

Challenges in Transfer Learning in NLP

MADRID NLP MEETUP

29th May 2019

Lara Olmos Camarena

*'Except as otherwise noted, this material
Olmos Camarena, Lara (2019). "Challenges in Transfer Learning in NLP"
is licenced under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0
International.*

Contents

Introduction

- Motivation
- Evolution
- First definitions, challenges, frontiers
- Formal definition

Word embeddings

- Word embedding: typical definition
- Pre-trained word embedding
- Conceptual architecture
- Training word embeddings
- Evaluation of word embeddings

Next steps!

- ¿Why? Word embedding limitations and new directions
- Universal Language Model Fine-tuning
- Transformer and BERT

Appendix

- Frameworks related
- NLP&DL: Education
- More references and papers

Introduction

Motivation

Natural Language Processing challenges

word enrichment and training meaningful representation requires so much *effort*

enterprise tasks face with *specific and domain data* of different sizes (difficulties with *small!*)

text annotation is *time-consuming*

Machine Learning models

suffer with vectorial spaces models (*VSM*) built with tokens (*bag of words*) because of it is *high dimensional* and *sparse*.

loss of *sequential and hierarchical information*

Loss of *relations between words and semantics*

Deep Learning and Transfer Learning!

Transferring *information* from one machine learning task to *another*. Not necessary to train from zero, save effort and time.

Transfer learning might involve *transferring knowledge* from the solution of a simpler task to a more complex one, or involve transferring knowledge from a task where there is *more data* to one where there is *less data*.

Most machine learning systems solve a *single* task. Transfer learning is a baby step towards artificial intelligence in which a *single program* can solve *multiple* tasks.

Transfer learning definition from: <https://developers.google.com/machine-learning/glossary/>

More on Transfer Learning! <https://towardsdatascience.com/a-comprehensive-hands-on-guide-to-transfer-learning-with-real-world-applications-in-deep-learning-212bf3b2f27a>

Evolution

- ▶ **Language modelling**, sometimes with **neural network methods**, to solve this problem: **low dimensional vectors** (**embeddings**) for textual units selected: *words, phrases or documents*.
 - ▶ **Words embeddings** obtained as output of different models and algorithms:
 - ▶ Random: *uniform, Xavier*
 - ▶ Predictive models: *word2vec (skip-gram, CBOW)*
 - ▶ Count-based or cooccurrences: *GloVe*
 - ▶ Deep neural network models with different blocks trained: *CNN, RNN, LSTM, biLSTM*
 - ▶ Recent and more complex ones: subword information *FastText*, *biLSTM based ELMo*
 - ▶ New perspectives in transfer learning in NLP:
 - ▶ Use **pre-trained word embedding** to initialize word vectors in embedding layers or other tasks
 - ▶ Use **pre-trained language models** to directly represent text. Prodigy (spaCy), Zero-shot learning for classifiers (Parallel Dots), Universal Sentence Encoder (Google), ULMFit, **BERT** and OpenAI Transformer
- to use later in more complex or **task specific** models/systems.

NEW!

First definitions, challenges, frontiers

Initial context: *Neural network models in NLP, catastrophic forgetting and frontiers*

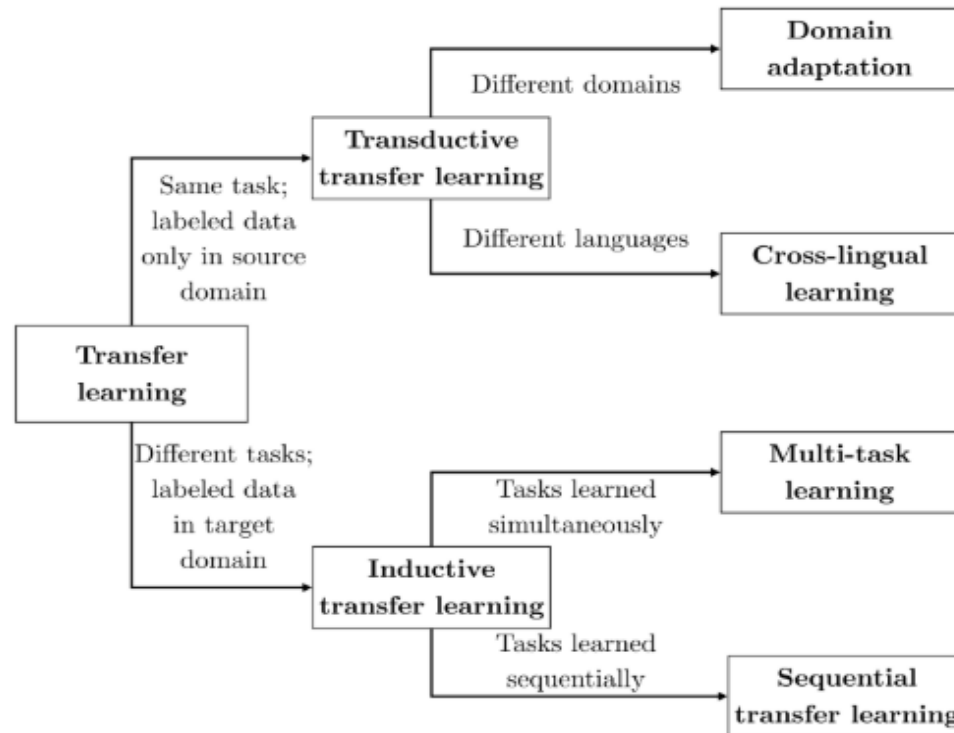
- [A review of the recent history of natural language processing](#) - Sebastian Ruder
 - Neural language models (2001), multi-task learning (2008), Word embeddings (2013), seq2seq (2014),
 - 2018-2019: pretrained language models
- [Fighting with catastrophic forgetting in NLP](#) - Matthew Honnibal (2017)
- [Frontiers of natural language processing](#)

Transfer learning in NLP divulgation

- [Rodríguez, Horacio. Embeddings, Deep Learning, Distributional Semantics en el Procesamiento de la lengua, ... ¿más que marketing?](#)
- [NLP's ImageNet moment](#)
 - **ULMFiT**, **ELMo**, and the **OpenAI** transformer is one key paradigm shift: going from just initializing the first layer of our models to pretraining the entire model with hierarchical representations. If learning word vectors is like only learning edges, these approaches are like learning the **full hierarchy of features**, from edges to shapes to **high-level semantic concepts**.
- [Effective transfer learning for NLP](#) - Madison May - Indico
- [Transfer learning with language models](#) - Sebastian Ruder
- [T3chFest - Transfer learning for NLP and Computer Vision](#) - Pablo Vargas Ibarra, Manuel López Sheriff - Kernel Analytics

Formal definition

- Neural Transfer Learning for Natural Language Processing - Sebastian Ruder



- Previous work found: “*Transfer learning for Speech and Language Processing*”
<https://arxiv.org/pdf/1511.06066.pdf>

Word embeddings

Word embedding: typical definitions

Word embedding is the collective name for a set of **language modeling** and **feature learning** techniques in natural language processing (NLP) where words or phrases from the **vocabulary are mapped to vectors of real numbers**. Conceptually it involves a mathematical embedding from a space with one dimension per word to a continuous vector space with a much **lower dimension**.

Methods to generate this mapping include **neural networks**, dimensionality reduction on the **word co-occurrence matrix**, **probabilistic models**, explainable **knowledge base** method, and explicit representation in terms of the **context** in which words appear.

Word and phrase embeddings, when used as the underlying input representation, have been shown to boost the performance in NLP tasks such as syntactic parsing and sentiment analysis.

(Ref.: Wikipedia)

A word embedding is a learned representation **from text** where words that have the **same meaning** have a **similar representation**.

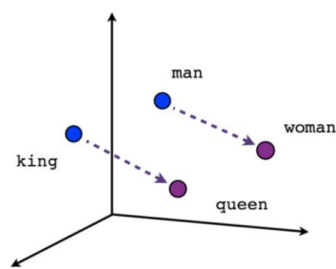
(Ref.: Machine Learning Mastering)

The **Distributional Hypothesis** is that words that **occur in the same contexts tend to have similar meanings** (Harris, 1954). The underlying idea that **"a word is characterized by the company it keeps"** was popularized by Firth (1957), and it is implicit in Weaver's (1955) discussion of **word sense disambiguation** (originally written as a memorandum, in 1949). The Distributional Hypothesis is the basis for **Statistical Semantics**. Although the Distributional Hypothesis originated in Linguistics, it is now receiving attention in Cognitive Science (McDonald and Ramscar, 2001). The origin and theoretical basis of the Distributional Hypothesis is discussed by Sahlgren (2008).

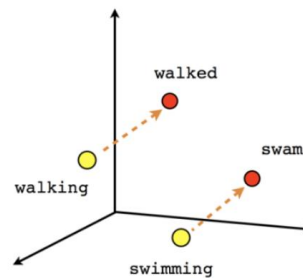
(Ref: ACL Wiki)

Pre-trained word embedding

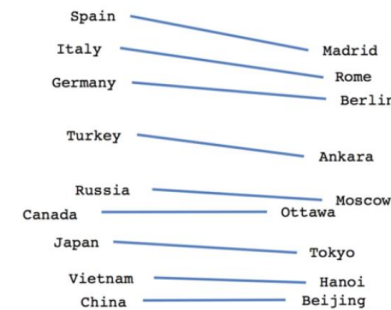
- ▶ **Word embedding** helps to **capture vocabulary** from a *corpus* (f.e.: *public* or *web corpora*) not seen in the task specific language or limited time training and used in general common language.
- ▶ The **representation** allows to establish *lineal relations* and capture *similarities* between words from *syntax and contexts* seen in the **language** and **domain** where the word embedding was trained. **Be careful with bias!**
- ▶ If pre-trained with public and external corpus, some words will be not recognized (**unknown word problem**), when fighting with **jargon** and **argots**! **Also ortographic errors detected!**



Male-Female



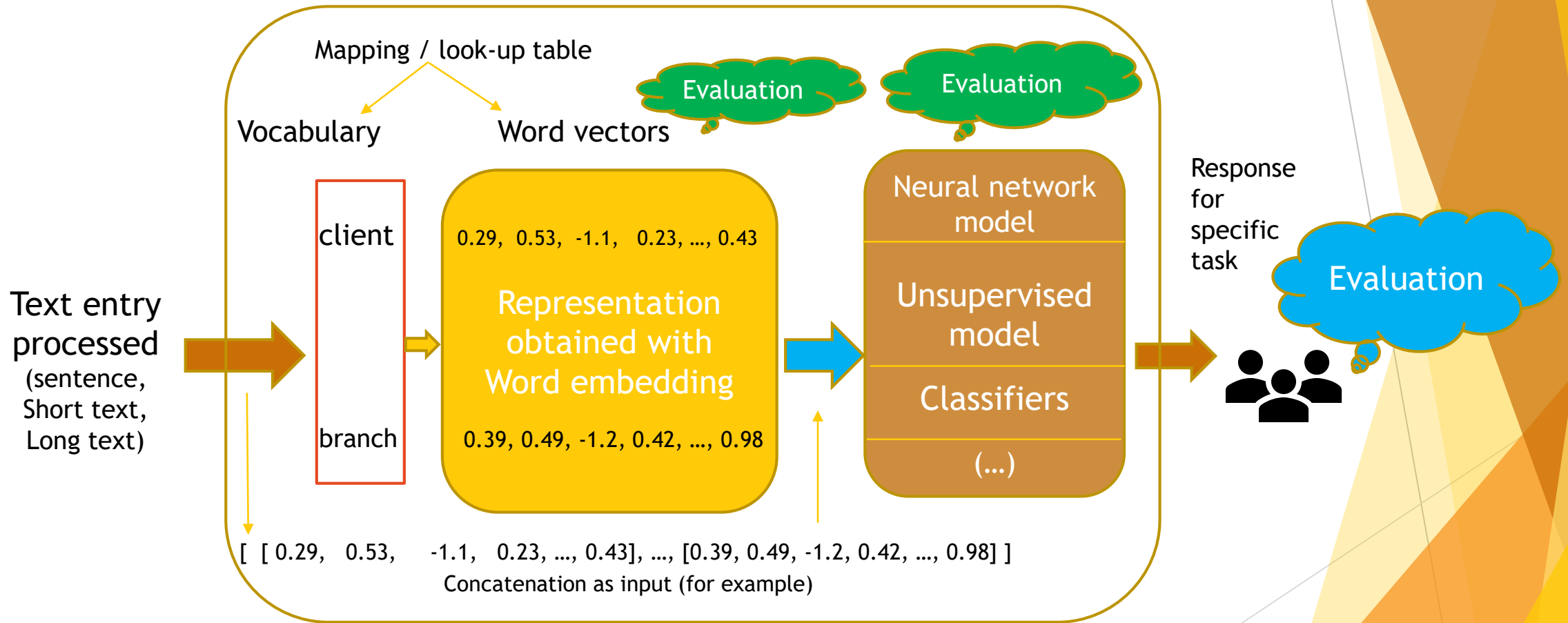
Verb tense



Country-Capital

Conceptual Architecture

(common offline/training and online/predict/execution time)



Training word embeddings

- ▶ Different Word embedding trained for *benchmarking* - perspective of pre-trained initialization of *embedding layer* for our neural network model
 - ▶ Initialization with **random distributions**: uniform, Xavier
 - ▶ Trained with **word2vec**
 - ▶ Trained with **GloVe**
 - ▶ Trained with **Tensorflow** with **NCE** estimation
- ▶ Using different *corpora* for training:
 - ▶ public corpora (example projects like Billion Words)
 - ▶ force-context phrases created for supervised synonyms
 - ▶ domain specific document corpus
 - ▶ real text users/industry dataset



gensim

glovepy 0.0.3



Evaluation of word embeddings

- ▶ **Intrinsic proofs:** discovering internal word synonyms. *Synonyms trained, are still synonyms after embedding? Which embedding improves similarity for our task?*
 - ▶ **Qualitative** evaluation: word neighbours clusters visualization (t-SNE and PCA), Brown clustering
 - ▶ **Quantitative** evaluation: similarity test with cosine function and threshold using
- ▶ **Extrinsic proofs:** testing Word embedding improvement when changing initialization for fix in specific task model.
 - ▶ Dataset, domain and task dependent
 - ▶ NLP processing pipeline must match the input sentence and word embedding construction (stopword removal, normalization, stemming...)

MORE?

Enso Benchmark - <https://github.com/IndicoDataSolutions/Enso>

Next steps!

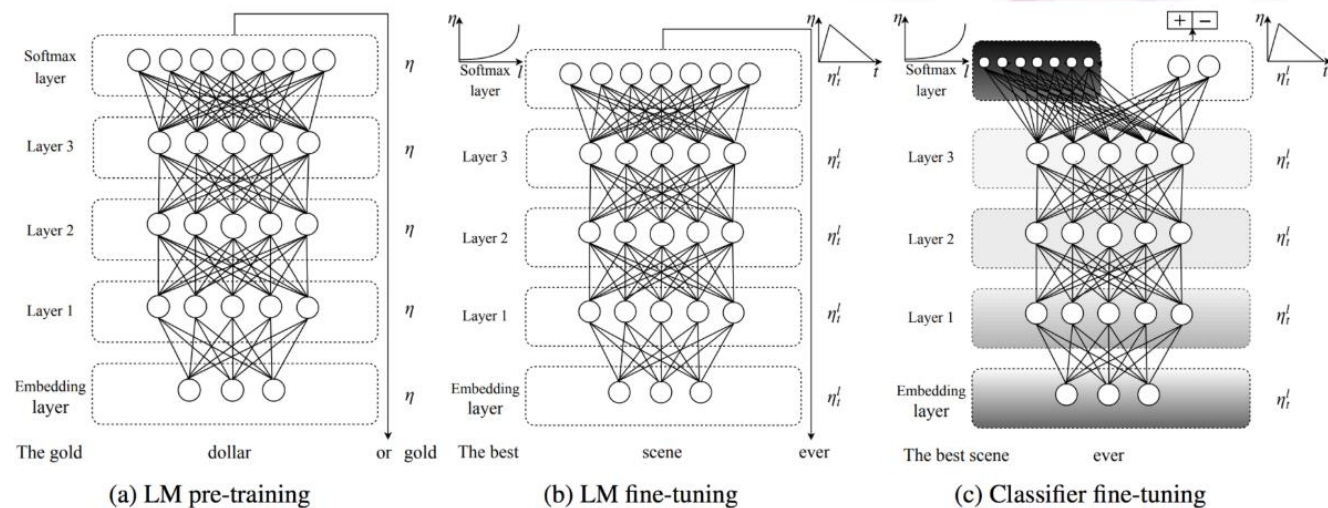
¿Why? Word embedding limitations

- ▶ Handling out of vocabulary words or unknown words
 - ▶ Solved with subword level embeddings (character and n-gran based based): [FastText embedding](#)
- ▶ Unique representation for word, problem with polysemy and word senses!
 - ▶ Solved with **Elmo**: [Deep contextualized word representations](#)
- ▶ Word embedding domain adaptation
- ▶ Lack of studies to solve bias and vectorial space studies, also in evaluation

New directions and challenges

- ▶ Research transition from *feature-based* to *fine-tuning* complex pre-trained Deep Learning architectures for NLP principal tasks: classification, inference, question answering, natural language generation...
 - ▶ Started in language modelling: sentence embedding like SpaCy, Universal Sentence Encoder (DAN and Transformer)
 - ▶ **State of the art**: ULMFit, Transformer, BERT, Open AI GPT...
- ▶ *Challenges with* Deep Learning difficulties: explainability of neural network models, its representation knowledge learned and to use and adapt in real business problems [ethically](#)!

Universal Language Model Fine-tuning



- ULMFiT consists of three phases:
 - a) Train language model (LM) on general domain data.
 - b) Fine-tune LM on unlabeled target data.
 - c) Train classifier on top of LM on labeled target data.

(Howard & Ruder, ACL 2018)

[Neural Transfer Learning for Natural Language Processing](#) - Sebastian Ruder

Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." arXiv preprint arXiv:1801.06146 (2018).

Transformer

BERT

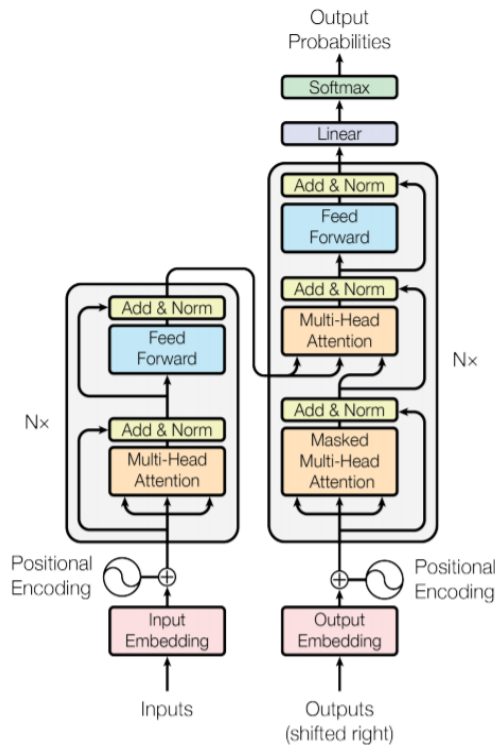


Figure 1: The Transformer - model architecture.

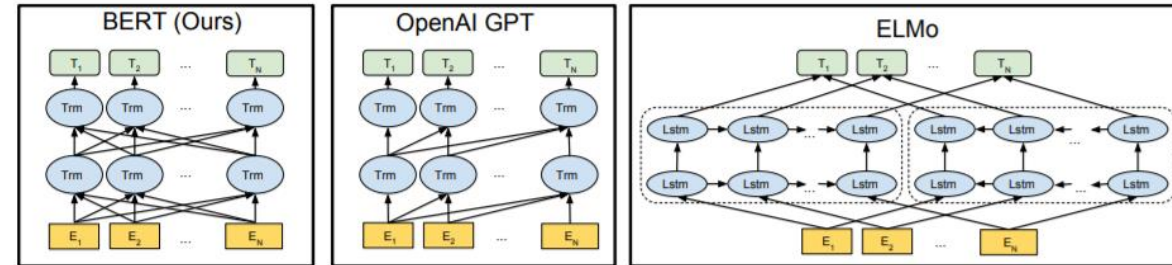


Figure 1: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks. Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

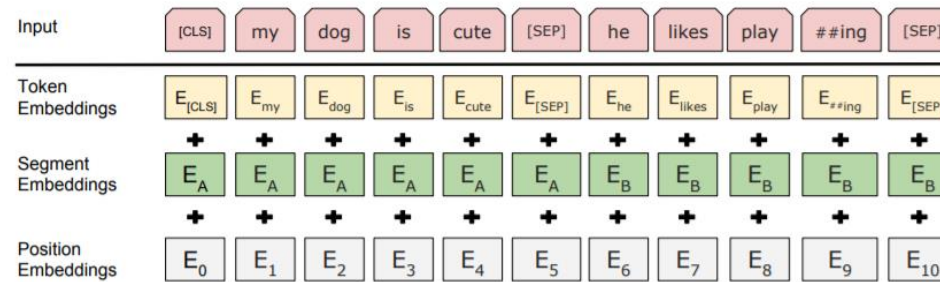


Figure 2: BERT input representation. The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

<https://ai.google/research/pubs/pub46201>

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008). <https://arxiv.org/pdf/1706.03762.pdf>

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. <https://arxiv.org/pdf/1810.04805.pdf>

TO START:

<http://jalammar.github.io/illustrated-bert/>

Conclusions

Thank you for your
attention!

Any questions?

Appendix

Frameworks related

- ▶ Gensim: <https://radimrehurek.com/gensim/models/word2vec.html>
- ▶ Glovepy: <https://pypi.org/project/glovepy/>
- ▶ SpaCy: <https://spacy.io/>
- ▶ FastText - Facebook Research: <https://github.com/facebookresearch/fastText>
- ▶ Universal Sentence Encoder - Tensorflow:
<https://tfhub.dev/google/universal-sentence-encoder/2>
- ▶ AllenNLP - ELMo - <http://aclweb.org/anthology/N18-1202>
 - ▶ <https://allennlp.org/elmo>
- ▶ Enso - <https://github.com/IndicoDataSolutions/Enso>
- ▶ BERT - <https://github.com/google-research/bert>

NLP&DL: Education

- ▶ **Stanford University** - Christopher Manning - Natural Language Processing with Deep Learning

- ▶ https://www.youtube.com/watch?v=OQQ-W_63UgQ&list=PL3FW7Lu3i5Jsnh1rnUwq_TcylNr7EkRe6
- ▶ <http://cs224d.stanford.edu/>
- ▶ <http://web.stanford.edu/class/cs224n/>



- ▶ **University of Oxford** - [Deep Natural Language Processing](#)

- ▶ <https://github.com/oxford-cs-deepnlp-2017/lectures>



- ▶ **Bar Ilan University's** - Yoav Goldberg - Senior Lecturer [Computer Science Department](#) NLP Lab - <http://u.cs.biu.ac.il/~nlp/>



- ▶ *Book: Neural Network Methods for Natural Language Processing*

- ▶ **Universitat Politècnica de Catalunya** - Horacio Rodríguez - Embeddings working notes - <http://www.cs.upc.edu/~horacio/docencia.html>
<http://www.lsi.upc.edu/~ageno/anlp/embeddings.pdf>

More references and papers

► Academic material:

- Ruder, Sebastian. On word embeddings. <http://ruder.io/word-embeddings-2017/>, <http://ruder.io/word-embeddings-1/>
- Jurafsky, Daniel, & Martin, James H (2017). *Speech and Language Processing. Draft 3^o Edition. Chapter 15-16. Vector Semantics and Semantics with Dense Vectors*
- Goldberg, Yoav (2017). *Neural Network Methods for Natural Language Processing (Book)*

► Word embedding papers:

- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." *Advances in neural information processing systems*. 2013.
- Le, Quoc, and Mikolov, Tomas. "Distributed representations of sentences and documents." *International conference on machine learning*. 2014.
- Pennington, Jeffrey, Richard Socher, and Christopher Manning. "GloVe: Global vectors for word representation." *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014.
- Mnih, Andriy, and Koray Kavukcuoglu. "Learning word embeddings efficiently with noise-contrastive estimation." *Advances in neural information processing systems*. 2013.
- Schnabel, Tobias and Labutov, Igor (2015). *Evaluation methods for unsupervised word embeddings*.
- Wang, Yanshan and Lui, Sijia (2018). *A Comparison of Word Embeddings for the Biomedical Natural Language Processing*. arXiv:1802.00400v3

► Next steps!

- <http://jalammar.github.io/illustrated-bert/>
- Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." arXiv preprint arXiv:1801.06146 (2018).
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.