

CS5740: Natural Language Processing

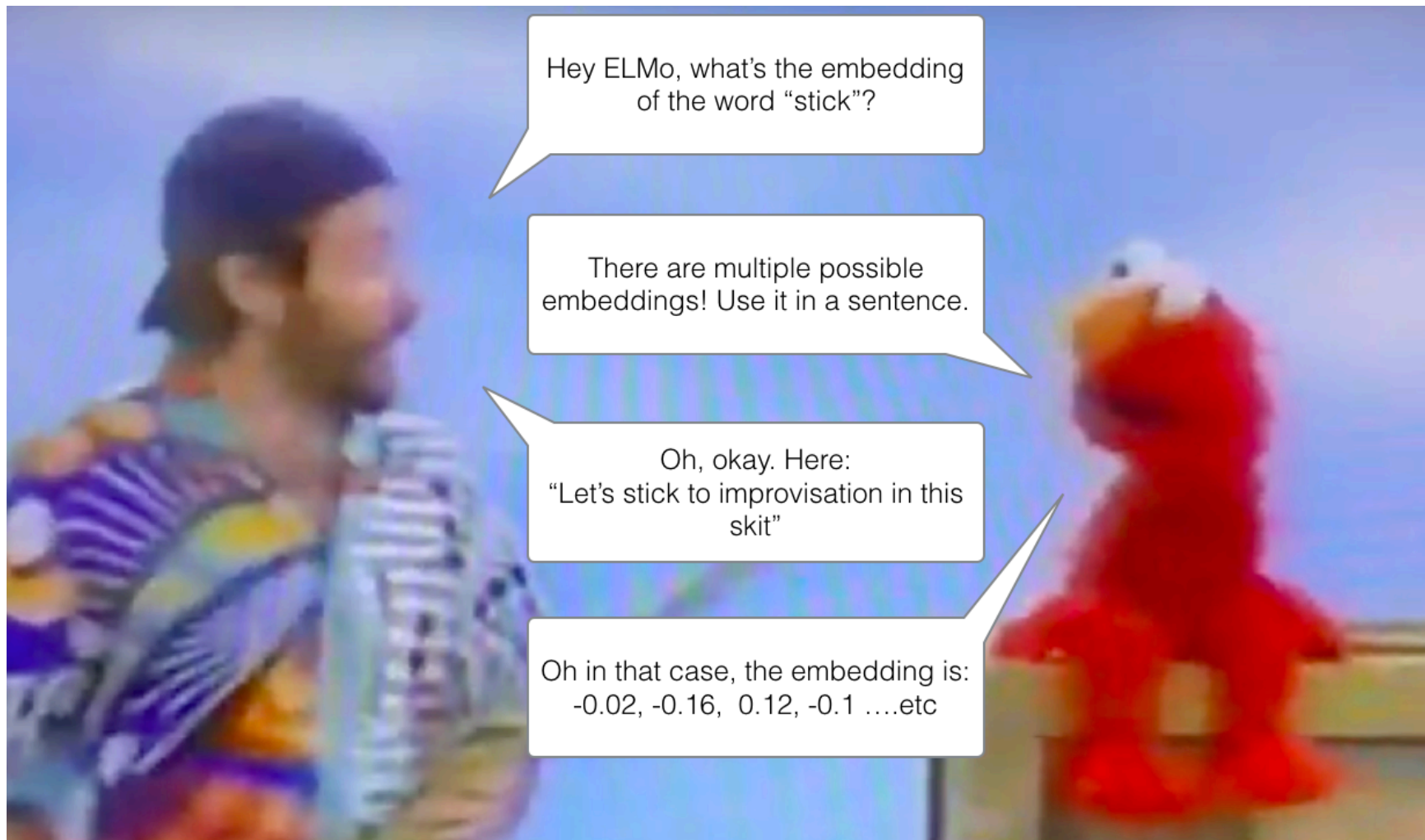
Contextualized Word Representations

Instructor: Yoav Artzi

Overview

- Contextualized word representations
- Models
 - context2vec
 - ELMo
 - BERT

Contextualized Word Representations

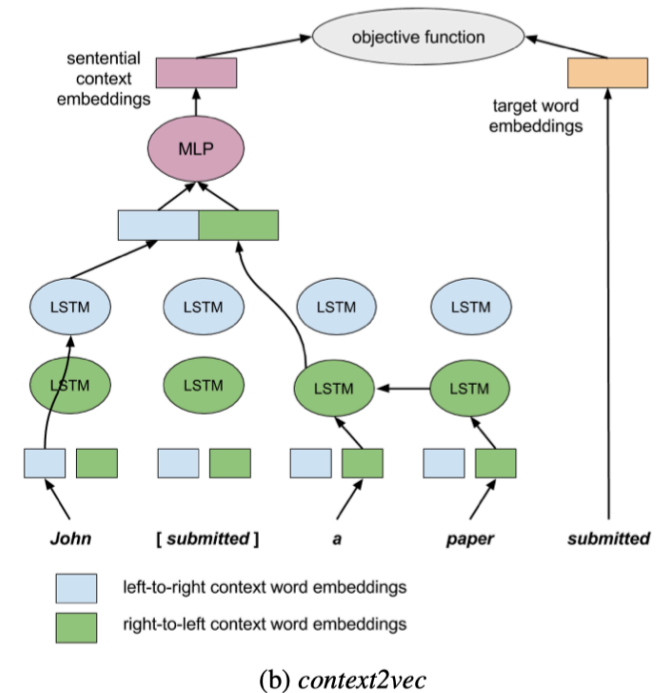
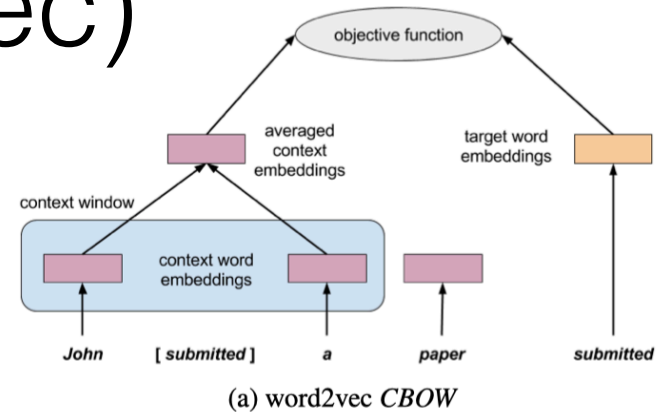


Contextualized Word Representations

- word2vec, GloVe
 - Learn a vector for every word type
 - Always the same vector
- Instead: learn a different vector for word type in every usage
 - But: how do we define the space of uses? Isn't it too large?
 - Solution: use sentence encoders to create a custom vector for every instance of a word

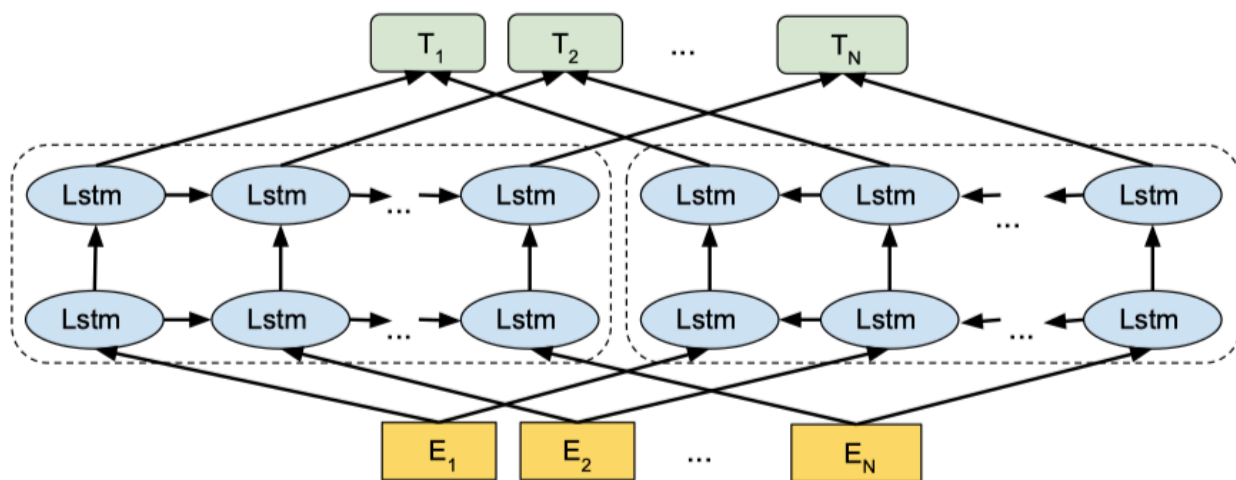
Central Word Prediction Objective (context2vec)

- Model: bi-directional LSTM
- Objective: predict the word given context
- Data: 2B word ukWaC (English data) corpus
- Downstream: use vectors for sentence completion, word-sense disambiguation, etc.



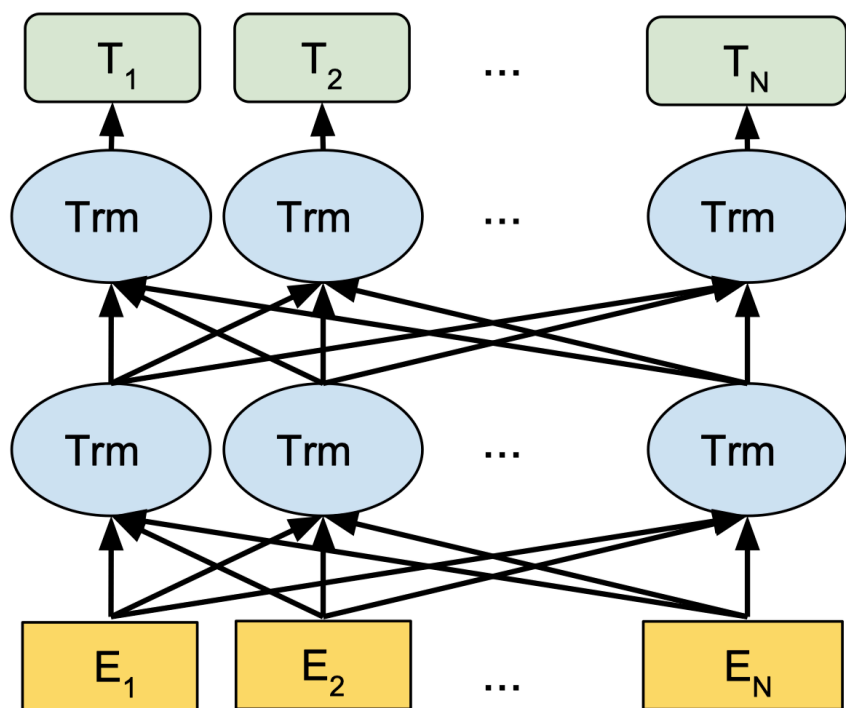
Bi-directional Language Modeling Objective (ELMo)

- Model: multi-layer bi-directional LSTM
- Objective: predict the next word left→right and next word right→left independently
- Data: 1B word benchmark LM dataset
- Downstream: fine-tune the weights of the linear combination of layers per task



Masked Word Prediction (BERT)

- Model: multi-layer self-attention (Transformer), input sentence (or pair w/[CLS] token) and subword representation
- Objective: masked word prediction + next-sentence prediction
- Data: BookCorpus + English Wikipedia
- Downstream: fine-tune weights per task



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\#ing}$	$E_{[SEP]}$
Segment Embeddings	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}



Masked Word Prediction

- Predict a masked word
 - 80%: substitute input word with “[MASK]”
 - 10%: substitute input word with random word
 - 10%: no change
- Like predicting the next word, but adapted for multi-layer self attention

Consecutive Sentence Prediction

- Classify two sentences as consecutive or not
 - 50% of training data is consecutive

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

he bought a gallon [MASK] milk [SEP]

Label = IsNext

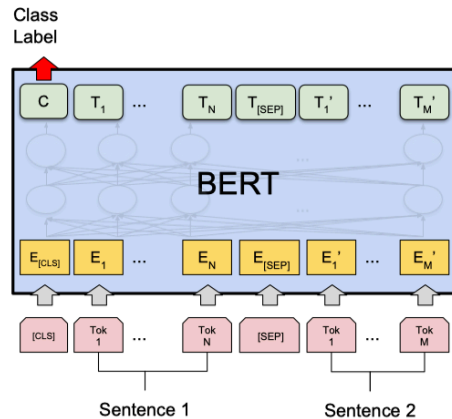
Input = [CLS] the man [MASK] to the store [SEP]

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

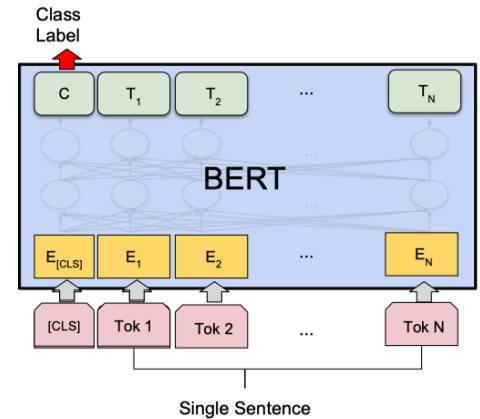
Label = NotNext

Using BERT for Downstream Tasks

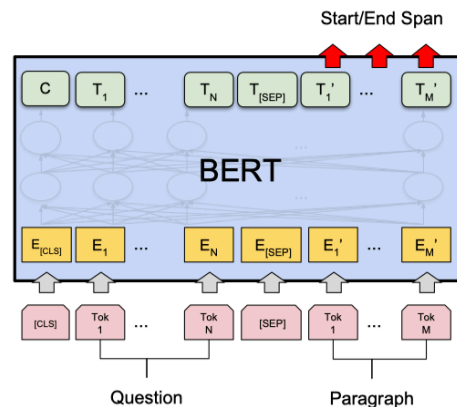
- Use the pre-trained model as the first “layer” of your final model
- Train with fine-tuning using your supervised data



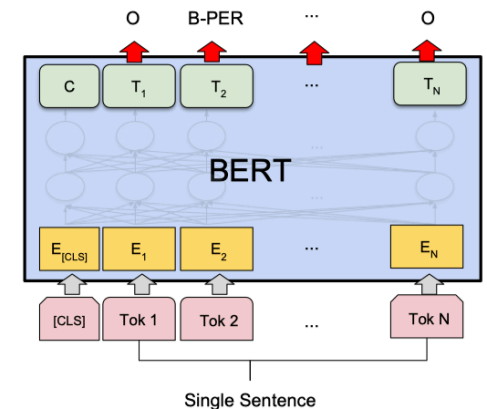
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



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

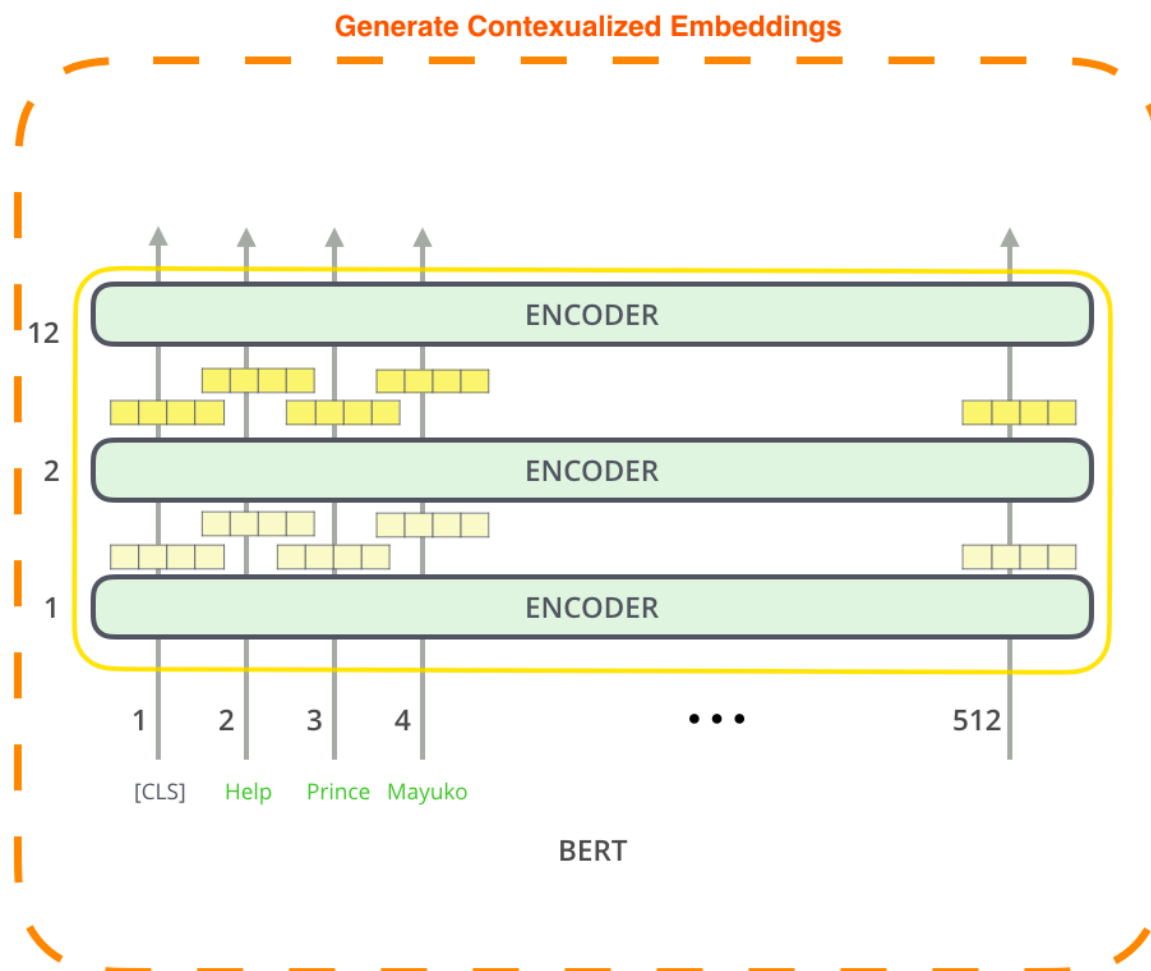


(c) Question Answering Tasks:
SQuAD v1.1

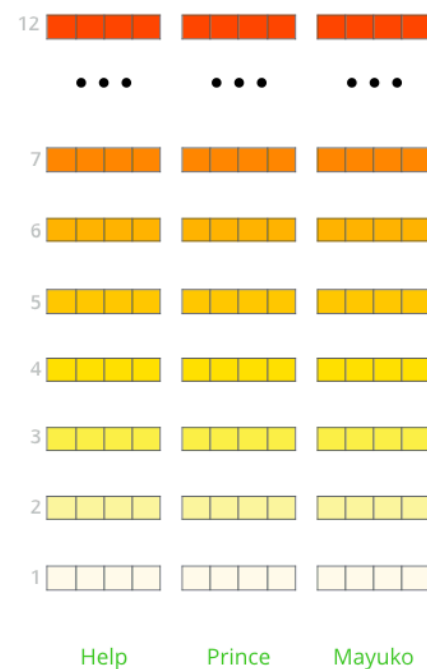


(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Using BERT for Feature Extraction



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?

Using BERT for Feature Extraction

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

