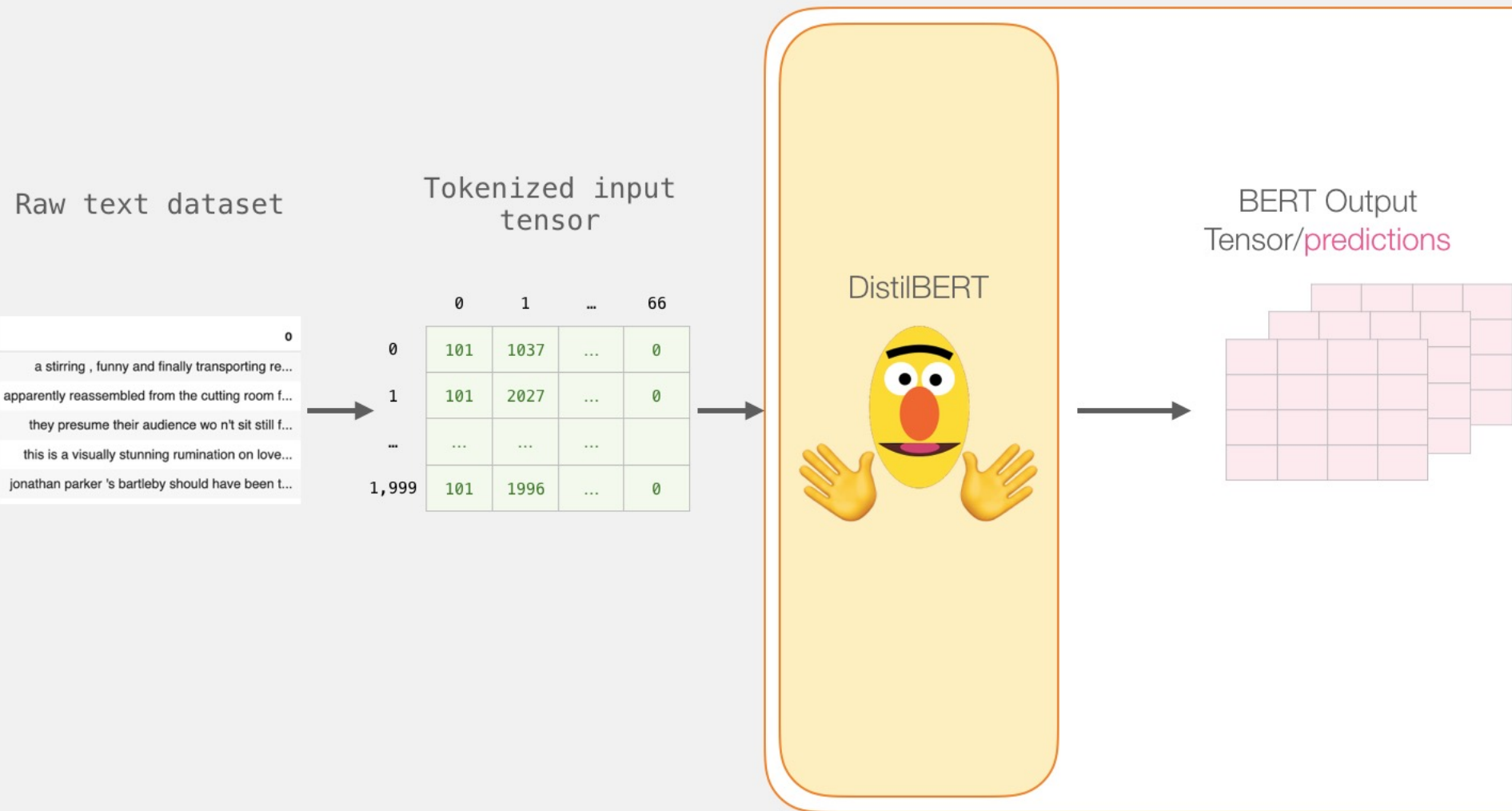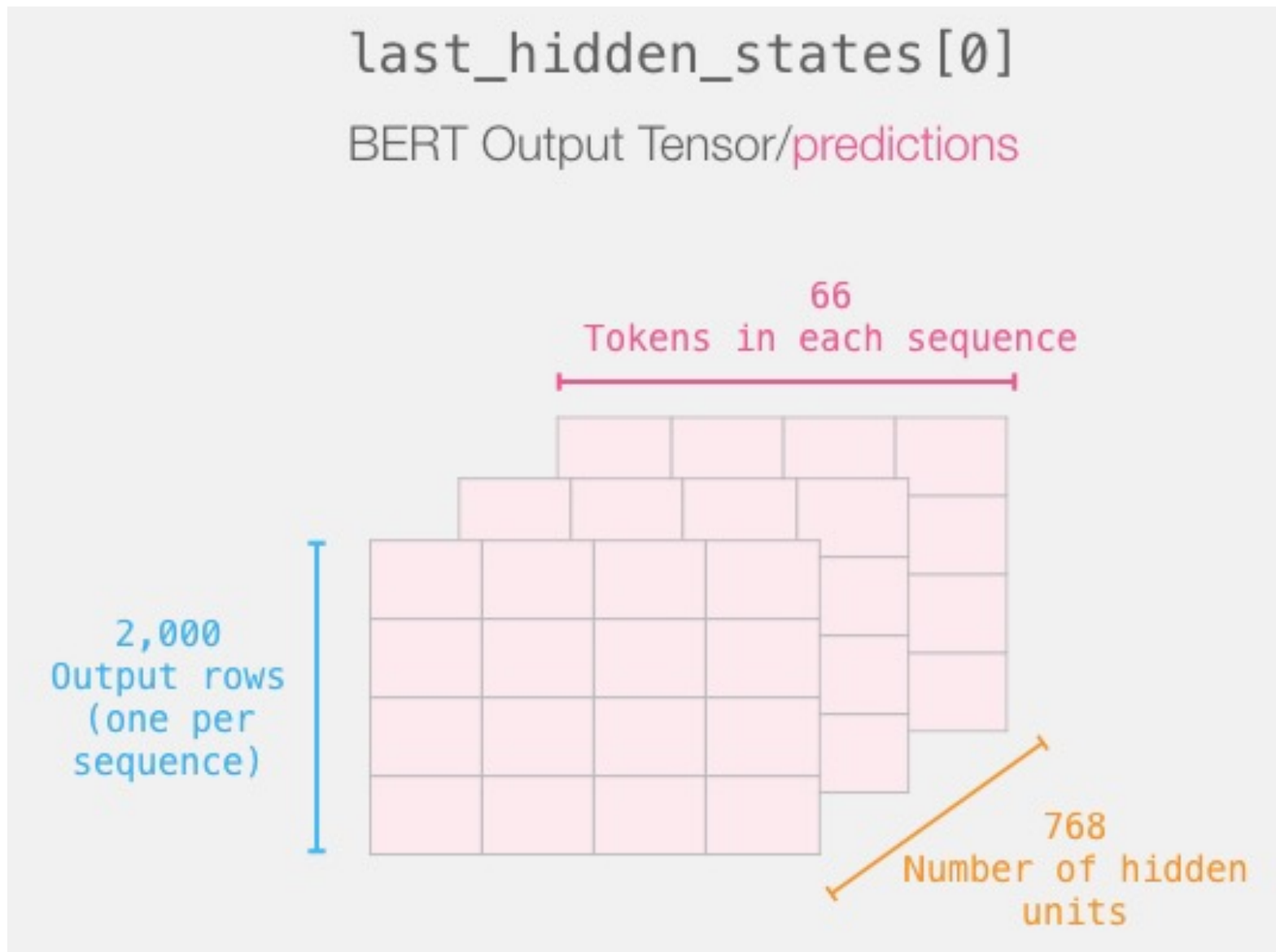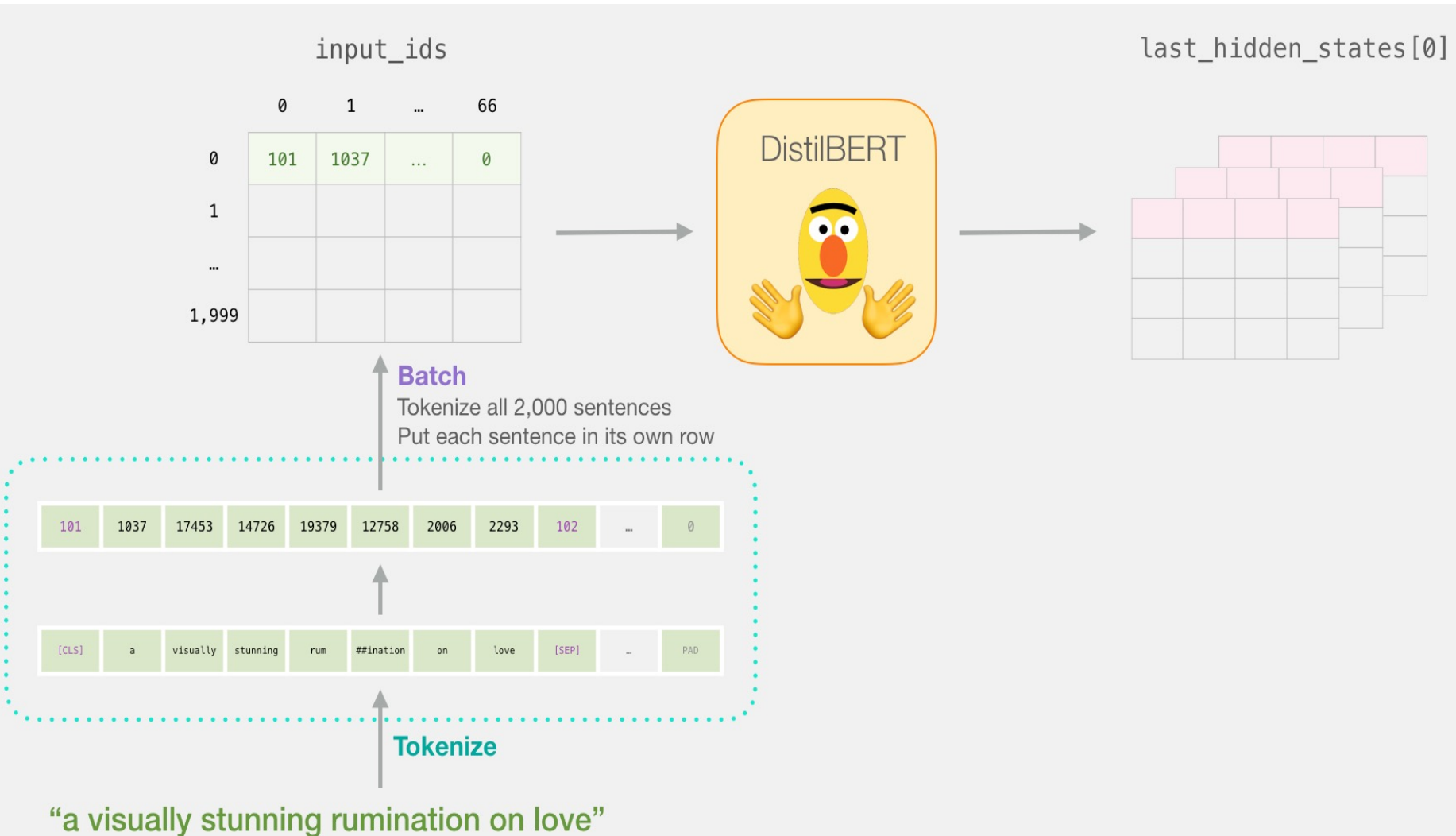| Raw Dataset | | Sequences of Token IDs |
|---|---|---|
| | **0** | |
| a stirring , funny and finally transporting re... | | [101, 1037, 18385, 1010, 6057, 1998, 2633, 182... |
| apparently reassembled from the cutting room f... | Tokenize → | [101, 4593, 2128, 27241, 23931, 2013, 1996, 62... |
| they presume their audience wo n't sit still f... | | [101, 2027, 3653, 23545, 2037, 4378, 24185, 10... |
| this is a visually stunning rumination on love... | | [101, 2023, 2003, 1037, 17453, 14726, 19379, 1... |
| jonathan parker 's bartleby should have been t... | | [101, 5655, 6262, 1005, 1055, 12075, 2571, 376... |

# BERT Input Tensor

# Processing with DistilBERT

```
input_ids = torch.tensor(np.array(padded))
last_hidden_states = model(input_ids)
```

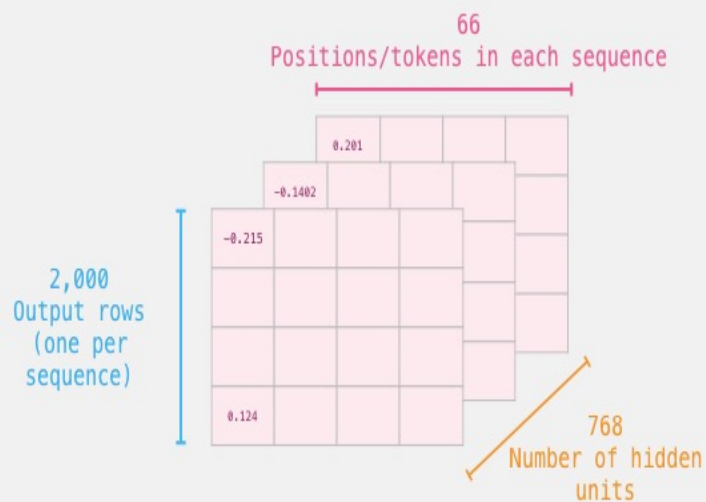# Unpacking the BERT output tensor

# Sentence to last_hidden_state[0]

# BERT's output for the [CLS] tokens

```
# Slice the output for the first position for all the
sequences, take all hidden unit outputs
features = last_hidden_states[0][:,0,:].numpy()
```
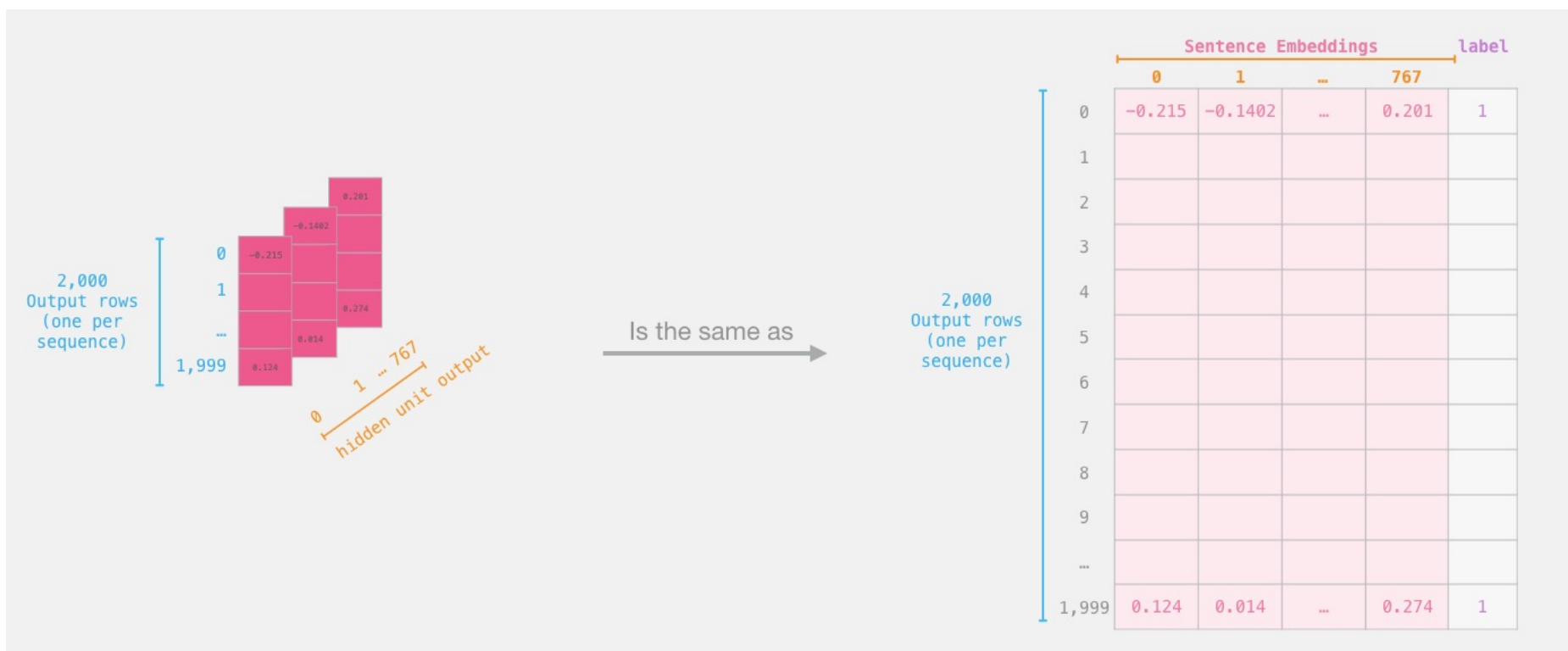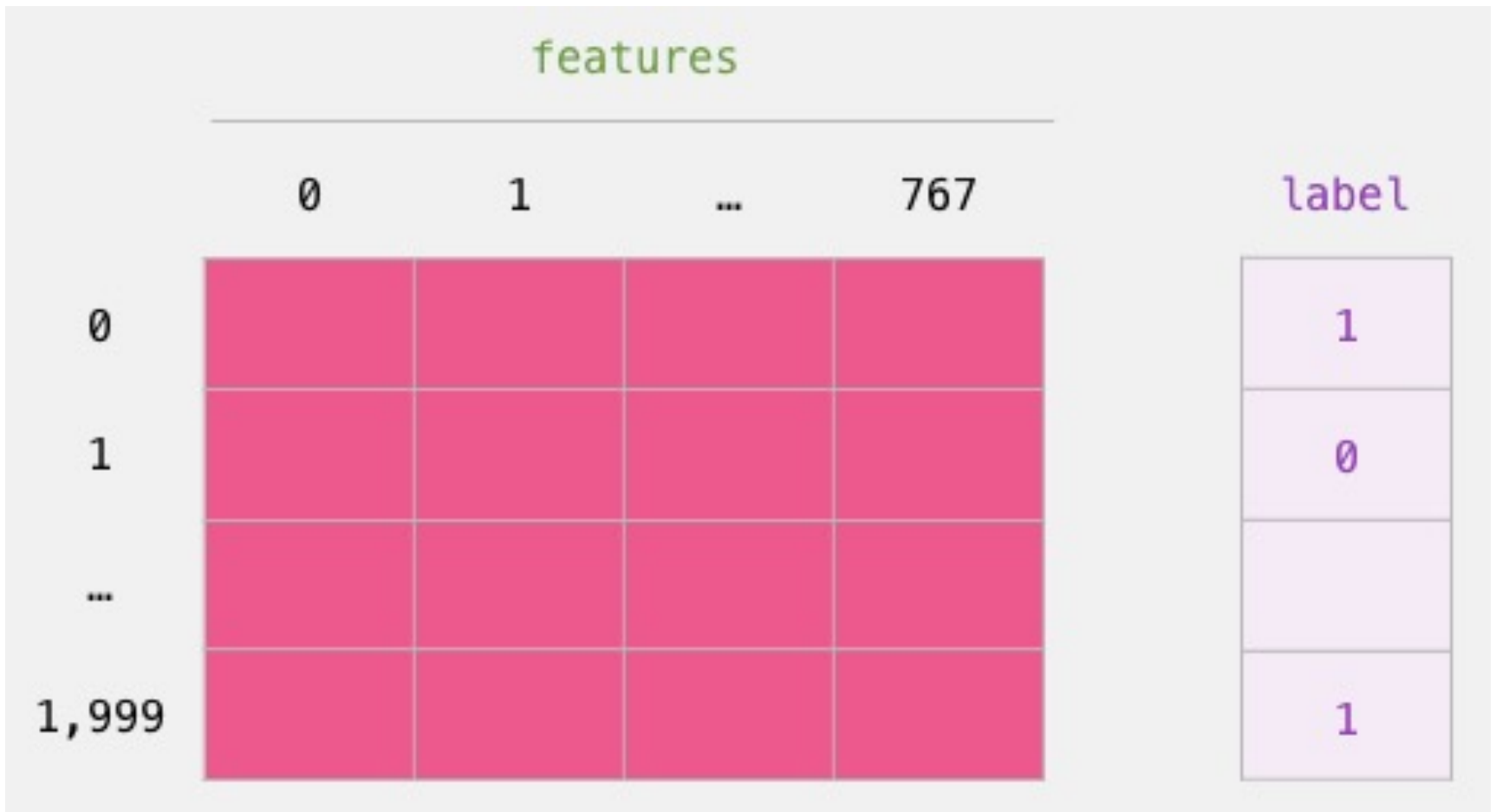
# The tensor sliced from BERT's output
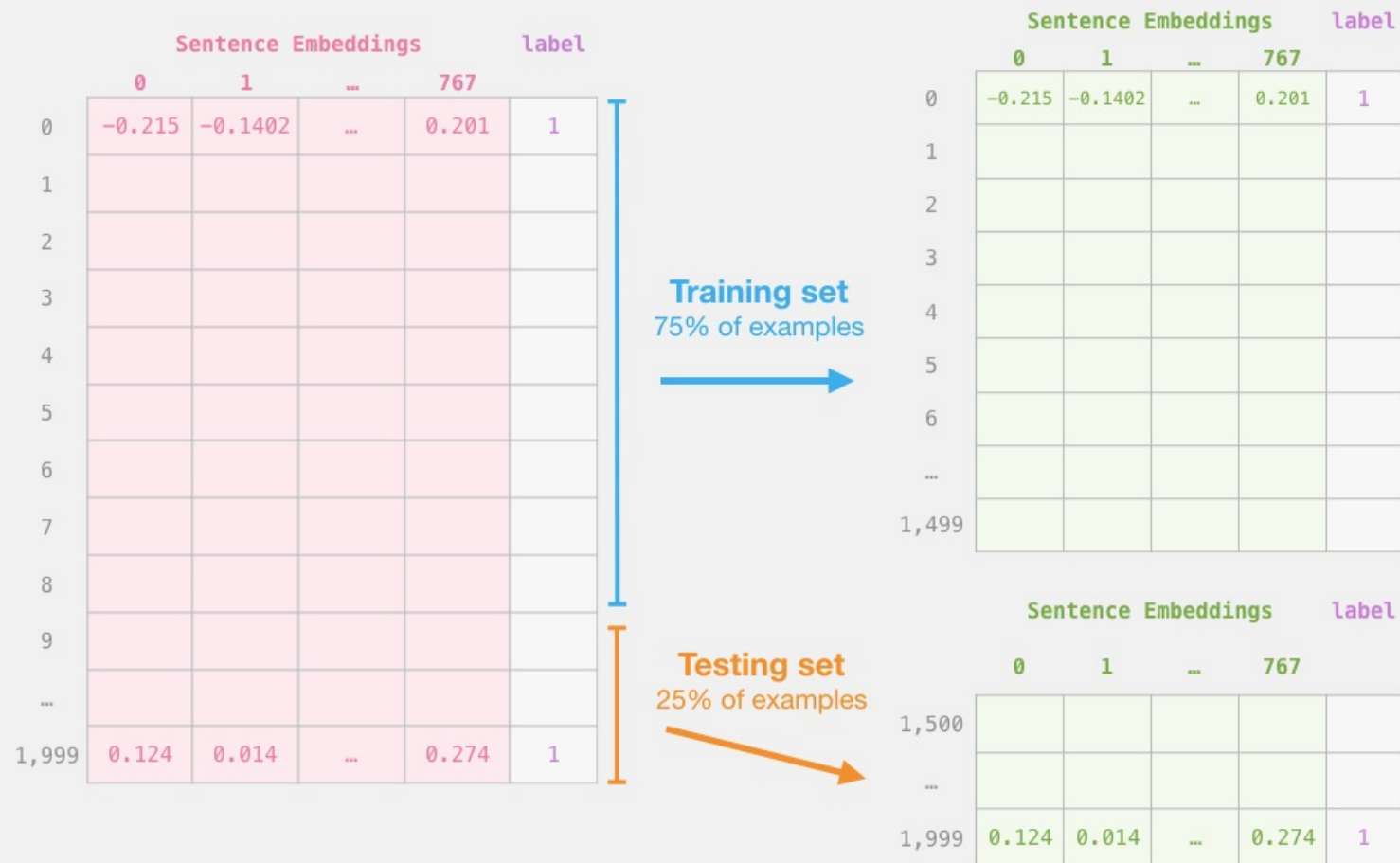## Sentence Embeddings

# Dataset for Logistic Regression (768 Features)

**The features are the output vectors of BERT for the [CLS] token (position #0)**

```
labels = df[1]
train_features, test_features, train_labels, test_labels =
train_test_split(features, labels)
```

Step #2: Test/Train Split for model #2, logistic regression

# Score Benchmarks
# Logistic Regression Model
# on SST-2 Dataset

```
# Training
lr_clf = LogisticRegression()
lr_clf.fit(train_features, train_labels)

#Testing
lr_clf.score(test_features, test_labels)

# Accuracy: 81%
# Highest accuracy: 96.8%
# Fine-tuned DistilBERT: 90.7%
# Full size BERT model: 94.9%
```

# Sentiment Classification: SST2 Sentences from movie reviews

| sentence | label |
|---|---|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

# A Visual Notebook to Using BERT for the First Time

# Pre-trained Language Model (PLM)



Semi-supervised Sequence Learning
context2Vec
Pre-trained seq2seq

ULMFiT — ELMo

GPT

Transformer

Bidirectional LM

Multi-lingual

Larger model
More data

MultiFiT

BERT

GPT-2 — Defense → Grover

Cross-lingual

XLM
UDify

Multi-task

MT-DNN

+ Generation

MASS
UniLM

Knowledge distillation

MT-DNN_KD

Span prediction
Remove NSP

Longer time
Remove NSP
More data

SpanBERT

RoBERTa

Permutation LM
Transformer-XL
More data

XLNet

+Knowledge Graph

ERNIE
(Tsinghua)

Neural entity linker

KnowBert

Cross-modal

VideoBERT
CBT
ViLBERT
VisualBERT
B2T2
Unicoder-VL
LXMERT
VL-BERT
UNITER

Whole Word Masking

ERNIE (Baidu)
BERT-wwm

By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# Turing Natural Language Generation (T-NLG)



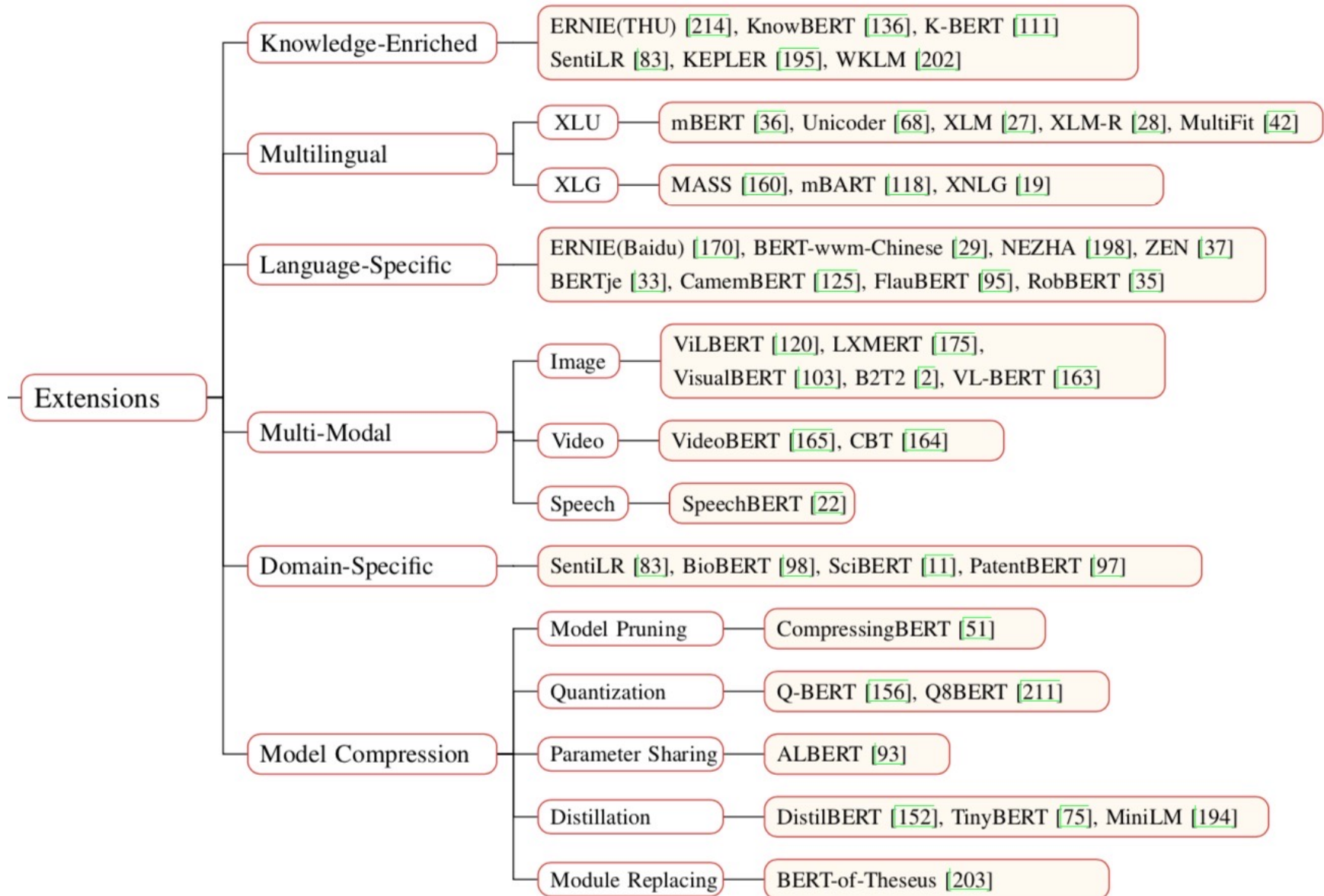**2018**　　　　**2019**　　　　**2020**

# Outline

- Word Embeddings

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- Pretraining and Transfer Learning

- State of the art (SOTA)

# Pre-trained Models (PTM)

# Pre-trained Models (PTM)



```
                    ┌─ Knowledge-Enriched ── ERNIE(THU) [214], KnowBERT [136], K-BERT [111]
                    │                         SentiLR [83], KEPLER [195], WKLM [202]
                    │
                    ├─ Multilingual ─┬─ XLU ── mBERT [36], Unicoder [68], XLM [27], XLM-R [28], MultiFit [42]
                    │                └─ XLG ── MASS [160], mBART [118], XNLG [19]
                    │
                    ├─ Language-Specific ── ERNIE(Baidu) [170], BERT-wwm-Chinese [29], NEZHA [198], ZEN [37]
                    │                        BERTje [33], CamemBERT [125], FlauBERT [95], RobBERT [35]
                    │
  Extensions ──────┤                 ┌─ Image ── ViLBERT [120], LXMERT [175],
                    ├─ Multi-Modal ──┤             VisualBERT [103], B2T2 [2], VL-BERT [163]
                    │                ├─ Video ── VideoBERT [165], CBT [164]
                    │                └─ Speech ── SpeechBERT [22]
                    │
                    ├─ Domain-Specific ── SentiLR [83], BioBERT [98], SciBERT [11], PatentBERT [97]
                    │
                    │                       ┌─ Model Pruning ── CompressingBERT [51]
                    │                       ├─ Quantization ── Q-BERT [156], Q8BERT [211]
                    └─ Model Compression ──┼─ Parameter Sharing ── ALBERT [93]
                                            ├─ Distillation ── DistilBERT [152], TinyBERT [75], MiniLM [194]
                                            └─ Module Replacing ── BERT-of-Theseus [203]
```

Source: Qiu, Xipeng, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. "Pre-trained Models for Natural Language Processing: A Survey." arXiv preprint arXiv:2003.08271 (2020).

# Transformers

## State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch

- Transformers
  - pytorch-transformers
  - pytorch-pretrained-bert
- provides state-of-the-art general-purpose architectures
  - (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, CTRL...)
  - for Natural Language Understanding (NLU) and
    Natural Language Generation (NLG)
    with over 32+ pretrained models
    in 100+ languages
    and deep interoperability between
    TensorFlow 2.0 and
    PyTorch.

Source: https://github.com/huggingface/transformers

# NLP Benchmark Datasets

| Task | Dataset | Link |
|------|---------|------|
| Machine Translation | WMT 2014 EN-DE<br>WMT 2014 EN-FR | http://www-lium.univ-lemans.fr/~schwenk/cslm_joint_paper/ |
| Text Summarization | CNN/DM<br>Newsroom<br>DUC<br>Gigaword | https://cs.nyu.edu/~kcho/DMQA/<br>https://summari.es/<br>https://www-nlpir.nist.gov/projects/duc/data.html<br>https://catalog.ldc.upenn.edu/LDC2012T21 |
| Reading Comprehension<br>Question Answering<br>Question Generation | ARC<br>CliCR<br>CNN/DM<br>NewsQA<br>RACE<br>SQuAD<br>Story Cloze Test<br>NarativeQA<br>Quasar<br>SearchQA | http://data.allenai.org/arc/<br>http://aclweb.org/anthology/N18-1140<br>https://cs.nyu.edu/~kcho/DMQA/<br>https://datasets.maluuba.com/NewsQA<br>http://www.qizhexie.com/data/RACE_leaderboard<br>https://rajpurkar.github.io/SQuAD-explorer/<br>http://aclweb.org/anthology/W17-0906.pdf<br>https://github.com/deepmind/narrativeqa<br>https://github.com/bdhingra/quasar<br>https://github.com/nyu-dl/SearchQA |
| Semantic Parsing | AMR parsing<br>ATIS (SQL Parsing)<br>WikiSQL (SQL Parsing) | https://amr.isi.edu/index.html<br>https://github.com/jkkummerfeld/text2sql-data/tree/master/data<br>https://github.com/salesforce/WikiSQL |
| Sentiment Analysis | IMDB Reviews<br>SST<br>Yelp Reviews<br>Subjectivity Dataset | http://ai.stanford.edu/~amaas/data/sentiment/<br>https://nlp.stanford.edu/sentiment/index.html<br>https://www.yelp.com/dataset/challenge<br>http://www.cs.cornell.edu/people/pabo/movie-review-data/ |
| Text Classification | AG News<br>DBpedia<br>TREC<br>20 NewsGroup | http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html<br>https://wiki.dbpedia.org/Datasets<br>https://trec.nist.gov/data.html<br>http://qwone.com/~jason/20Newsgroups/ |
| Natural Language Inference | SNLI Corpus<br>MultiNLI<br>SciTail | https://nlp.stanford.edu/projects/snli/<br>https://www.nyu.edu/projects/bowman/multinli/<br>http://data.allenai.org/scitail/ |
| Semantic Role Labeling | Proposition Bank<br>OneNotes | http://propbank.github.io/<br>https://catalog.ldc.upenn.edu/LDC2013T19 |

114

# Question Answering

# (QA)

# SQuAD

**Stanford Question Answering Dataset**

# SQuAD

# SQuAD2.0
## The Stanford Question Answering Dataset

## What is SQuAD?

**S**tanford **Qu**estion **A**nswering **D**ataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or *span*, from the corresponding reading passage, or the question might be unanswerable.

SQuAD2.0 combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

## Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph.

| Rank | Model | EM | F1 |
|---|---|---|---|
|  | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Apr 06, 2020 | SA-Net on Albert (ensemble)<br>*QIANXIN* | **90.724** | **93.011** |
| 2<br>May 05, 2020 | SA-Net-V2 (ensemble)<br>*QIANXIN* | 90.679 | 92.948 |
| 2 | Retro-Reader (ensemble) | 90.578 | 92.978 |

# SQuAD

## SQuAD: 100,000+ Questions for Machine Comprehension of Text

**Pranav Rajpurkar** and **Jian Zhang** and **Konstantin Lopyrev** and **Percy Liang**
{pranavsr,zjian,klopyrev,pliang}@cs.stanford.edu
Computer Science Department
Stanford University

### Abstract

We present the Stanford Question Answering Dataset (SQuAD), a new reading comprehension dataset consisting of 100,000+ questions posed by crowdworkers on a set of Wikipedia articles, where the answer to each question is a segment of text from the corresponding reading passage. We analyze the dataset to understand the types of reasoning required to answer the questions, leaning heavily on dependency and constituency trees. We build a strong logistic regression model, which achieves an F1 score of 51.0%, a significant improvement over a simple baseline (20%). However, human performance (86.8%) is much higher, indicating that the dataset presents a good challenge problem for future research. The dataset is freely available at https://stanford-qa.com.

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
**gravity**

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
**graupel**

Where do water droplets collide with ice crystals to form precipitation?
**within a cloud**

**Figure 1:** Question-answer pairs for a sample passage in the

117

# SQuAD (Question Answering)

**Q:** What causes precipitation to fall?

## Precipitation

From Wikipedia, the free encyclopedia

*For other uses, see Precipitation (disambiguation).*

In meteorology, **precipitation** is any product of the condensation of atmospheric water vapor that falls under gravity from clouds.[2] The main forms of precipitation include drizzle, rain, sleet, snow, ice pellets, graupel and hail. Precipitation occurs when a portion of the atmosphere becomes saturated with water vapor (reaching 100% relative humidity), so that the water condenses and "precipitates". Thus, fog and mist are not precipitation but suspensions, because the water vapor does not condense sufficiently to precipitate. Two processes, possibly acting together, can lead to air becoming saturated: cooling the air or adding water vapor to the air. Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers."[3]

https://en.wikipedia.org/wiki/Precipitation

# SQuAD (Question Answering)

Paragraph

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail… Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

**A:** gravity

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

**A:** graupel

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** Where do water droplets collide with ice crystals to form precipitation?

**A:** within a cloud

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

**A:** gravity

**Q:** What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

**A:** graupel

**Q:** Where do water droplets collide with ice crystals to form precipitation?

**A:** within a cloud

# Natural Language Processing with Python
## – Analyzing Text with the Natural Language Toolkit

← → C  ⓘ www.nltk.org/book/

# Natural Language Processing with Python

## – Analyzing Text with the Natural Language Toolkit

**Steven Bird, Ewan Klein, and Edward Loper**

*This version of the NLTK book is updated for Python 3 and NLTK 3. The first edition of the book, published by O'Reilly, is available at http://nltk.org/book_1ed/. (There are currently no plans for a second edition of the book.)*

http://www.nltk.org/book/

# spaCy



## Industrial-Strength Natural Language Processing
### in Python

**Fastest in the world**

spaCy excels at large-scale information extraction tasks. It's written from the ground up in carefully memory-managed Cython. Independent research has confirmed that spaCy is the fastest in the world. If your application needs to process entire web dumps, spaCy is the library you want to be using.

**Get things done**

spaCy is designed to help you do real work — to build real products, or gather real insights. The library respects your time, and tries to avoid wasting it. It's easy to install, and its API is simple and productive. I like to think of spaCy as the Ruby on Rails of Natural Language Processing.

**Deep learning**

spaCy is the best way to prepare text for deep learning. It interoperates seamlessly with TensorFlow, Keras, Scikit-Learn, Gensim and the rest of Python's awesome AI ecosystem. spaCy helps you connect the statistical models trained by these libraries to the rest of your application.

https://spacy.io/

# gensim

# TextBlob

## TextBlob: Simplified Text Processing

Release v0.12.0. ([Changelog](#))

*TextBlob* is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

```python
from textblob import TextBlob

text = '''
The titular threat of The Blob has always struck me as the ultimate movie
monster: an insatiably hungry, amoeba-like mass able to penetrate
virtually any safeguard, capable of--as a doomed doctor chillingly
describes it--"assimilating flesh on contact.
Snide comparisons to gelatin be damned, it's a concept with the most
devastating of potential consequences, not unlike the grey goo scenario
proposed by technological theorists fearful of
artificial intelligence run rampant.
'''

blob = TextBlob(text)
blob.tags            # [('The', 'DT'), ('titular', 'JJ'),
                     #  ('threat', 'NN'), ('of', 'IN'), ...]

blob.noun_phrases    # WordList(['titular threat', 'blob',
                     #           'ultimate movie monster',
                     #           'amoeba-like mass', ...])

for sentence in blob.sentences:
    print(sentence.sentiment.polarity)
# 0.060
```

https://textblob.readthedocs.io

# Polyglot

Search docs

Installation

Language Detection

Tokenization

Command Line Interface

Downloading Models

Word Embeddings

Part of Speech Tagging

Named Entity Extraction

Morphological Analysis

Transliteration

Sentiment

polyglot

Docs » Welcome to polyglot's documentation!                    ○ Edit on GitHub

## Welcome to polyglot's documentation!

## polyglot

`downloads` `17k/month`  `pypi package` `16.7.4`  `build` `passing`  `docs` `passing`

Polyglot is a natural language pipeline that supports massive multilingual applications.

- Free software: GPLv3 license
- Documentation: http://polyglot.readthedocs.org.

## Features

- Tokenization (165 Languages)
- Language detection (196 Languages)
- Named Entity Recognition (40 Languages)
- Part of Speech Tagging (16 Languages)
- Sentiment Analysis (136 Languages)
- Word Embeddings (137 Languages)
- Morphological analysis (135 Languages)
- Transliteration (69 Languages)

https://polyglot.readthedocs.io/

128

# Text Classification with TF Hub

# Text Classification with Pre Text

# Text Classification
# IMDB Movie Reviews

https://colab.research.google.com/drive/1x16h1GhHsLIrLYtPCvCHaoO1W-i_gror

# Papers with Code
# State-of-the-Art (SOTA)

## Browse State-of-the-Art

📈 1509 leaderboards • 1327 tasks • 1347 datasets • 17810 papers with code

Follow on 🐦 Twitter for updates

## Computer Vision

| | | | | |
|---|---|---|---|---|
| Semantic Segmentation | Image Classification | Object Detection | Image Generation | Pose Estimation |
| 📈 33 leaderboards | 📈 52 leaderboards | 📈 54 leaderboards | 📈 51 leaderboards | 📈 40 leaderboards |
| 667 papers with code | 564 papers with code | 467 papers with code | 231 papers with code | 231 papers with code |

▸ See all 707 tasks

## Natural Language Processing

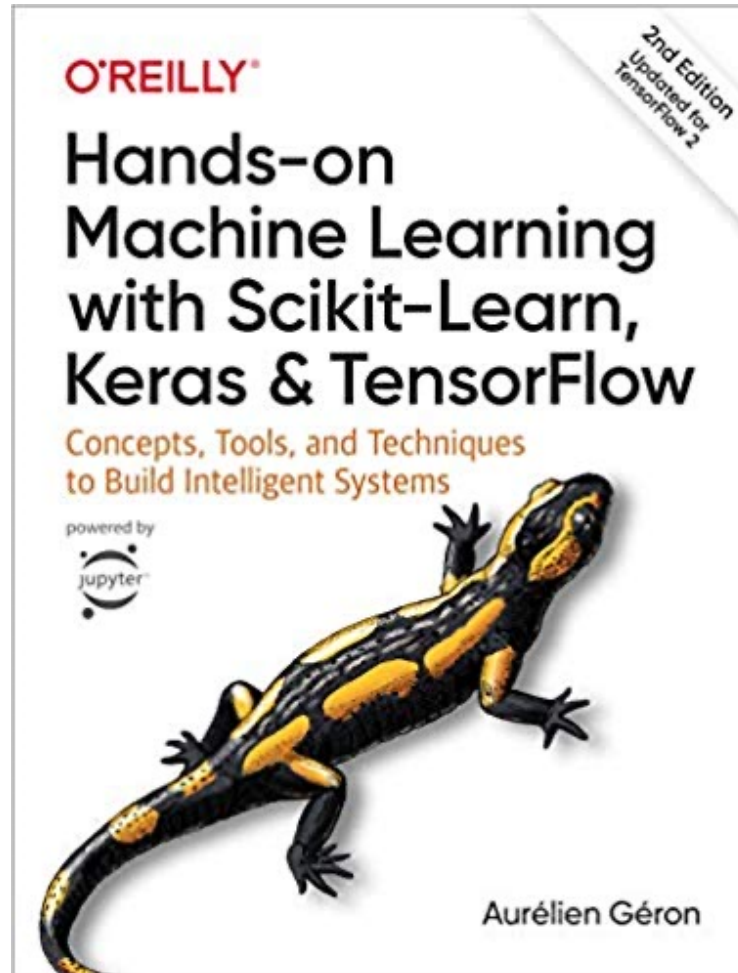| | | | | |
|---|---|---|---|---|
| Machine Translation | Language Modelling | Question Answering | Sentiment Analysis | Text Generation |

https://paperswithcode.com/sota

# Aurélien Géron (2019),
# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition
# O'Reilly Media, 2019



https://github.com/ageron/handson-ml2

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

**Notebooks**
1. The Machine Learning landscape
2. End-to-end Machine Learning project
3. Classification
4. Training Models
5. Support Vector Machines
6. Decision Trees
7. Ensemble Learning and Random Forests
8. Dimensionality Reduction
9. Unsupervised Learning Techniques
10. Artificial Neural Nets with Keras
11. Training Deep Neural Networks
12. Custom Models and Training with TensorFlow
13. Loading and Preprocessing Data
14. Deep Computer Vision Using Convolutional Neural Networks
15. Processing Sequences Using RNNs and CNNs
16. Natural Language Processing with RNNs and Attention
17. Representation Learning Using Autoencoders
18. Reinforcement Learning
19. Training and Deploying TensorFlow Models at Scale

https://github.com/ageron/handson-ml2

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

CO   🔺 python101.ipynb ☆

💬 Comment   👥 Share ⚙ A

+ Code   + Text     ✓ RAM ▭ ⌄   ✏ Editing ⌃

## Text Classification: IMDB Movie Review

- Source: https://www.tensorflow.org/tutorials/keras/text_classification_with_hub

```
[1]  1 !pip install -q tensorflow-hub
     2 !pip install -q tensorflow-datasets
```

```
[2]  1 import os
     2 import numpy as np
     3
     4 import tensorflow as tf
     5 import tensorflow_hub as hub
     6 import tensorflow_datasets as tfds
     7
     8 print("Version: ", tf.__version__)
     9 print("Eager mode: ", tf.executing_eagerly())
    10 print("Hub version: ", hub.__version__)
    11 print("GPU is", "available" if tf.config.list_physical_devices("GPU") else "NOT AVAILABLE")
```

```
Version:  2.4.1
Eager mode:  True
Hub version:  0.12.0
GPU is available
```

```
[3]  1 # Split the training set into 60% and 40% to end up with 15,000 examples
     2 # for training, 10,000 examples for validation and 25,000 examples for testing.
     3 train_data, validation_data, test_data = tfds.load(
     4     name="imdb_reviews",
```

https://tinyurl.com/aintpupython101

135

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

python101.ipynb

File  Edit  View  Insert  Runtime  Tools  Help    All changes saved

+ Code    + Text

- Huggingface Transformers: https://github.com/huggingface/transformers

```
[18]  1 !pip install transformers
```

```
1 from transformers import pipeline
2 classifier = pipeline('sentiment-analysis')
3 classifier('We are very happy to introduce pipeline to the transformers repository.')
```

Downloading: 100%    629/629 [00:00<00:00, 1.31kB/s]

Downloading: 100%    268M/268M [00:05<00:00, 46.9MB/s]

Downloading: 100%    232k/232k [00:01<00:00, 159kB/s]

Downloading: 100%    48.0/48.0 [00:00<00:00, 522B/s]

```
[{'label': 'POSITIVE', 'score': 0.9996980428695679}]
```

```
[11]  1 classifier('This movie is very good.')
```

```
[{'label': 'POSITIVE', 'score': 0.9998621940612793}]
```

```
[12]  1 classifier('This movie is very boring.')
```

```
[{'label': 'NEGATIVE', 'score': 0.999795138835907}]
```

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

```python
# from transformers import pipeline
question_answerer = pipeline('question-answering')
question_answerer({'question': 'What is the name of the repository ?',
                   'context': 'Pipeline has been included in the huggingface/transformers repository'})
```

```
{'answer': 'huggingface/transformers',
 'end': 58,
 'score': 0.309702068567276,
 'start': 34}
```

```python
context = '''In meteorology, precipitation is any product of the condensation of atmospheric water vapor
context
```

'In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".'

```python
question_answerer({'question': 'Where do water droplets collide with ice crystals to form precipitation?',
                   'context': context})
```

```
{'answer': 'within a cloud',
 'end': 321,
 'score': 0.5175967812538147,
 'start': 307}
```

```python
question_answerer({'question': 'What causes precipitation to fall?',
                   'context': context})
```

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

python101.ipynb

File   Edit   View   Insert   Runtime   Tools   Help   All changes saved

💬 Comment   👥 Share   ⚙   A

+ Code   + Text

RAM ▭
Disk ▭
✏ Editing   ⌃

**Table of contents**   ✕

## Deep Learning and Universal Sentence-Embedding Models

↑  ↓  🔗  💬  ✏  🗑  ⋮

### Universal Sentence Encoder (USE)

- Source: Universal Sentence Encoder: https://tfhub.dev/google/universal-sentence-encoder/4

```python
import tensorflow as tf
import tensorflow_hub as hub
import numpy as np
import pandas as pd
import os
import re
import matplotlib.pyplot as plt
import seaborn as sns

module_url = "https://tfhub.dev/google/universal-sentence-encoder/4"
#"https://tfhub.dev/google/universal-sentence-encoder-large/5"
model = hub.load(module_url)
print ("module %s loaded" % module_url)
def embed(input):
  return model(input)
```

module https://tfhub.dev/google/universal-sentence-encoder/4 loaded

```python
word = "Elephant"
sentence = "I am a sentence for which I would like to get its embedding."
```

138

# Python in Google Colab (Python101)

Semantic Textual Similarity

# Summary

- Word Embeddings

- Recurrent Neural Networks for NLP

- Sequence-to-Sequence Models

- The Transformer Architecture

- Pretraining and Transfer Learning

- State of the art (SOTA)

# References

- Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson.
- Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media.
- Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress. https://github.com/Apress/text-analytics-w-python-2e
- Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python, O'Reilly Media. https://www.oreilly.com/library/view/applied-text-analysis/9781491963036/
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Céspedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, Ray Kurzweil (2018). Universal Sentence Encoder. arXiv:1803.11175.
- Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez Abrego , Steve Yuan, Chris Tar, Yun-hsuan Sung, Ray Kurzweil (2019). Multilingual Universal Sentence Encoder for Semantic Retrieval.
- Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang (2020). "Pre-trained Models for Natural Language Processing: A Survey." arXiv preprint arXiv:2003.08271.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.
- Jay Alammar (2019), The Illustrated Transformer, http://jalammar.github.io/illustrated-transformer/
- Jay Alammar (2019), A Visual Guide to Using BERT for the First Time, http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/
- Christopher Olah, (2015) Understanding LSTM Networks, http://colah.github.io/posts/2015-08-Understanding-LSTMs/
- HuggingFace (2020), Transformers Notebook, https://huggingface.co/transformers/notebooks.html
- The Super Duper NLP Repo, https://notebooks.quantumstat.com/
- Min-Yuh Day (2021), Python 101, https://tinyurl.com/aintpupython101