

DisSent: Learning Sentence Representations from Explicit Discourse Relations

Allen Nie ----- Department of Computer Science
Erin D. Bennett ----- Department of Computer Science
Noah D. Goodman ----- Department of Psychology
Stanford University

Presented by: Pere-Lluís Huguet Cabot

Motivation

We have models like InferSent which worked well for sentence embeddings and are simple and straightforward.

But they need a lot of annotated data.

There are also models that don't need annotated data.

But they need a lot of data and can be complex and slow to train.

Using a self-supervised approach can take the best from both worlds. How?

Discourse markers

Keys

- An explicit discourse marker dataset generator.
- Use an unsupervised method to get sentence embeddings.
- Be good at it.
- Provide a new task and dataset.

What is an Explicit Discourse Relation?

[I wore a jacket] because [it was cold outside].
S1 marker S2

Because [it was cold outside], [I wore a jacket].
marker S2 S1

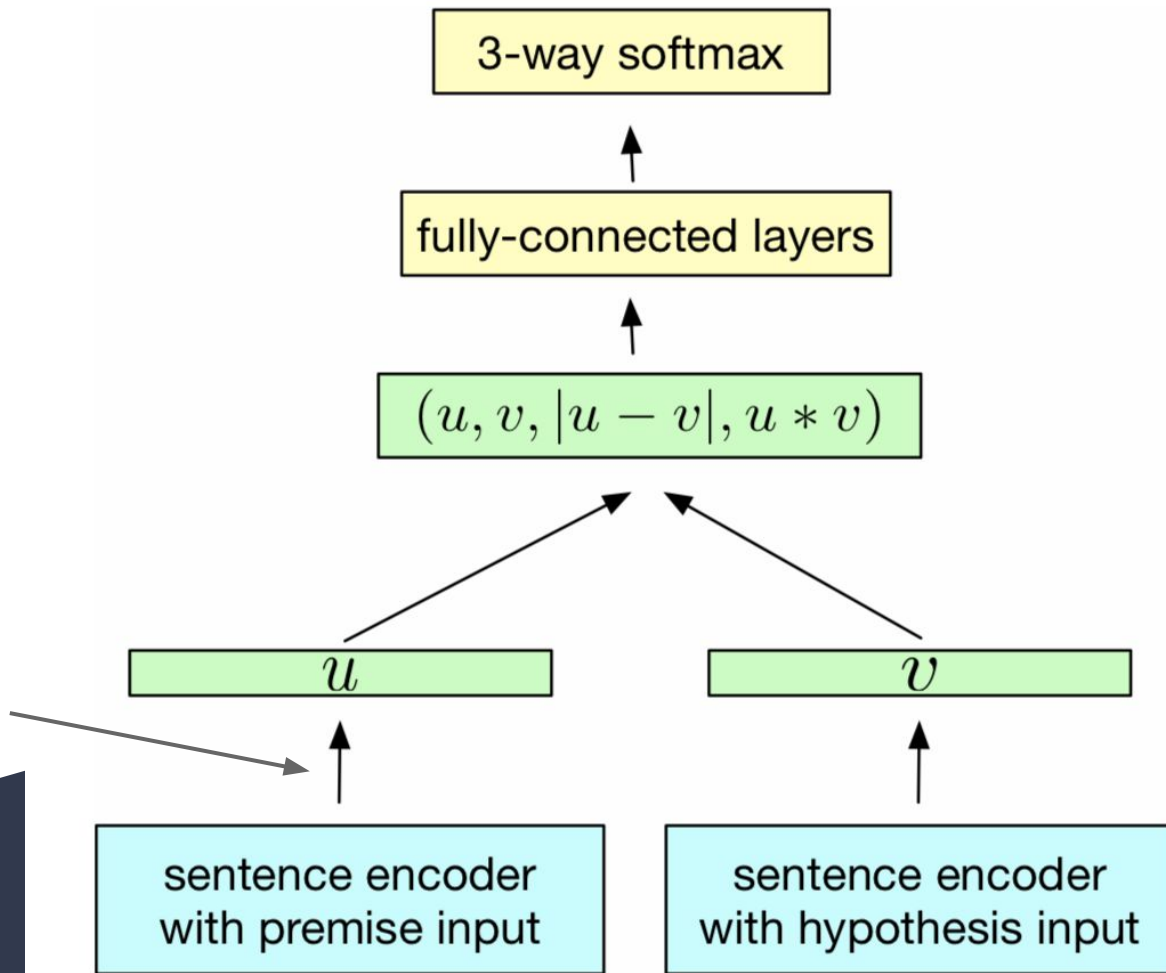
Bidirectional LSTM:

$$\vec{h}_t = \text{LSTM}_t(w_1, \dots, w_T | \theta_1)$$

$$\overleftarrow{h}_t = \text{LSTM}_t(w_T, \dots, w_1 | \theta_2)$$

$$h_t = [\vec{h}_t; \overleftarrow{h}_t]$$

$$s_i = \text{MaxPool}(h_1, \dots, h_T)$$

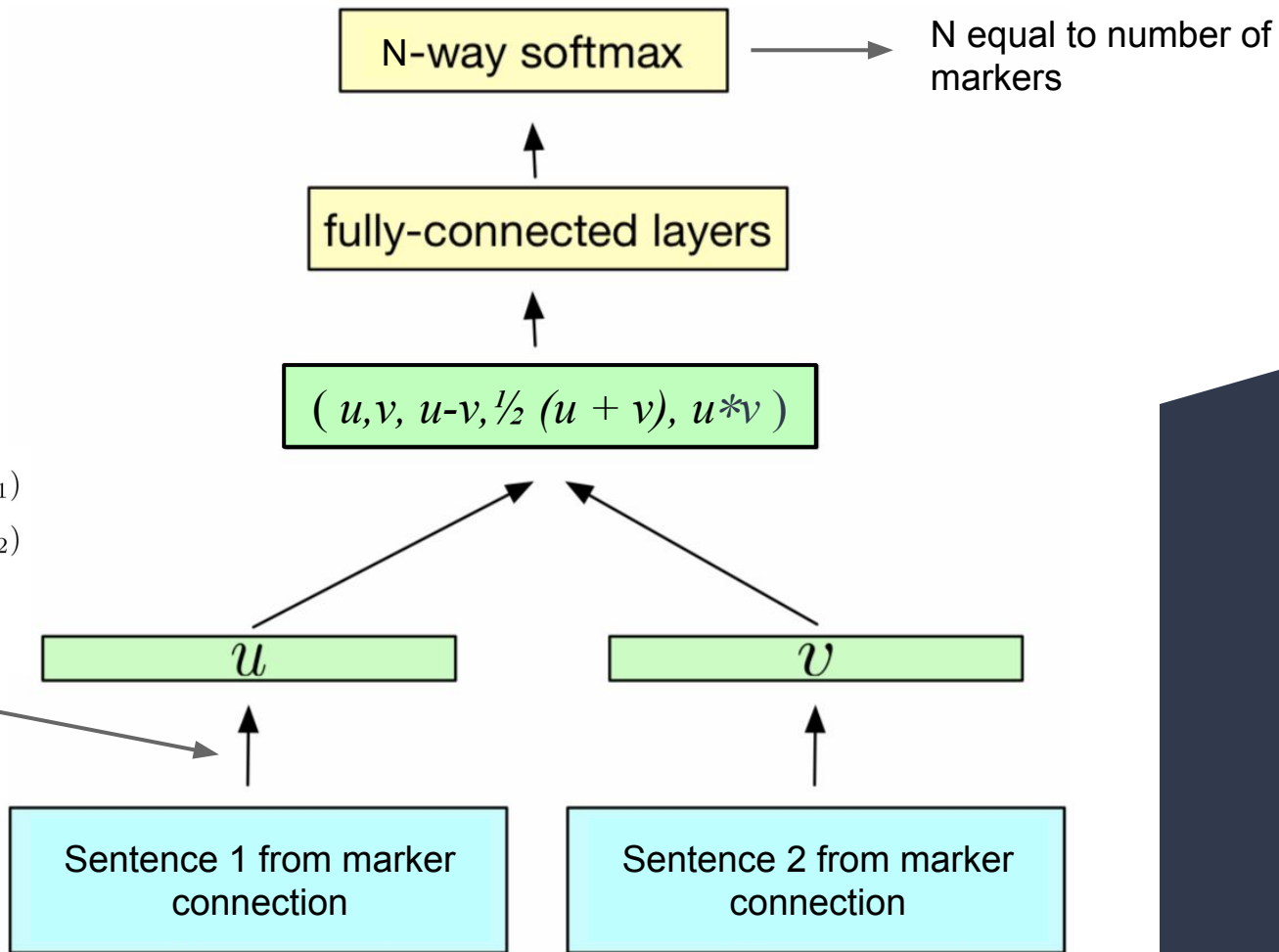


$$\vec{h}_t = \text{LSTM}_t(w_1, \dots, w_T | \theta_1)$$

$$\overleftarrow{h}_t = \text{LSTM}_t(w_T, \dots, w_1 | \theta_2)$$

$$h_t = [\vec{h}_t; \overleftarrow{h}_t]$$

$$s_i = \text{MaxPool}(h_1, \dots, h_T)$$



Dataset (BookCorpus)

Same Sentence (**SS**):

- [I wore a jacket] because [it was cold outside].

Immediate Predecessor (**IPS**):

- Because [it was cold outside], [I wore a jacket].

Non-adjacent previous sentence (**NAPS**):

- [I wore a jacket], which my sister gave to me, because [it was cold outside].

Sentence that follows (**FS**):

- [It was cold outside]. So [I wore a jacket]

1. Choice of Discourse Markers
 - a. At least 1% in the overall corpus
2. Stanford CoreNLP dependency parser
 - a. 91% of instances
3. Length-based Filtering
 - a. Less than 5
 - b. More than 50
 - c. Ratio of more than 5

Outcome:

4,706,292 pairs of sentences for 15
discourse markers.

Train: 90%

Validation: 5%

Test: 5%

With no supervision!

=)

Marker	Extracted Pairs	Percent (%)
and	818,634	21.1
as	761,330	19.6
when	552,540	14.2
but	508,648	13.1
if	491,394	12.6
before	268,787	6.9
while	120,231	3.1
because	116,444	3.0
after	84,330	2.2
though	61,023	1.6
so	57,816	1.5
although	13,933	0.4
still	11,125	0.3
also	10,026	0.3
then	8,414	0.2
Total	4,706,292	100.0

Table 1: Number of pairs of sentences extracted from BookCorpus for each discourse marker and percent of each marker in the resulting dataset.

Experiments

- Implicit vs. Explicit Prediction Task
- Implicit Relation Prediction Task
- *SentEval* (Conneau et al. (2017)):
 - sentiment analysis (MR, SST)
 - question-type (TREC)
 - product reviews (CR)
 - subjectivity-objectivity (SUBJ)
 - opinion polarity (MPQA)
 - entailment (SICK-E)
 - relatedness (SICK-R)
 - Paraphrase detection (MRPC)
- New Task: DIS, with a new dataset.

SDG, learning rate 0.1 with factor 5 annealing, 20 epochs, no dropout. Max pooling. 4096 hidden state size.

Marker	All	Books 8	Books 5
and	0.78 / 0.72	0.78 / 0.78	0.79 / 0.81
but	0.73 / 0.71	0.79 / 0.72	0.80 / 0.75
because	0.36 / 0.45	0.37 / 0.50	0.38 / 0.55
if	0.75 / 0.79	0.80 / 0.78	0.81 / 0.81
when	0.62 / 0.61	0.74 / 0.71	0.77 / 0.77
so	0.48 / 0.49	0.46 / 0.56	—
though	0.30 / 0.48	0.39 / 0.61	—
before	0.61 / 0.65	0.64 / 0.77	—
as	0.77 / 0.68	—	—
while	0.36 / 0.46	—	—
after	0.42 / 0.55	—	—
although	0.07 / 0.24	—	—
still	0.21 / 0.42	—	—
also	0.14 / 0.36	—	—
then	0.12 / 0.31	—	—
Overall	67.5	73.5	77.3

Table 4: **Training task performance:** Test recall / precision for each discourse marker on the classification task, and we report overall accuracy.

Model	IMP	IVE
DisSent Books 5 [†]	40.7	86.5
DisSent Books 8 [†]	41.4	87.9
DisSent Books ALL [†]	42.9	87.6
InferSent (Conneau et al., 2017)	38.4	84.5
Patterson and Kehler (2013)	—	86.6
Word Vectors (Qin et al., 2017)	36.9	74.8
Lin et al. (2009) + Brown Cluster	40.7	—
Adversarial Net (Qin et al., 2017)	46.2	—

Table 5: **Discourse Generalization Tasks using PDTB:** Following the metric used in these literature, we report overall test accuracy for sentence embedding models, as well as baselines and state of the art for these task.

Model	MR	CR	SUBJ	MPQA	SST	TREC	SICK-R	SICK-E	MRPC	DIS
Self-supervised training methods										
DisSent Books 5 [†]	80.2	85.4	93.2	90.2	82.8	91.2	0.845	83.5	76.1	75.7
DisSent Books 8 [†]	79.8	85.0	93.4	90.5	83.9	93.0	0.854	83.8	76.1	80.2
DisSent Books ALL [†]	80.1	84.9	93.6	90.1	84.1	93.6	0.849	83.7	75.0	79.9
Disc BiGRU	—	—	88.6	—	—	81.0	—	—	71.6	—
Unsupervised training methods										
FastSent	70.8	78.4	88.7	80.6	—	76.8	—	—	72.2	—
FastSent + AE	71.8	76.7	88.8	81.5	—	80.4	—	—	71.2	—
Skipthought	76.5	80.1	93.6	87.1	82.0	92.2	0.858	82.3	73.0	70.1
Skipthought-LN	79.4	83.1	93.7	89.3	82.9	88.4	0.858	79.5	—	—
Supervised training methods										
DictRep (bow)	76.7	78.7	90.7	87.2	—	81.0	—	—	—	—
InferSent	81.1	86.3	92.4	90.2	84.6	88.2	0.884	86.1	76.2	65.4
Multi-task training methods										
LSMTL	82.5	87.7	94.0	90.9	83.2	93.0	0.888	87.8	78.6	—

Conclusion and Personal thoughts

Self-supervised methods provide a huge advantage in terms of getting a dataset.

Provides a great tool to use in other tasks.

Use in other languages.

Limitations of evaluation.

Provides a new task and dataset to evaluate.

This is how we got models for word embeddings, which means this can be a good path for sentence embeddings.


Future

We could check for more insight on the embeddings by checking which neurons activates each marker.

Use other relations for the same purpose. Sentence order has been used. Maybe use other relations like Punctuation marks. Although it may not comprise semantic meaning it may show other interesting structure relations.

Thank you very much.

Feel free to ask any question
(but don't make it too hard)



ELMo

Deep Contextualized Word Representations

Victor Zuanazzi - MSc. AI
Universiteit van Amsterdam
15/04/2019

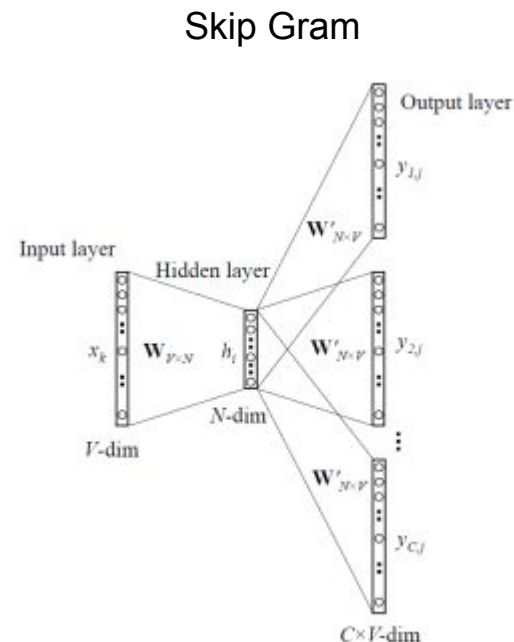
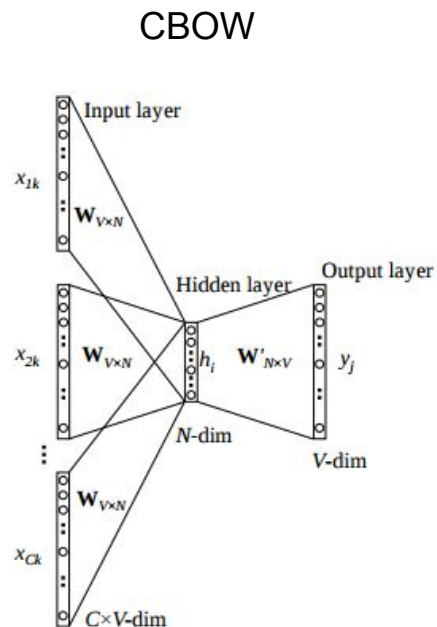


Agenda

- Word Embeddings
- (Deep) Contextualized Word Representations
- Embeddings from Language Models
- ELMo
- Experiments
- Results
- Personal Analysis
- Future Research
- Resources
- Q & A

Word embeddings

- Fixed sized vectors that represent words;
- They encode some semantic and syntactic aspects of the words;
- Built using word windows or dependencies.

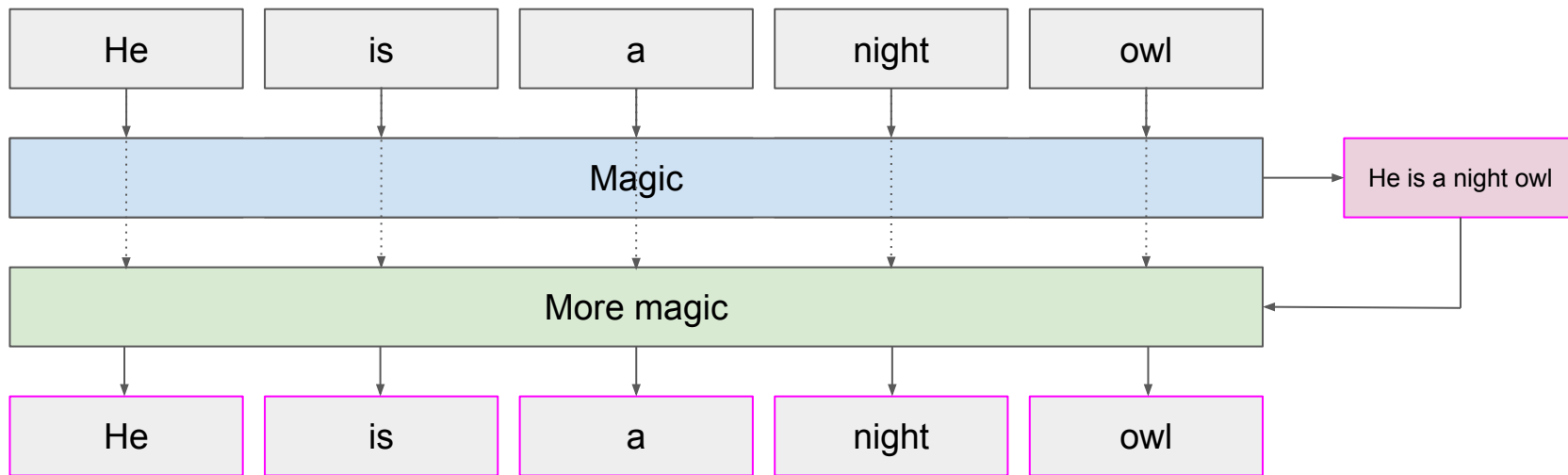


Word embeddings

- The embeddings have problems capturing polysemy;
- It is not straightforward how to use embeddings to capture sentence level semantics;
- In practice, different embeddings serve different applications.

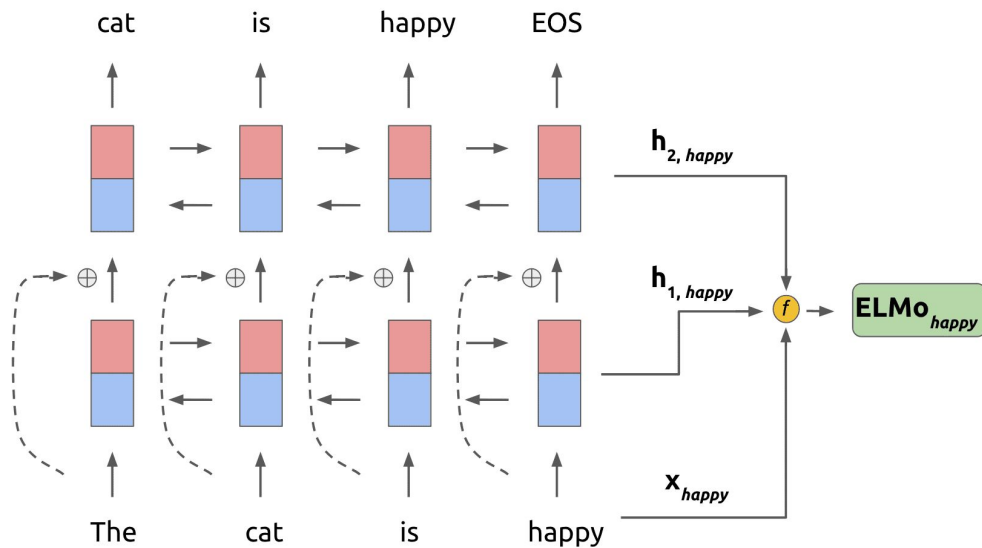
Contextualized Word Representations

- Fixed sized vectors that represent words in context;
- The same word has different representations in different contexts;
- The contextualized word representation is a function of the word embedding and the latent representation of the sentence.



Deep Contextualized Word Representations

- There is no reason for not using a deep neural net architecture and take many layers to create the (deep) contextualized word representation.



Questions

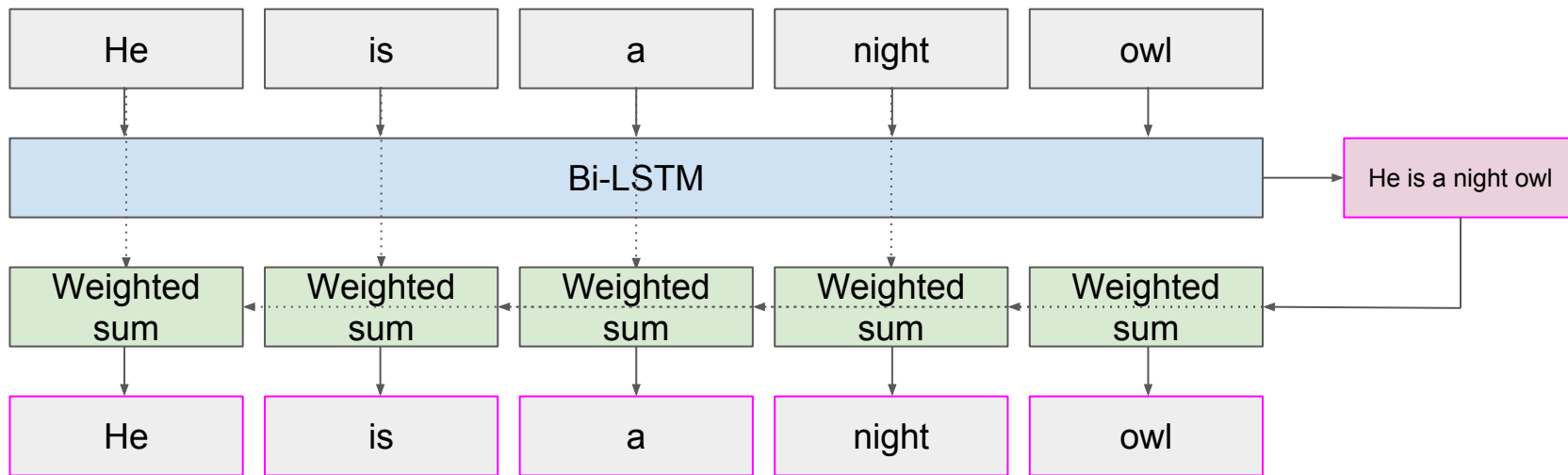
- Word Embeddings
- (Deep) Contextualized Word Representations
- Embeddings from Language Models
- ELMo
- Experiments
- Results
- Personal Analysis
- Future Research
- Resources
- Q & A

Embedding from Language Models - ELMo



Embedding from Language Models - ELMo

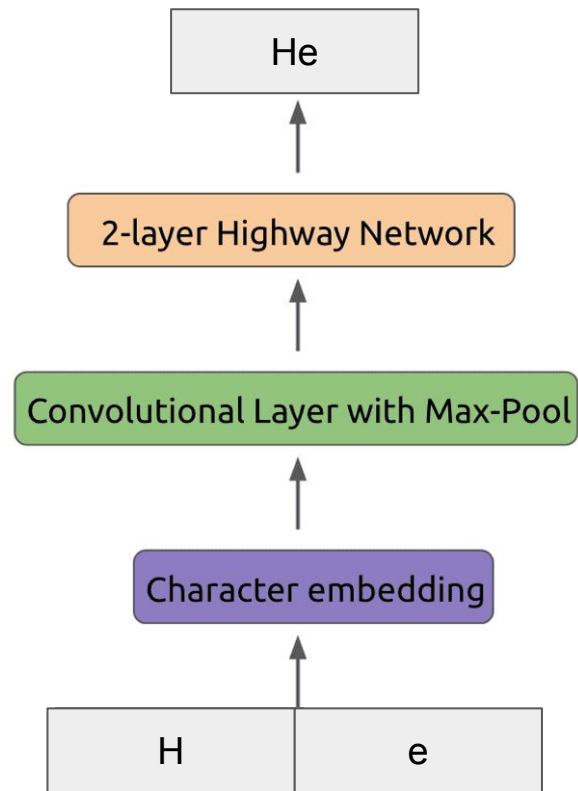
- ELMo learns word representations that are dependent on context
- Learned language model: BI-LSTM





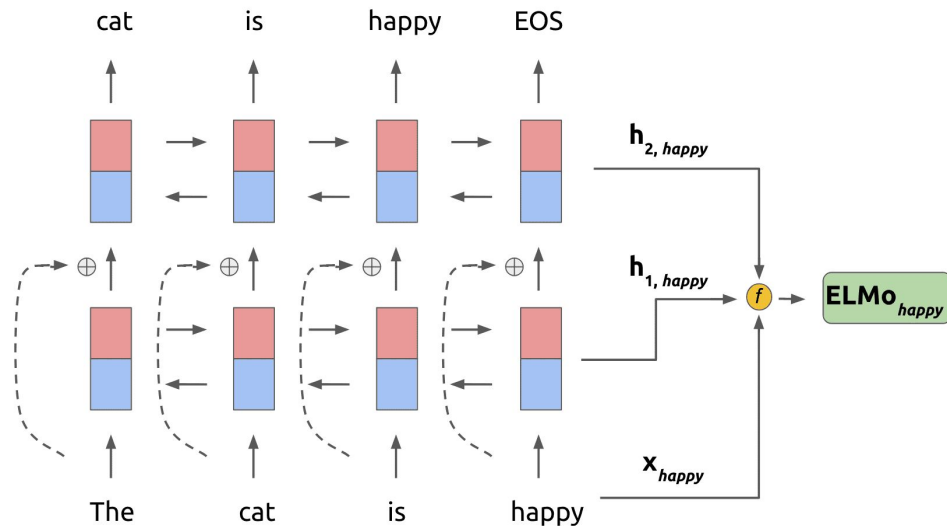
ELMo - Input

- Takes character embedding
- Convolutional Layer
- Max Pool Layer
- Highway Network



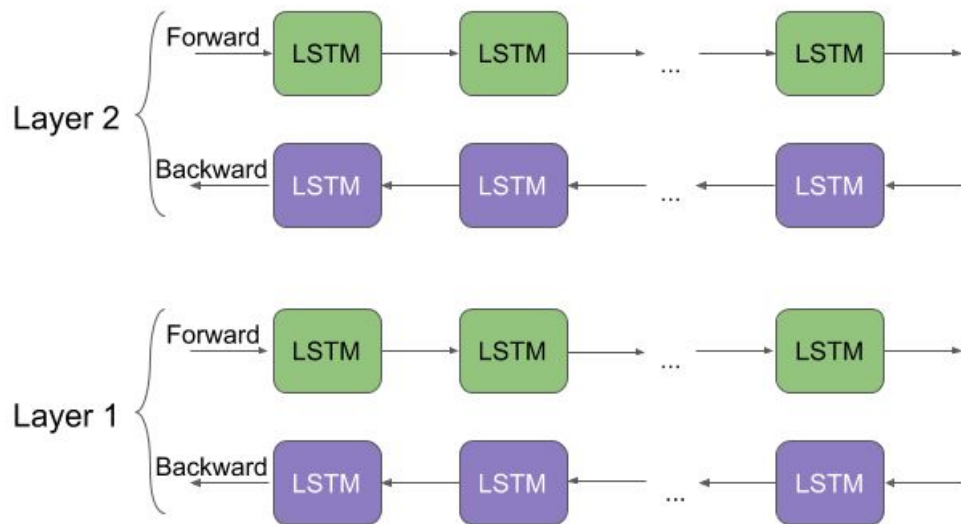
ELMo - Language Model

- The word embeddings are fed to a bi-Language Model modeled by a 2 layer bi-LSTM
- The hidden state of all layers is used together with the word embedding to create the contextualized word vector
- Trained using 30 million sentences.

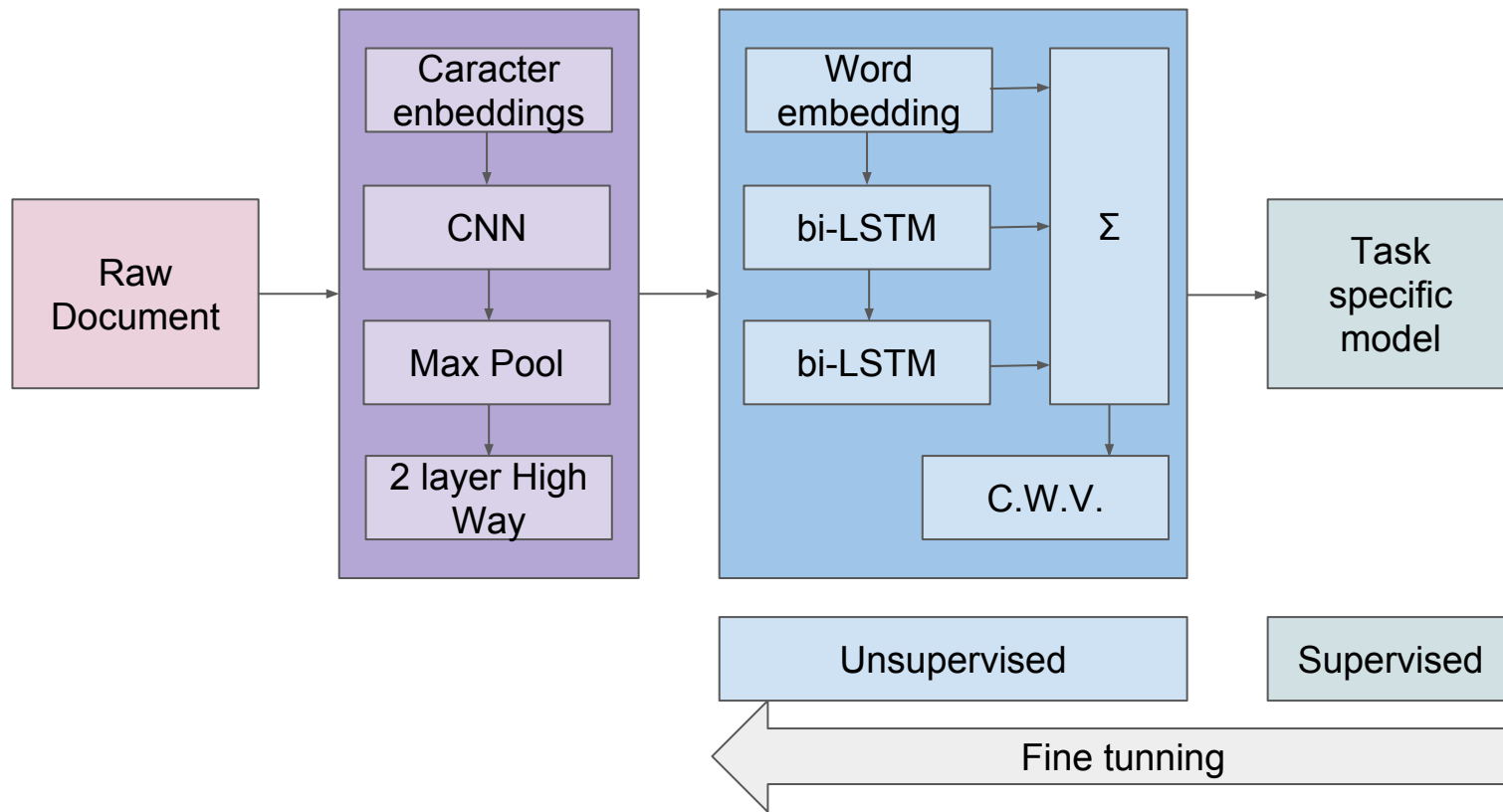


$$p(t_1, t_2, \dots, t_N) = \sum_{k=1}^N (\log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) \\ + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s)).$$

ELMo - Language Model



ELMo - supervised tasks



Questions

- Word Embeddings
- (Deep) Contextualized Word Representations
- Embeddings from Language Models
- ELMo
- Experiments
- Results
- Personal Analysis
- Future Research
- Resources
- Q & A

Experiments

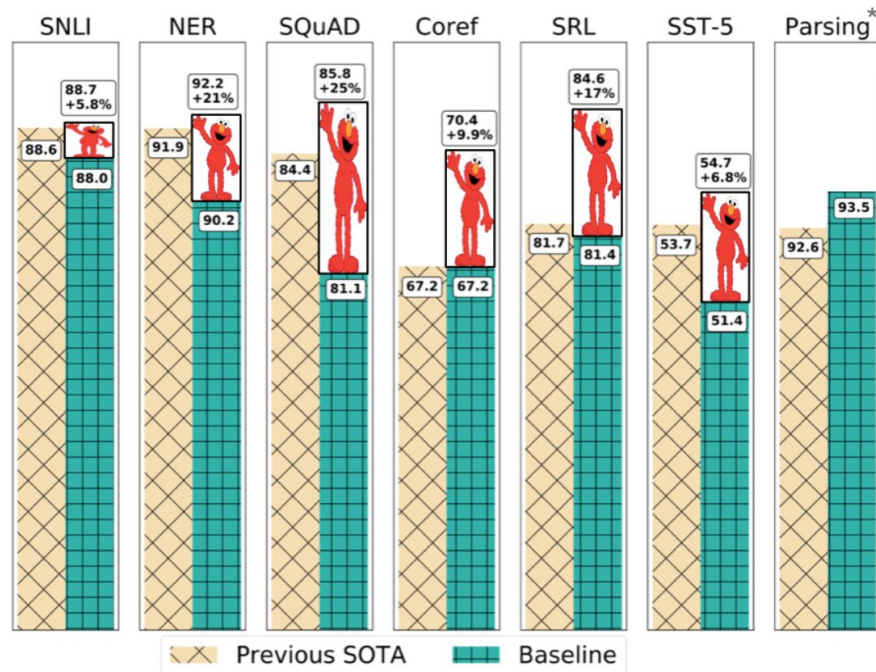
- Question Answering (SQuAD)
- Textual Entailment (SNLI)
- Semantic Role Labeling (SRL)
- Coreference Resolution (Coref)
- Named Entity Recognition (NER)
- Sentiment Analysis(SST-5)

Experiments

- Question Answering (SQuAD)
- **Textual Entailment (SNLI)**
- **Semantic Role Labeling (SRL)**
- Coreference Resolution (Coref)
- Named Entity Recognition (NER)
- Sentiment Analysis(SST-5)

Results

TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	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%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

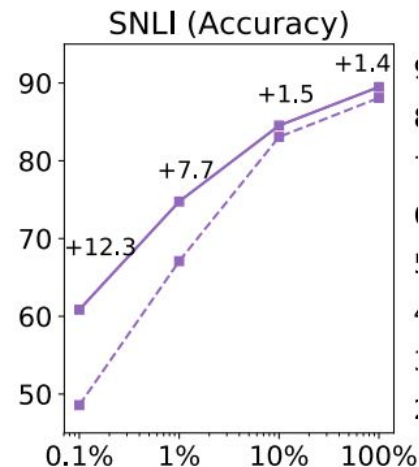




Textual Entailment - SNLI

TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%

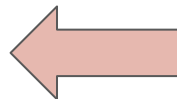
- Base line: ESIM;
- Small improvement;
- No significance test was reported;
- Performance is consistently better for different training set sizes.



Textual Entailment - SNLI

- Not current state of the art

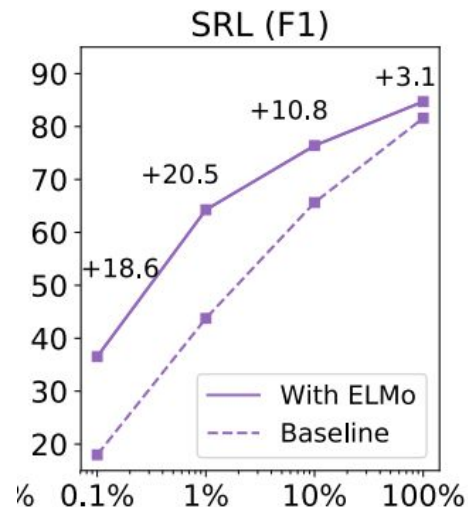
Publication	Model	Parameters	Train (% acc)	Test (% acc)
Qian Chen et al. '16	600D ESIM + 300D Syntactic TreeLSTM (code)	7.7m	93.5	88.6
Peters et al. '18	ESIM + ELMo	8.0m	91.6	88.7
Boyuan Pan et al. '18	300D DMAN	9.2m	95.4	88.8
Zhiguo Wang et al. '17	BiMPM Ensemble	6.4m	93.2	88.8
Yichen Gong et al. '17	448D Densely Interactive Inference Network (DIIN, code) Ensemble	17m	92.3	88.9
Seonhoon Kim et al. '18	Densely-Connected Recurrent and Co-Attentive Network	6.7m	93.1	88.9
Zhuosheng Zhang et al. '18	SLRC	6.1m	89.1	89.1
Qian Chen et al. '17	KIM Ensemble	43m	93.6	89.1
Ghaeini et al. '18	450D DR-BiLSTM Ensemble	45m	94.8	89.3
Peters et al. '18	ESIM + ELMo Ensemble	40m	92.1	89.3
Yi Tay et al. '18	300D CAFE Ensemble	17.5m	92.5	89.3
Chuanqi Tan et al. '18	150D Multiway Attention Network Ensemble	58m	95.5	89.4
Boyuan Pan et al. '18	300D DMAN Ensemble	79m	96.1	89.6
Radford et al. '18	Fine-Tuned LM-Pretrained Transformer	85m	96.6	89.9
Seonhoon Kim et al. '18	Densely-Connected Recurrent and Co-Attentive Network Ensemble	53.3m	95.0	90.1
Xiaodong Liu et al. '19	MT-DNN	110m	96.8	91.1



Semantic Role Labeling - SRL

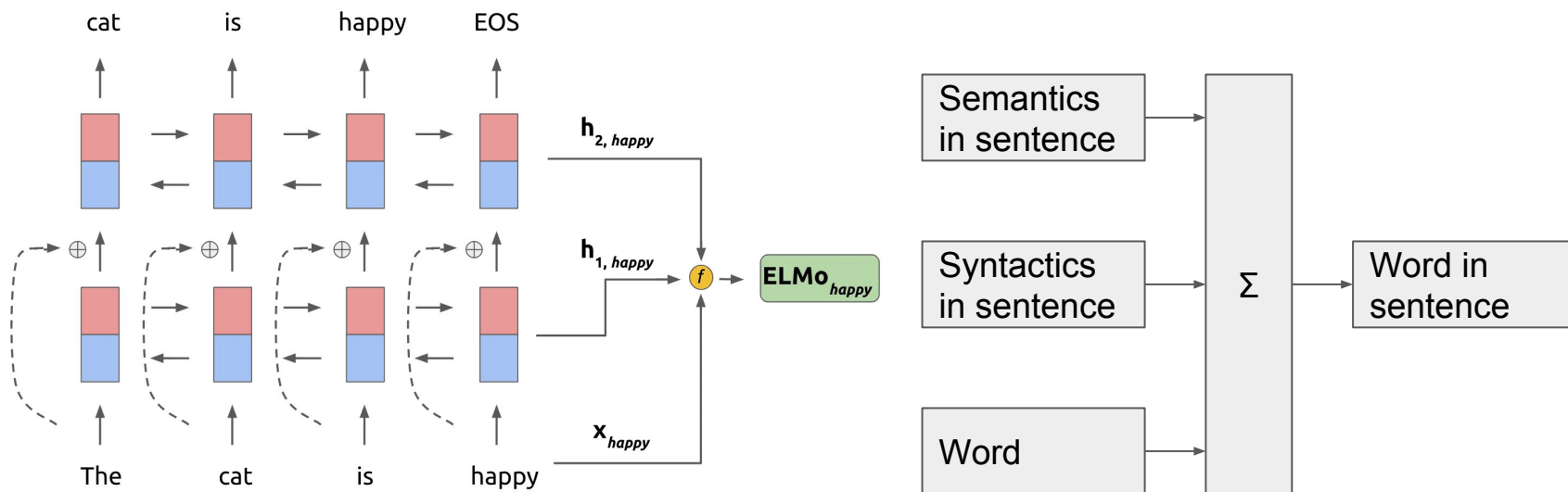
TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%

- Strong improvement;
- No significance test was reported;
- Performance is consistently better for different training set sizes.



Semantic Role Labeling - SRL

- First LSTM layer captures syntactic roles;
- Second LSTM layer captures semantic roles;



Semantic Role Labeling - SRL

Model	F1	Paper / Source
He et al., (2018) + ELMO	85.5	Jointly Predicting Predicates and Arguments in Neural Semantic Role Labeling
(He et al., 2017) + ELMo (Peters et al., 2018)	84.6	Deep contextualized word representations
Tan et al. (2018)	82.7	Deep Semantic Role Labeling with Self-Attention
He et al. (2018)	82.1	Jointly Predicting Predicates and Arguments in Neural Semantic Role Labeling
He et al. (2017)	81.7	Deep Semantic Role Labeling: What Works and What's Next

Questions

- Word Embeddings
- (Deep) Contextualized Word Representations
- Embeddings from Language Models
- ELMo
- Experiments
- Results
- Future Research
- Personal Analysis
- Resources
- Q & A

Future Research

	Base model	pre-training	Downstream tasks	Downstream model	Fine-tuning
CoVe	seq2seq NMT model	supervised	feature-based	task-specific	/
ELMo	two-layer biLSTM	unsupervised	feature-based	task-specific	/
CVT	two-layer biLSTM	semi-supervised	model-based	task-specific / task-agnostic	/
ULMFiT	AWD-LSTM	unsupervised	model-based	task-agnostic	all layers; with various training tricks
GPT	Transformer decoder	unsupervised	model-based	task-agnostic	pre-trained layers + top task layer(s)
BERT	Transformer encoder	unsupervised	model-based	task-agnostic	pre-trained layers + top task layer(s)
GPT-2	Transformer decoder	unsupervised	model-based	task-agnostic	pre-trained layers + top task layer(s)

Personal Analysis

- Quite a complex architecture;
- Why bi-LSTMs? Why two layers?
- No statistical significance study made on the results;
- No need for annotated data
- “Plug and play” embeddings = Transfer Learning
- Interesting findings on syntactic and semantic modeling.

Resources

1. <https://arxiv.org/abs/1802.05365>
2. <https://www.mihaileric.com/posts/deep-contextualized-word-representations-elmo/>
3. <https://www.analyticsvidhya.com/blog/2019/03/learn-to-use-elmo-to-extract-features-from-text/>
4. <https://towardsdatascience.com/review-highway-networks-gating-function-to-highway-image-classification-5a33833797b5>
5. <https://towardsdatascience.com/besides-word-embedding-why-you-need-to-know-character-embedding-6096a34a3b10>
6. <https://arxiv.org/pdf/1609.06038.pdf>
7. <https://lilianweng.github.io/lil-log/2019/01/31/generalized-language-models.html>
8. <http://jalammar.github.io/illustrated-bert/>
9. <https://nlpprogress.com/>
10. <http://runder.io/nlp-imagenet/>

Q & A

#FALLONTONIGHT

