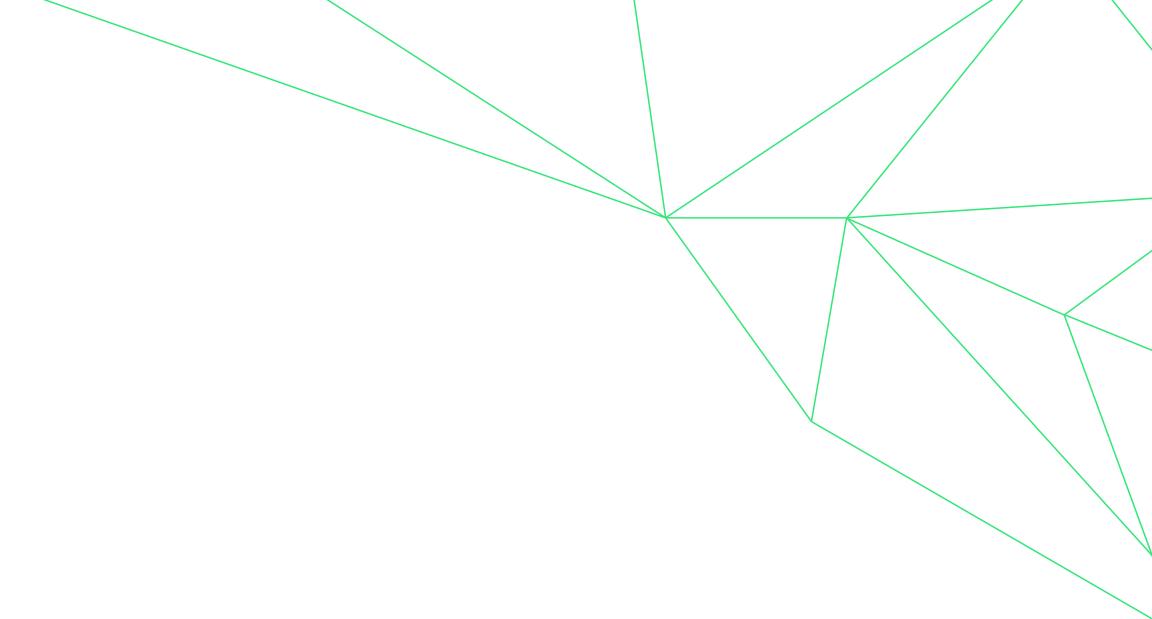




prometeia



How do I embed a sequence of words? From context-free to contextual embeddings

Fabio Fumarola

Milan, July 18 2019



- **Language and Machine Learning**
- **Distributional Semantics**
- **Context-free word embeddings:**
 - **Count based**
 - **Context based**
- **Contextual embeddings**
- **Final Considerations & Conclusions**



- **Language and Machine Learning**
- Distributional Semantics
- Context-free word embeddings:
 - Count based
 - Context based
- Contextual embeddings
- Final Considerations & Conclusions



prometeia

- Language and Machine Learning
- **Distributional Semantics**
- Context-free word embeddings
 - Count based
 - Context based
- Contextual embeddings
- Final Considerations & Conclusions

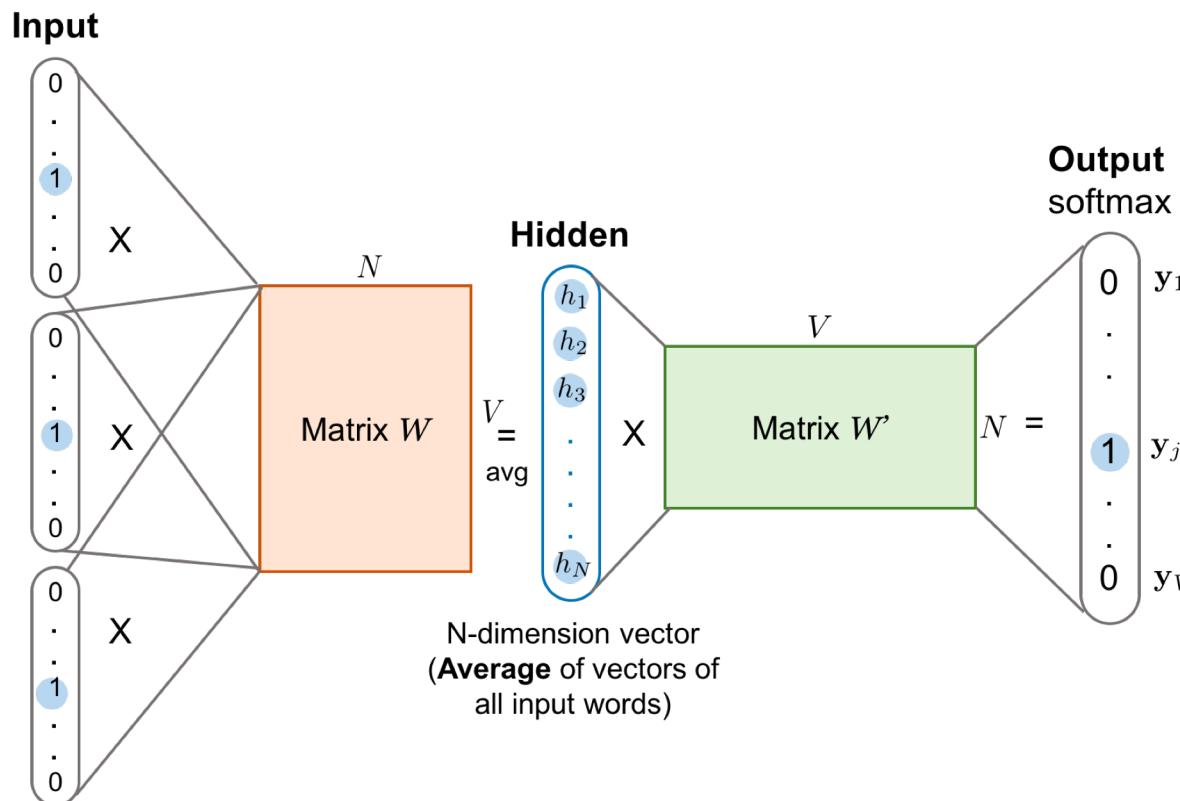


prometeia

- Language and Machine Learning
- Distributional Semantics
- **Context-free word embeddings:**
 - Count based
 - Context based
- Contextual embeddings
- Final Considerations & Conclusions

Continuous Bag-of-Words (CBOW)

- The Continuous Bag-of-Words (Mikolov et al., 2013) predicts the target word from the surrounding window



CBOW model is better for small dataset

Loss function: Negative Sampling

- Considering an input word w_I , and the correct word w
- We sample N other words from a noise distribution $\hat{w}_1, \hat{w}_2, \dots, \hat{w}_N \sim Q$.

$$\mathcal{L}_{\theta} = -[\log \sigma({v'_w}^T v_{w_I}) + \sum_{\substack{i=1 \\ \tilde{w}_i \sim Q}}^N \log \sigma(-{v'_{\tilde{w}_i}}^T v_{w_I})]$$

Distance
input target embedding Distance
input noise embedding

1. Get close to 0 when input and target embedding are equal
2. Get close to zero when the input embedding is far from noise word embeddings

Global Vectors (GloVe)

<https://nlp.stanford.edu/pubs/glove.pdf>

- Proposed by Pennington et al. in 2014, combines count-based matrix factorization to skip-gram
- Word meanings are captured by the ratios of co-occurrence probabilities

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

- large values correlate well with ice, and small values correlate well with steam
- *The ratio of probabilities encodes some crude form of meaning associated with the abstract concept*

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2$$

FastText: Enriching Word Vectors with Sub-word Information

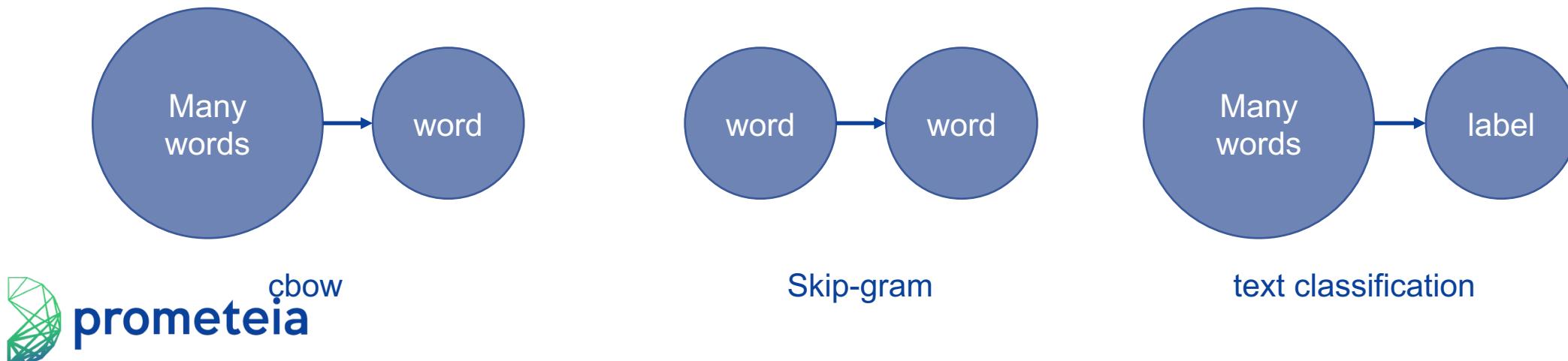
<https://arxiv.org/pdf/1607.04606v1.pdf>

- Word2Vec has several drawbacks:
 - **No sentence representations:** taking the average pre-trained word vector is popular but does not work always well
 - **Not exploiting morphology:** words with the same radicals are learned as different vectors

disastrous / disaster

mangera / mangerai

FastText is an extension of word2vec to take into account the morphology of words



FastText: Idea

Represent each word as a bag of *character n-grams*

1. A vector representation is associated to each character n-gram
2. Words are represented as the sum of these representations

Example: “matter”, with n=3 is represented as

^ma, mat, att, tte, ter, er\$

Where ^ and \$ are special symbols to distinguish the n-grams from words in the vocabulary

- in this way compounds name and out of vocabulary words are easy to model

FastText: Features

- Features of a word computed using n-grams

$$h_w = \sum_{g \in w} x_g$$



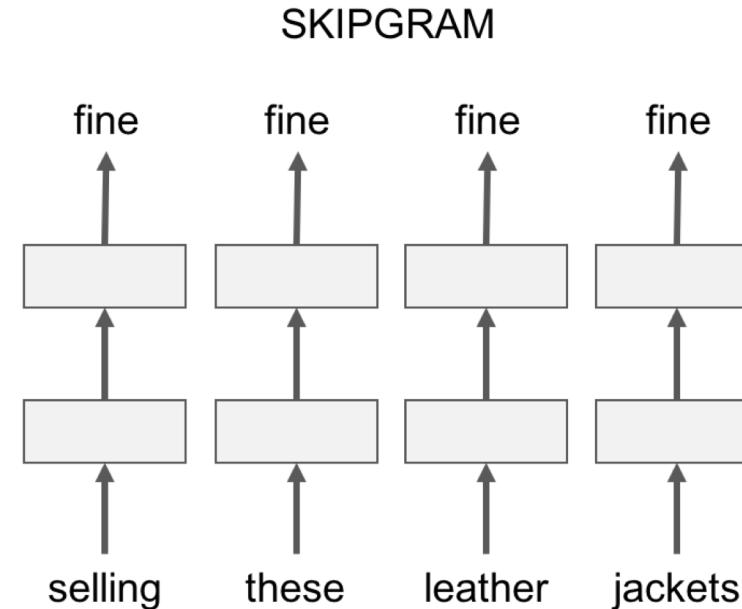
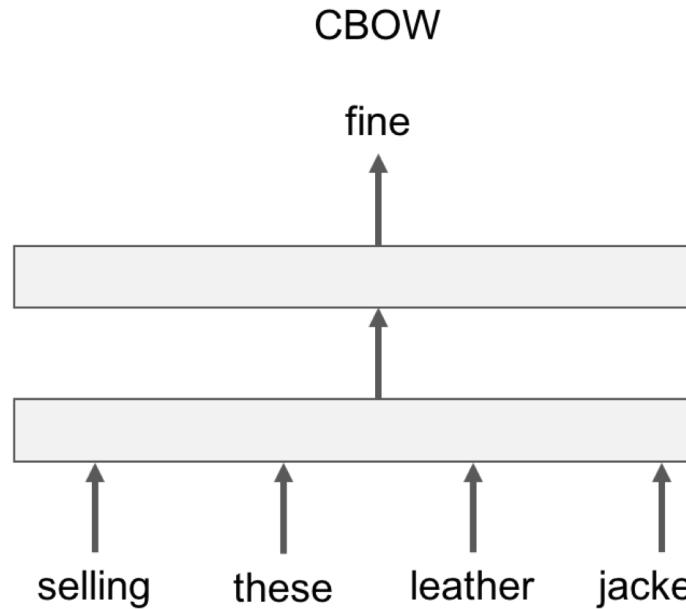
- Out Of Vocabulary (OOV) words representations

$$h_w = \sum_{g \in w} x_g$$



FastText: Model

- the features preprocessed can be used to train a CBOW or Skip-Gram (this works better) model



I am selling these fine leather jackets

FastText: Extensions

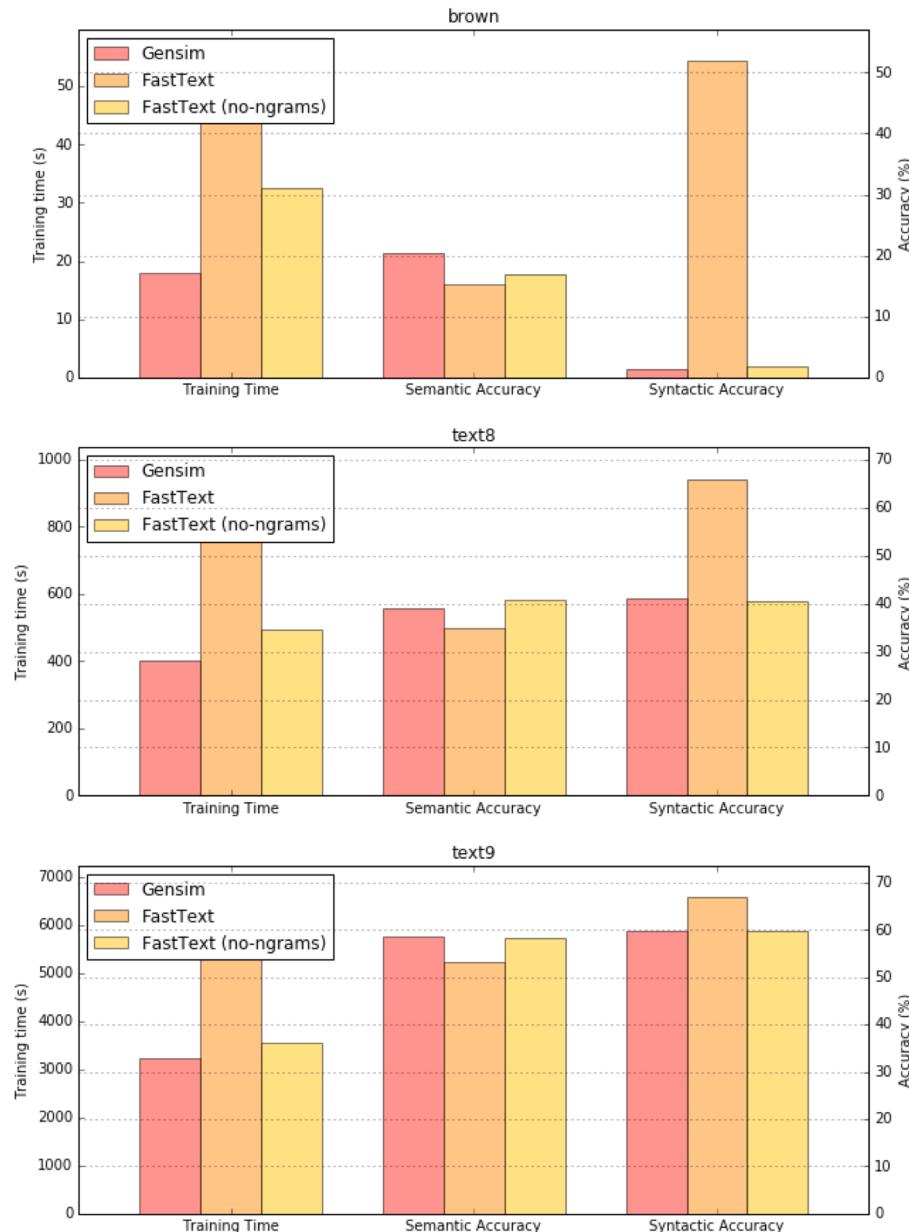
It can be used for:

1. **Classification**: A. Joulin, E. Grave, P. Bojanowski, T. Mikolov, Bag of Tricks for Efficient Text Classification, <https://arxiv.org/abs/1607.01759>

2. **Translation**: A. Joulin, P. Bojanowski, T. Mikolov, H. Jegou, E. Grave, Loss in Translation: Learning Bilingual Word Mapping with a Retrieval Criterion, <https://aclweb.org/anthology/D18-1330>

Word2Vec vs. FastText: Results Comparison

1. fastText with n-grams do better on syntactic tasks
2. Gensim word2vect and fastText with no n-grams do slightly better on semantic task (due to the test dataset)
3. In general, the performance of the models seems to get closer increasing the corpus size
4. The semantic accuracy increase with the corpus size
5. The syntactic accuracy slowly increase with the corpus size



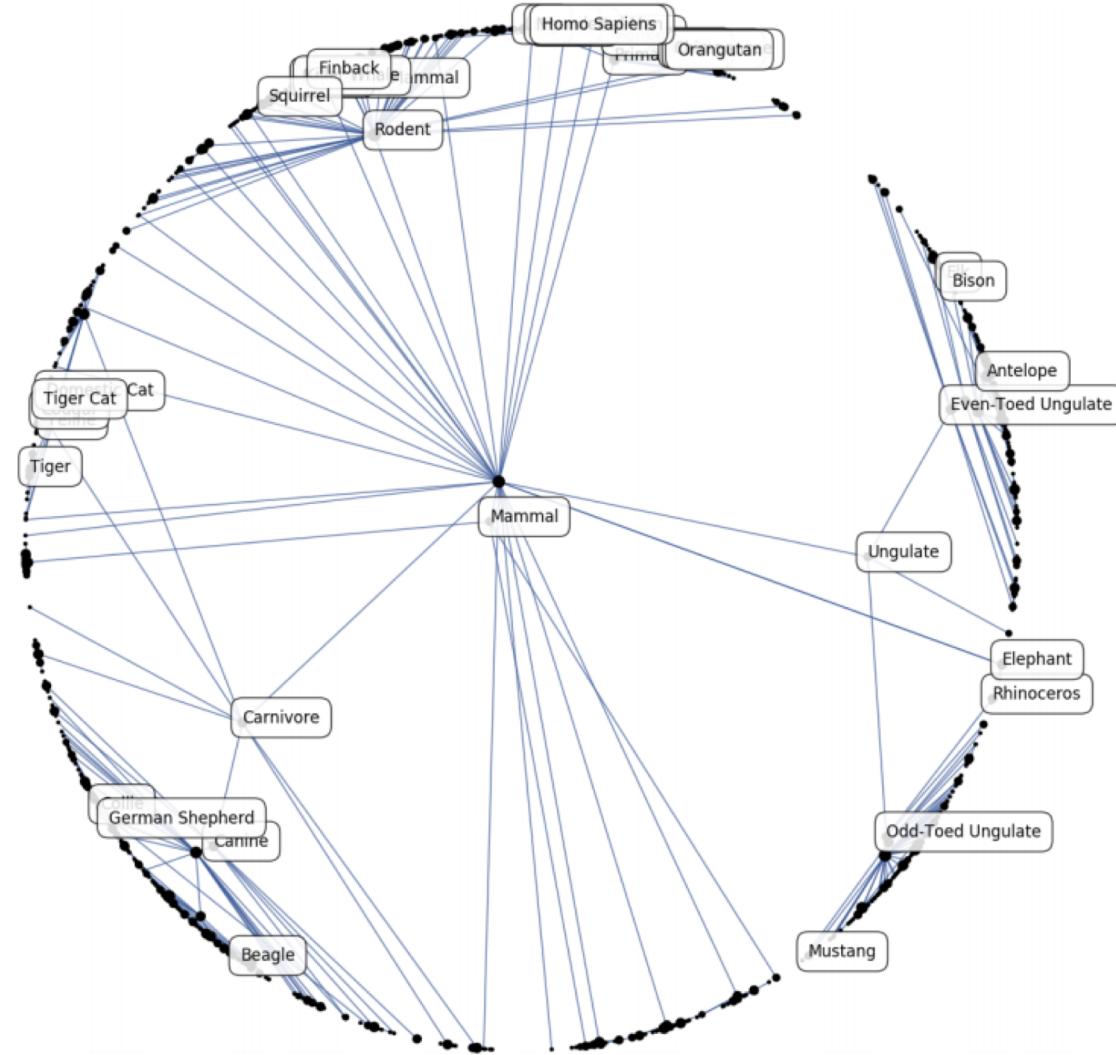
Recap: Count vs. Context

The main difference between the approaches is :

- **Documents as context:** for LSA and topic models (legacy from Information Retrieval)
 - *The document-based:* semantic relatedness (e.g. ‘boat’ – ‘water’)
- **Words as context:** for word skip-gram, CBOW, glove and FastText (linguistic and cognitive perspective)
 - *The word-based:* semantic similarity (e.g. ‘boat’ - ‘ship’)

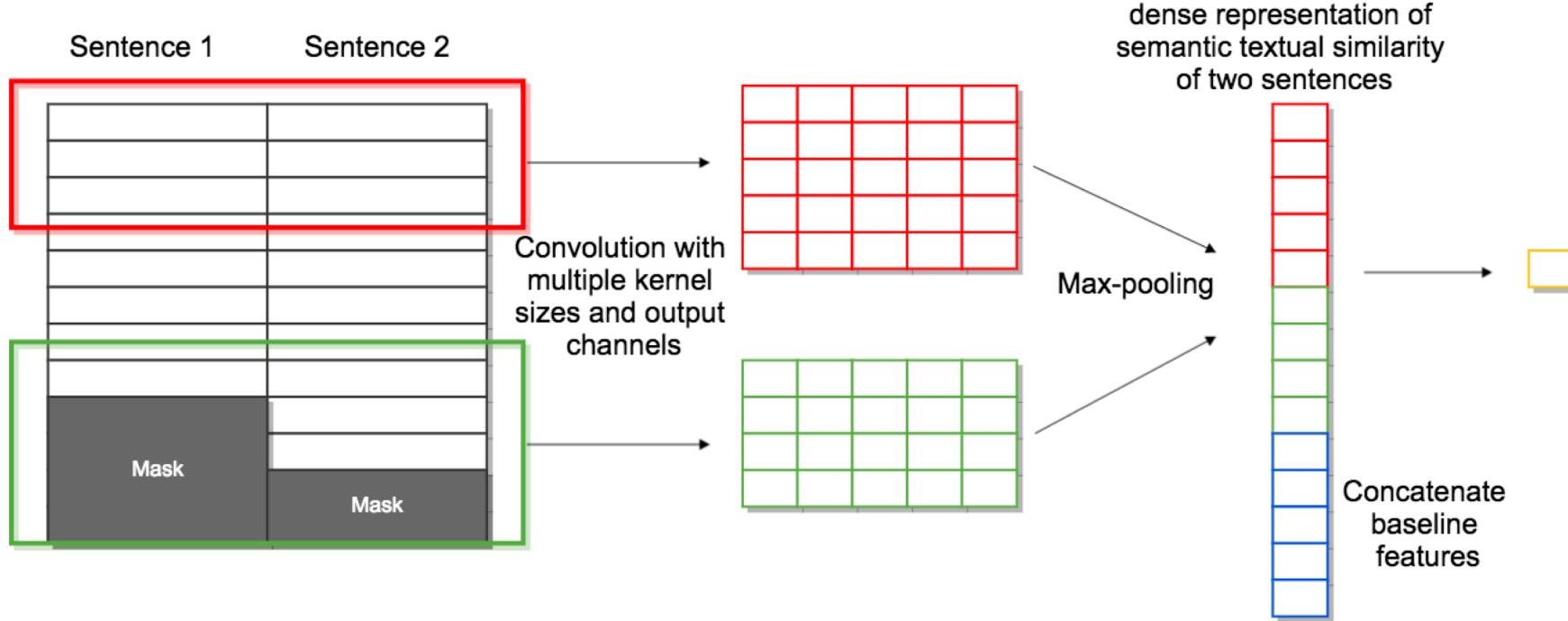
Interesting extensions not covered

- Poincaré Embeddings
(<https://arxiv.org/pdf/1705.08039.pdf>)



Interesting extensions not covered

- Sent2Vec Embeddings NAACL 2018 (<https://aclweb.org/anthology/N18-1049>)



- An extension of FastText and CBOW to sentences

- Language and Machine Learning
- Distributional Semantics
- Context-free word embeddings:
 - Count based
 - Context based
- **Contextual embeddings**
- Final Considerations & Conclusions

Contextualized Embeddings

- The introduced embeddings suffer of a major problem they cannot model ***polysemy***

- In the sentences:

- “*I have an Apple phone*”, and
 - “*I am eating an apple*”

the two “apple” words refer to very **different things** but share the **same embedding** vector

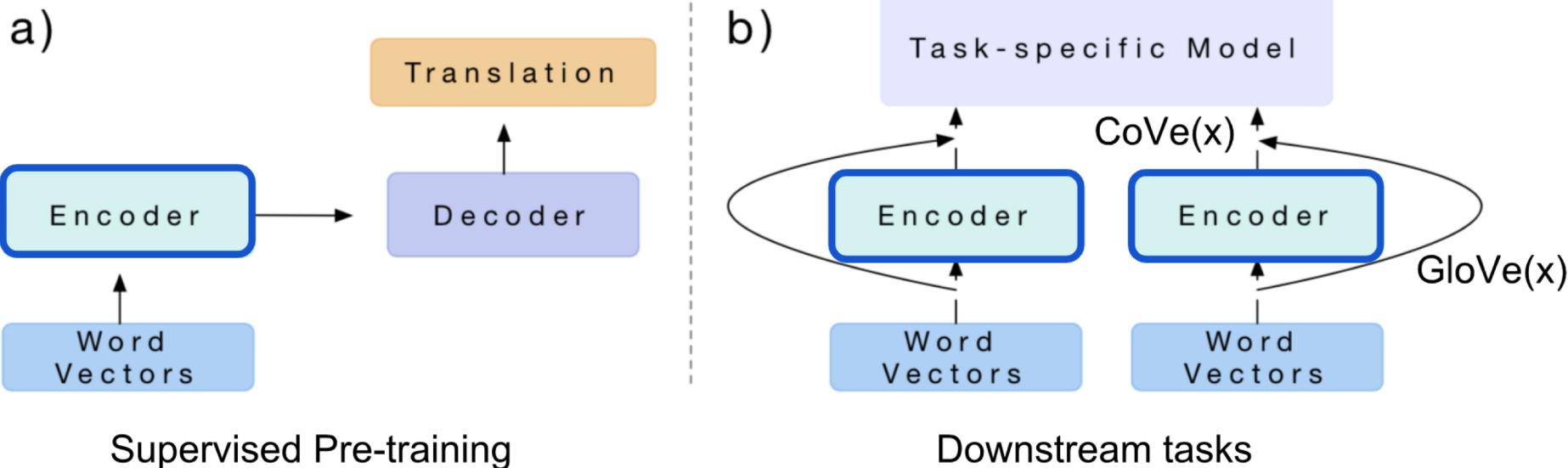
- **Solution:** Contextual Embeddings

Contextualized Embeddings

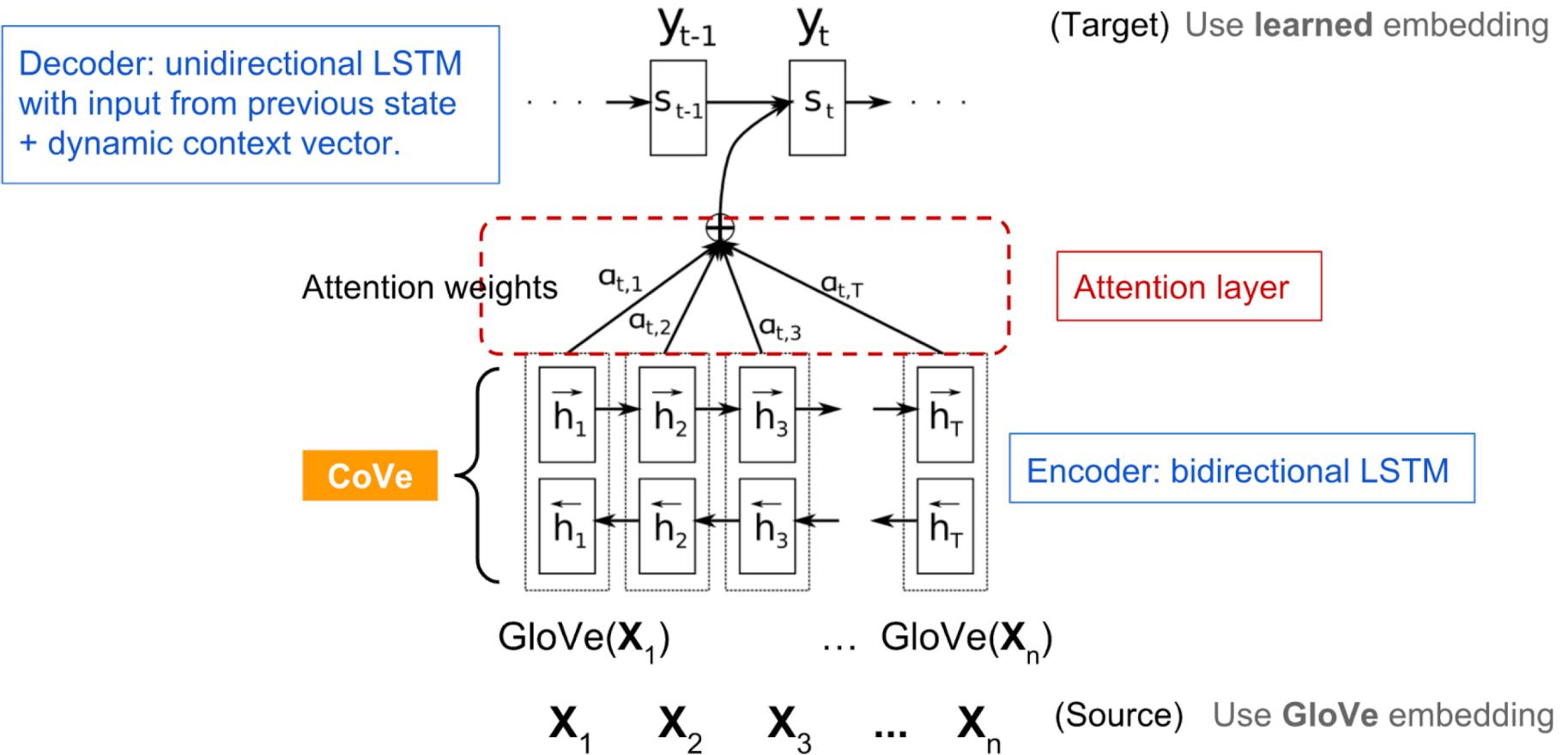
- CoVe
- ELMO
- ULMFit
- OpenAI GPT
- BERT
- OpenAI GPT 2

Contextual Word Vectors (CoVe)

- Learns contextual word embeddings by using the encoder in an attentional sequence to sequence model trained for machine translation
- CoVe word representations are function of the entire input sequence



Contextual Word Vectors (CoVe)



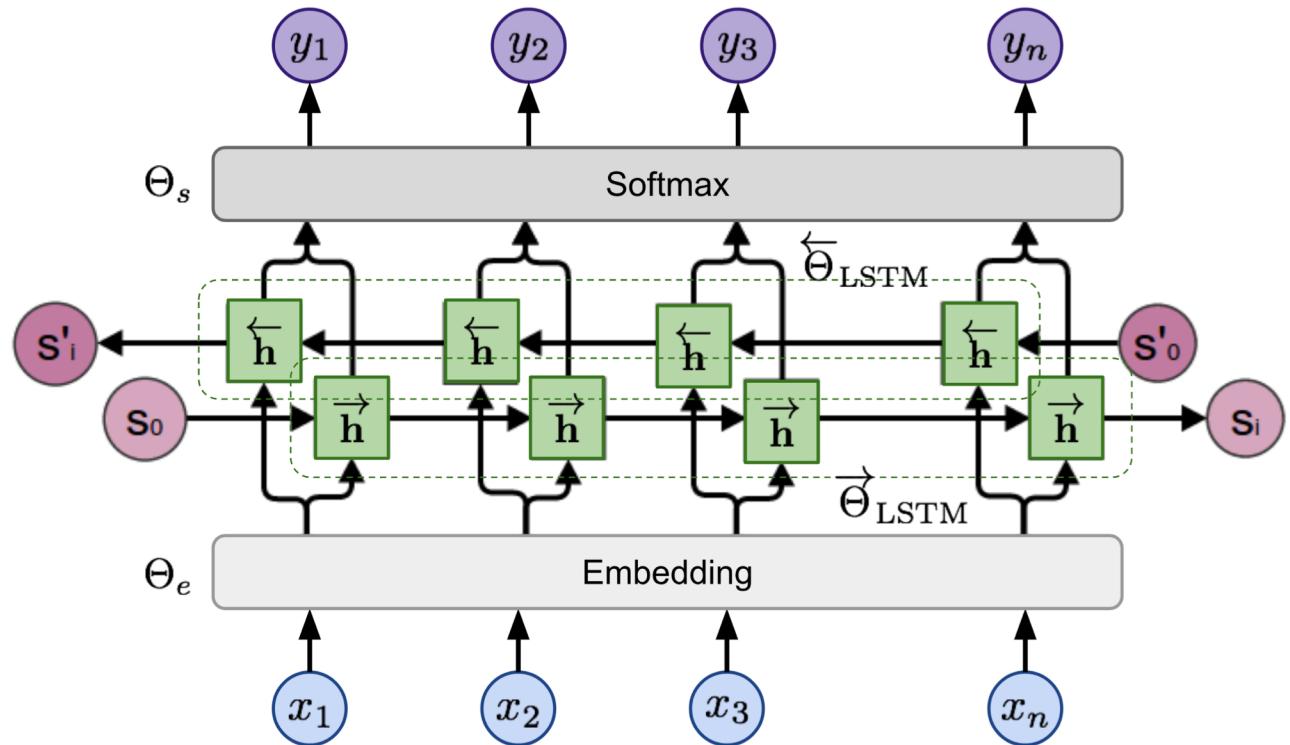
CoVe: Limitations

1. Pre-training is bounded by datasets for supervised translation
2. Limited contribution of CoVe when used as layer for other tasks

Embeddings from Language Model (ELMo)

Peters, Matthew E. and Neumann, Mark and Iyyer, Mohit and Gardner, Matt and Clark, Christopher and Lee, Kenton and Zettlemoyer, Luk
Deep contextualized word representations NAACL 2018

- Learns contextualized embeddings by pre-training a language model that learns to predict the next token given the history



ELMo Representations

- **Contextual:** the representation of each word depends on the sentence in which it is used
- **Deep:** the word representation is computed by stacking the internal states of all the layers
bidirectional language model

$$R_i = \{\mathbf{h}_{i,\ell} \mid \ell = 0, \dots, L\} \text{ where } \mathbf{h}_{0,\ell} \text{ is the embedding layer output and } \mathbf{h}_{i,\ell} = [\vec{\mathbf{h}}_{i,\ell}; \tilde{\mathbf{h}}_{i,\ell}].$$

Different layers encodes different kind of information (e.g. PoS tagging from lower layers and word-sense disambiguation in higher levels)

- **Character based:** the representation is based on characters rather than words (can take advantage of sub-words to compute ...)

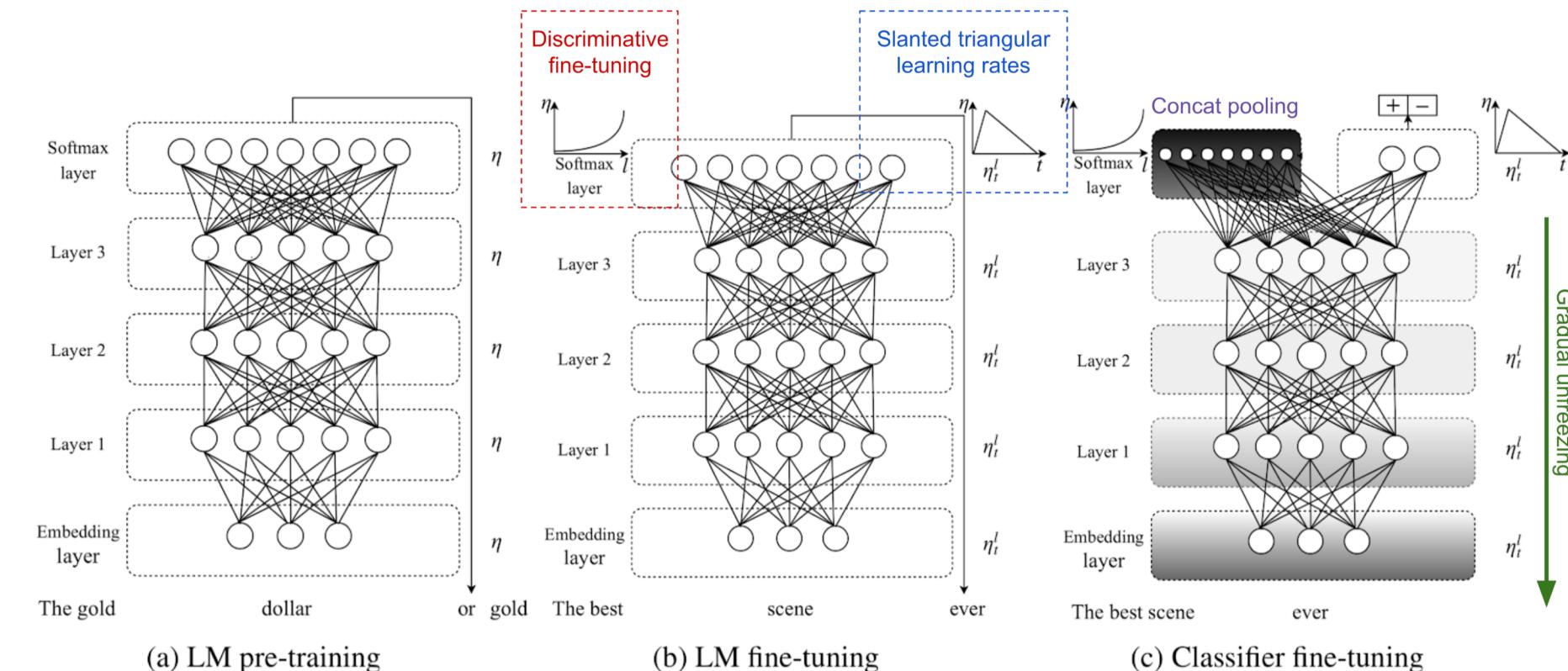
ELMo: Key Results

Task	Previous SOTA		Our baseline	ELMo + Baseline	Increase (Absolute/Relative)
SQuAD	SAN	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al (2017)	88.6	88.0	88.7 +/- 0.17	0.7 / 5.8%
SRL	He et al (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al (2017)	91.93 +/- 0.19	90.15	92.22 +/- 0.10	2.06 / 21%
Sentiment (5-class)	McCann et al (2017)	53.7	51.4	54.7 +/- 0.5	3.3 / 6.8%

Tasks and Datasets

- SQuAD: Stanford Question Answering Dataset
- SNLI: Stanford Natural Language inference (entailment, contradiction, neutral)
- SRL: Semantic Role Labeling
- Coref: Clustering mentions in text that refer to the same underlying real world entities
- NER: Named Entity Extraction
- Sentiment: The fine-grained sentiment classification task in the Stanford Sentiment Treebank

Universal Language Model Fine-tuning (ULMFiT)



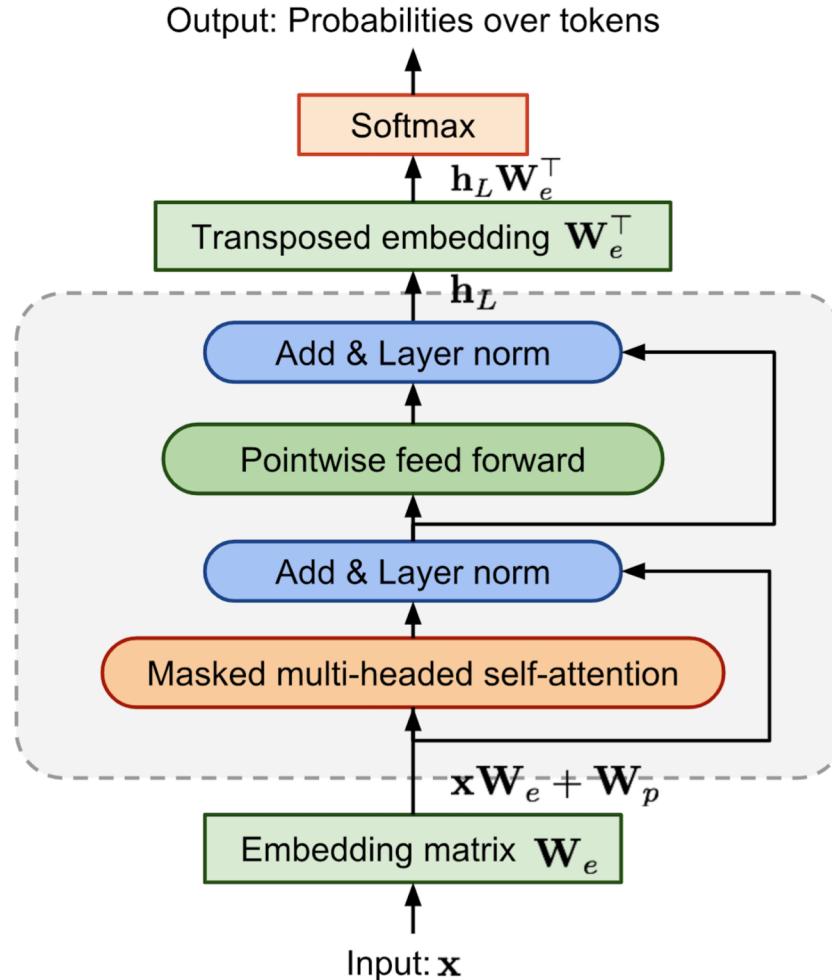
Fine-tuning: (i) different learning rate per layer (ii) scheduled increase decrease of the learning rate

OpenAI Generative Pre-training Transformer

GPT has two major differences from ELMo:

Model	ELMo	GPT
Architecture	BiLSTM	Transformer
Type of Embeddings	Customized per task	Same embeddings fine tuned for al tasks

GPT: Transformer Architecture



Transformer Block
Repeat $\times L=12$

$$\begin{aligned}\mathbf{h}_\ell &= \text{transformer_block}(\mathbf{h}_{\ell-1}) \\ \ell &= 1, \dots, L\end{aligned}$$

GPT: Byte Pair Encoding

<https://arxiv.org/abs/1508.07909>

- Is a data compression technique that iteratively replaces the most frequent pair of symbols (originally bytes) in a given dataset with a single unused symbol.
- In each iteration, the algorithm finds the most frequent (adjacent) pair of symbols, each can be constructed of a single character or a sequence of characters, and merged them to create a new symbol
- It adopted to solve the open-vocabulary issue

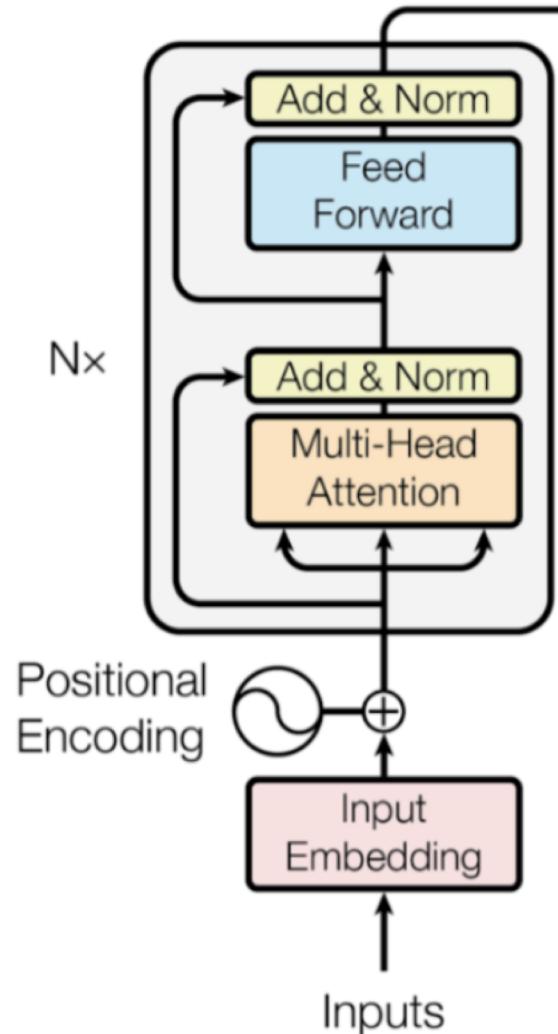
Bidirectional Encoder Representations from Transformer (BERT)

- Compared to GPT the largest difference is that BERT is Bidirectional

Contextual Embedding example: “*I made a bank deposit*”

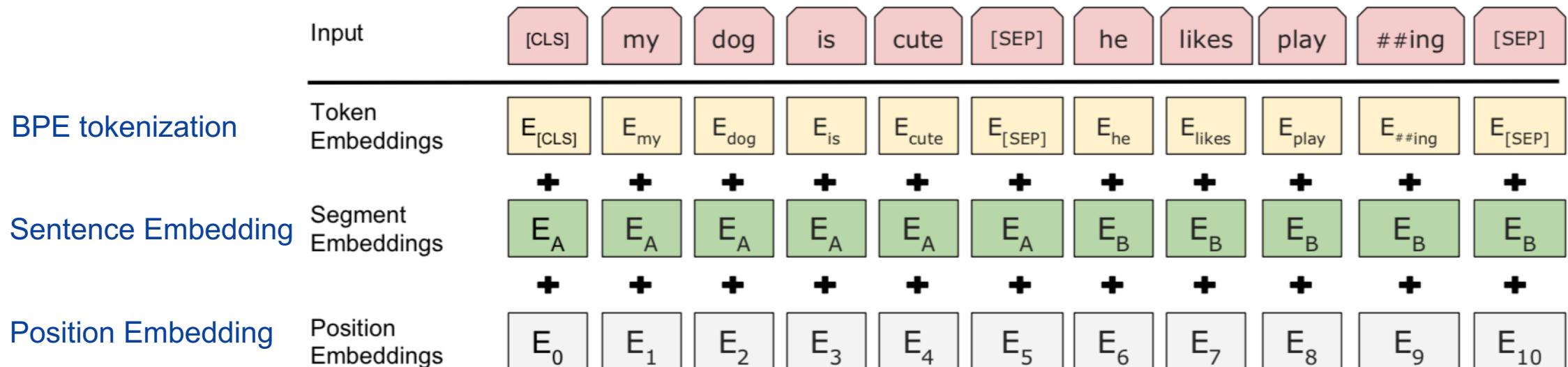
1. **Unidirectional** : “bank” is represented only based on “*I made*” but not “*deposit*”
2. **Bidirectional**: “bank” is represented by using the left and right context

Moreover, BERT is available with **multilingual** where sentences from different language are decoded in the same latent space



BERT input representation

- The input embedding is the sum of three parts:



OpenAI GPT-2

Extension of GPT (1.5B parameters, 10x more) that:

1. Zero-Shot transfer: no need for transfer learning. All the downstream task are modeled as predicting conditional probabilities

2. Model Modifications

- Layer normalization was moved to the input of each sub-block
- Additional layer normalization after the final self-attention block
- The weights of residual layers are scaled by $1/\sqrt{N}$ where N is the number of residual layers
- Use large vocabulary size and context

- Language and Machine Learning
- Distributional Semantics
- Context-free word embeddings:
 - Count based
 - Context based
- Contextual embeddings
- **Final Considerations & Conclusions**

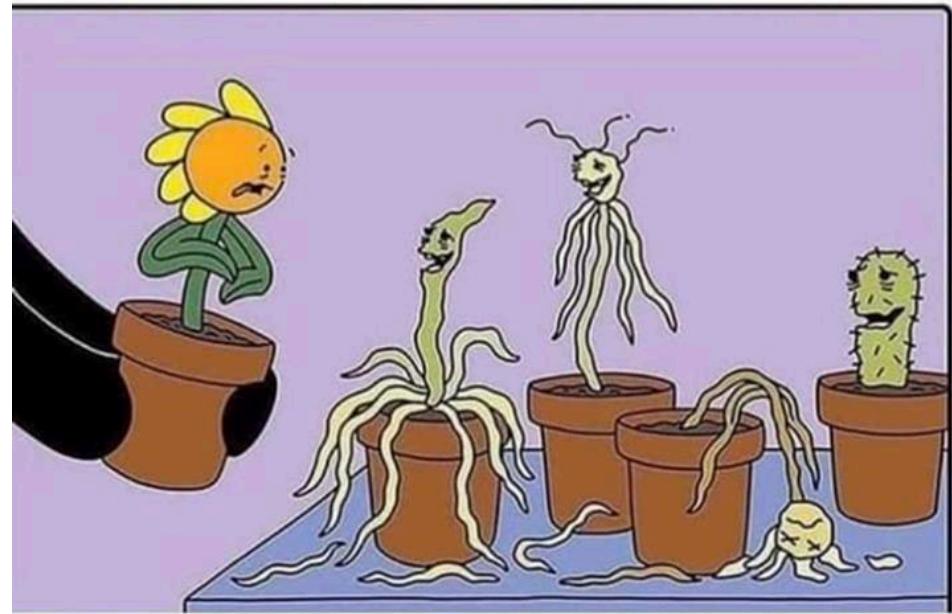
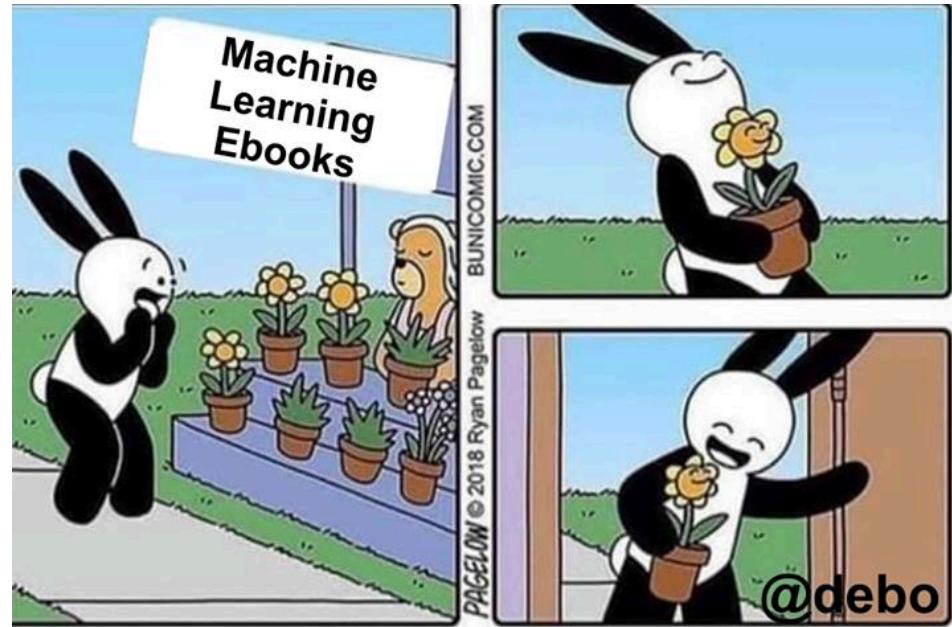
Final considerations & Conclusions

- We evaluated the evolution of context-free word embeddings to deal with:
 - Subword-level
 - OOV handling
- Then we considered contextual embeddings to address:
 - Polysemy and multi-sense word representation
 - Sentence embeddings
 - Embeddings for multiple languages
- What we do not covered are:
 - Temporal dimension
 - bias

Confidentiality

Any partial or total reproduction of its content is prohibited without written consent by Prometeia.

Copyright © 2019 Prometeia



Contacts

Bologna

Via Guglielmo Marconi, 43
+39 051 6480911
italy@prometeia.com

Milan

Via Brera, 18
Viale Monza, 265
+39 02 80505845
italy@prometeia.com

Rome

Viale Regina Margherita, 279
italy@prometeia.com

London

Dashwood House 69 Old Broad Street
EC2M 1QS
+44 (0) 207 786 3525
uk@prometeia.com

Istanbul

River Plaza, Kat 19
Büyükdere Caddesi Bahar Sokak
No. 13, 34394
| Levent | Istanbul | Turkey
+ 90 212 709 02 80 – 81 – 82
turkey@prometeia.com

 Prometeia

 @PrometeiaGroup

 Prometeiagroup

 Prometeia

Cairo

Smart Village - Concordia Building, B2111
Km 28 Cairo Alex Desert Road
6 of October City, Giza
egypt@prometeia.com

Moscow

ul. Ilyinka, 4
Capital Business Center Office 308
+7 (916) 215 0692
russia@prometeia.com

www.prometeia.com