# Lexicon, Syntax, Semantics IIb: Modeling Meaning

## Contextualized word embeddings

Eva Maria Vecchi

Center for Information and Language Processing
LMU Munich

July 16, 2020

# bank

The man was accused of robbing a bank.

The man went fishing by the bank of the river.

# *bank*

*The man was accused of robbing a bank.*

*The man went fishing by the bank of the river.*

- **word2vec**: same word embedding for the word "bank" in both sentences
  - each word has a fixed representation under Word2Vec regardless of the context within which the word appears
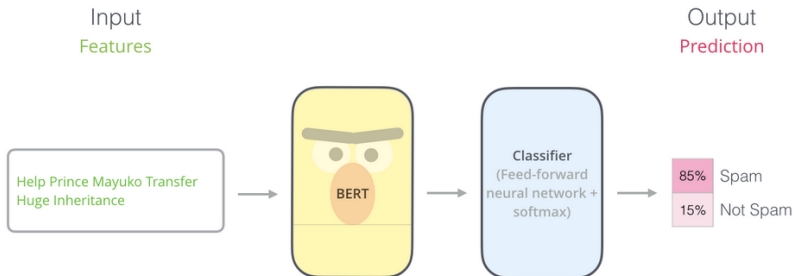
# *bank*

*The man was accused of robbing a bank.*

*The man went fishing by the bank of the river.*

- **word2vec**: same word embedding for the word "bank" in both sentences
  - each word has a fixed representation under Word2Vec regardless of the context within which the word appears

- **contextualized embedding** (e.g. BERT): word embedding for "bank" would be different for each sentence
  - produces word representations that are dynamically informed by the words around them

# Task Example: Sentence Classification



Input
Features

Output
Prediction

Help Prince Mayuko Transfer
Huge Inheritance

BERT

Classifier
(Feed-forward
neural network +
softmax)

85% Spam

15% Not Spam

# Other examples

- **Sentiment Analysis**
    - Input: Movie/Product review. Output: is the review positive or negative?
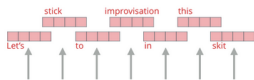    - Example dataset: SST
- **Fact-checking**
    - Input: sentence. Output: "Claim" or "Not Claim"
    - More ambitious/futuristic example:
        - Input: Claim sentence. Output: "True" or "False"
    - Full Fact is an organization building automatic fact-checking tools for the benefit of the public. Part of their pipeline is a classifier that reads news articles and detects claims (classifies text as either "claim" or "not claim") which can later be fact-checked (by humans now, by with ML later, hopefully).

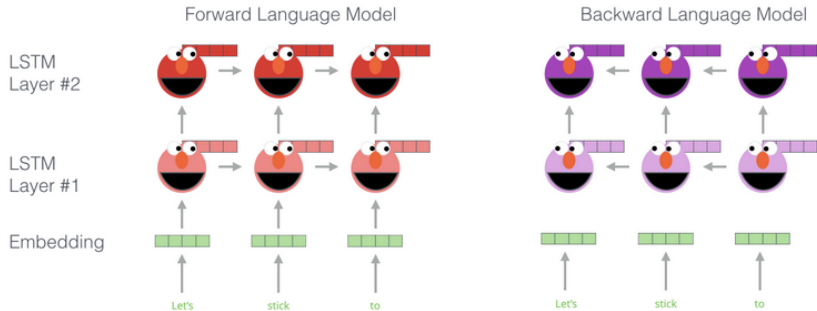# ELMo: Context Matters

# ELMo Embeddings



- Instead of using a fixed embedding for each word, ELMo looks at the entire sentence before assigning each word in it an embedding.
- Uses a bi-directional LSTM (e.g. language model) trained on a specific task to be able to create those embeddings

# ELMo Embeddings



Embedding of "stick" in "Let's stick to" - Step #1

# ELMo Embeddings

Embedding of "stick" in "Let's stick to" - Step #2

1- Concatenate hidden layers
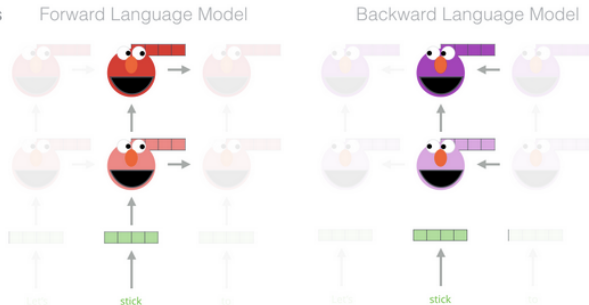
2- Multiply each vector by a weight based on the task

x  s₂
x  s₁
x  s₀

3- Sum the (now weighted) vectors

ELMo embedding of "stick" for this task in this context

Forward Language Model

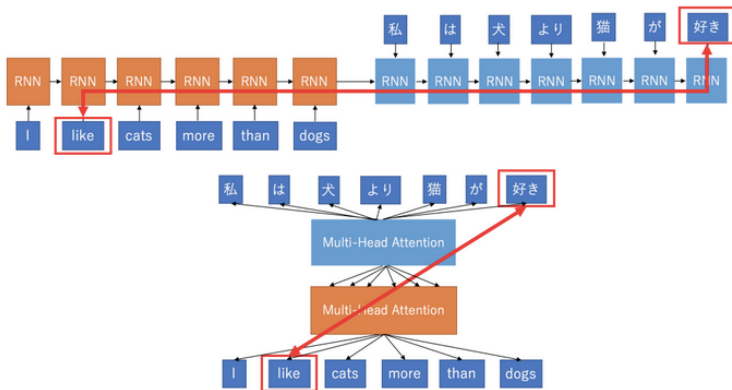Backward Language Model

stick

stick

# The Transformer: Going beyond LSTMs

- Google AI presented Transformer architecture that showed great benefit over conventional sequence models (LSTM, RNN, GRU)

  (Vaswani et al, 2017)

- Advantages:
  - more effective modeling of long term dependencies among tokens in a temporal sequence
  - more efficient training of the model in general by eliminating the sequential dependency on previous tokens
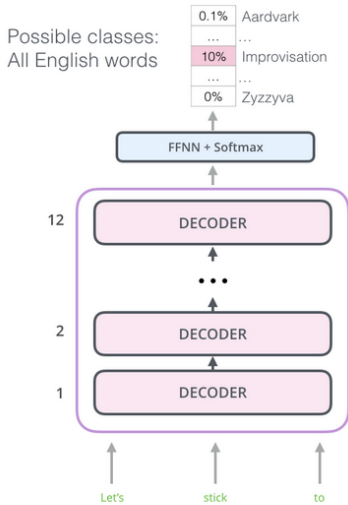
# The Transformer: Going beyond LSTMs

- In a nutshell... a transformer is an **encoder-decoder** architecture model which uses **attention mechanisms** to forward a more complete picture of the whole sequence to the decoder at once rather than sequentially
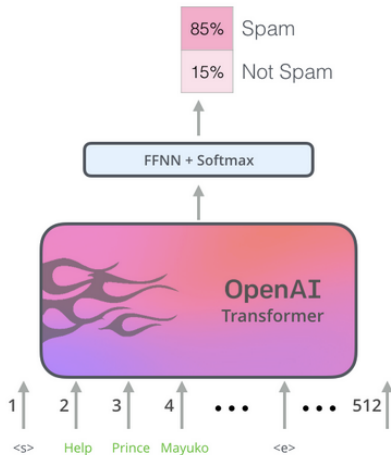
# OpenAI Transformer

- Turns out, you don't need an entire transformer – all you need is the decoder!
- **Decoder**: a natural choice for language modeling (predicting the next word) as it's built to mask future tokens – a valuable feature when it's generating a translation word by word

# OpenAI Transformer



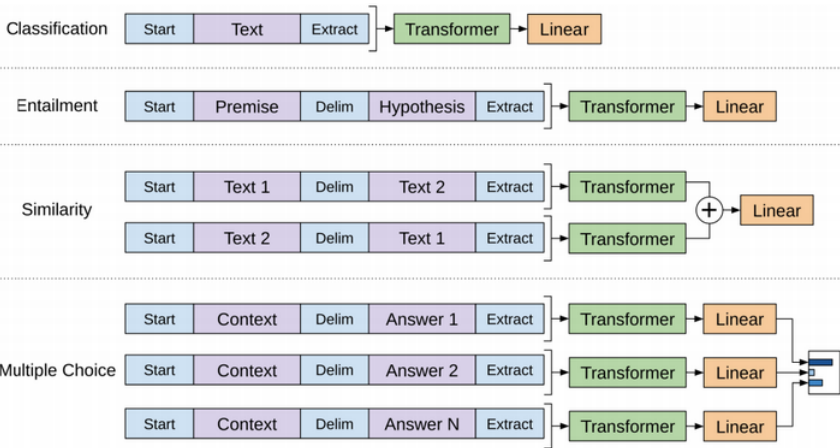The OpenAI Transformer is now ready to be trained to predict the next word on a dataset made up of 7,000 books.

# OpenAI Transformer



How to use a pre-trained OpenAI transformer to do sentence clasification

# OpenAI Transformer

Various structures of the models and input transformations to carry out different downstream tasks

# Something's missing

- Something went missing in transition from LSTM architecture of ELMo...

- Could we build a transformer-based model whose language model looks both forward and backwards (conditioned on both left and right context)?

# Something's missing

- Something went missing in transition from LSTM architecture of ELMo. . .

- Could we build a transformer-based model whose language model looks both forward and backwards (conditioned on both left and right context)?

- "Hold my beer", said BERT, aka Bidirectional Encoder Representations from Transformer
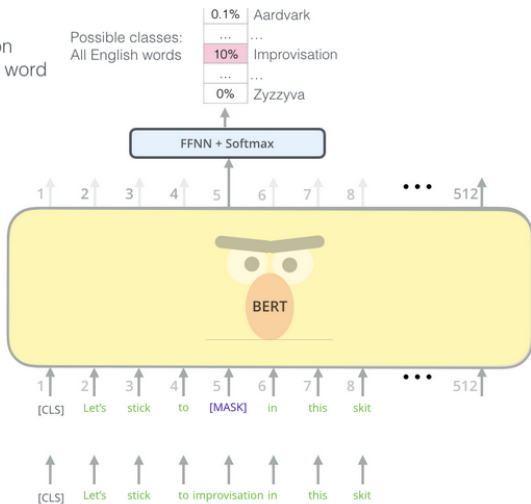
# Masked Language Model

- "We'll use transformer encoders", said BERT.

- "This is madness", replied Ernie, "Everybody knows bidirectional conditioning would allow each word to indirectly see itself in a multi-layered context."

- "We'll use masks", said BERT confidently.

# BERT

Delvin et al, 2018 (arXiv)



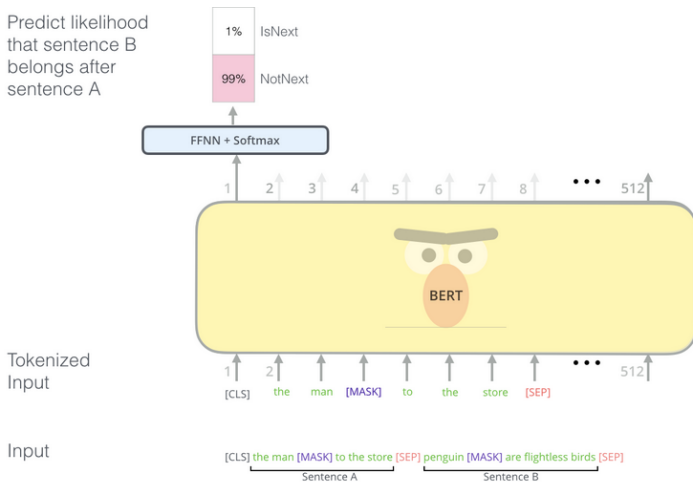Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| | |
|---|---|
| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  ...  512

BERT

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  ...  512

[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

# Two-sentence Tasks with BERT
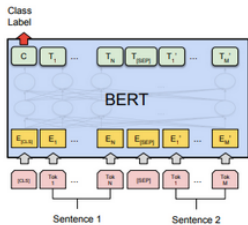


Predict likelihood that sentence B belongs after sentence A

| | |
|---|---|
| 1% | IsNext |
| 99% | NotNext |

FFNN + Softmax

1 2 3 4 5 6 7 8 ••• 512

BERT

Tokenized Input

1 2 3 4 5 6 7 8 ••• 512

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

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]
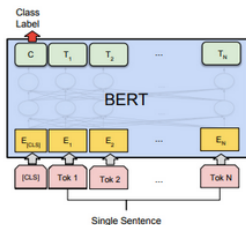
Sentence A        Sentence B

The second task BERT is pre-trained on is a two-sentence classification task. The tokenization is oversimplified in this graphic as BERT actually uses WordPieces as tokens rather than words --- so some words are broken down into smaller chunks.
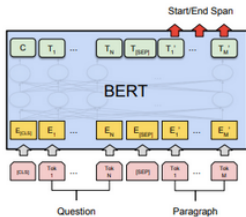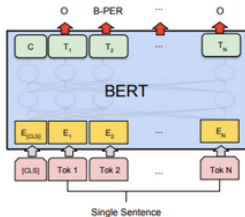
# Task Specific Models



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

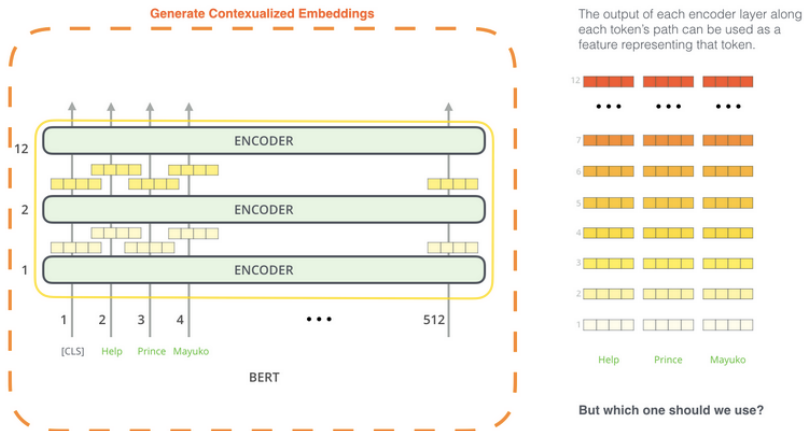(b) Single Sentence Classification Tasks:
SST-2, CoLA

(c) Question Answering Tasks:
SQuAD v1.1

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

# BERT for feature extraction



Use the pre-trained BERT to create contextualized word embeddings. Then you can feed these embeddings to your existing model, e.g. named entity recognition.

# Which layer to choose?



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

| | Dev F1 Score |
|---|---|
| First Layer | 91.0 |
| Last Hidden Layer | 94.9 |
| Sum All 12 Layers | 95.5 |
| Second-to-Last Hidden Layer | 95.6 |
| Sum Last Four Hidden | 95.9 |
| Concat Last Four Hidden | 96.1 |

# BERT Unveiled

## Delvin et al, 2018 (arXiv)



**Thang Luong**
@lmthang

Follow

A new era of NLP has just begun a few days ago: large pretraining models (Transformer 24 layers, 1024 dim, 16 heads) + massive compute is all you need. BERT from @GoogleAI: SOTA results on everything arxiv.org/abs/1810.04805. Results on SQuAD are just mind-blowing. Fun time ahead!

### SQuAD1.1 Leaderboard

Since the release of SQuAD1.0, the community has made rapid progress, with the best models now rivaling human performance on the task. Here are the ExactMatch (EM) and F1 scores evaluated on the test set of v1.1.

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar et al. '16) | 82.304 | 91.221 |
| 1<br>Oct 05, 2018 | BERT (ensemble)<br>*Google A.I.* | **87.433** | **93.160** |
| 2<br>Oct 05, 2018 | BERT (single model)<br>*Google A.I.* | 85.083 | 91.835 |
| 2<br>Sep 09, 2018 | nlnet (ensemble)<br>*Microsoft Research Asia* | 85.356 | 91.202 |