

# Deep contextualized word representations

Khoklin I., Kovrigin A., Shibaev E.

# Introduction: limitations of previous models

Pre-trained word representations are vectorized forms of words that capture their meanings, relationships, and contexts. Developed models like Word2Vec and GloVe have significantly enhanced our capability to understand and process human languages in various NLP applications.

However, these models have their limitations:

- Words are assigned a single vector regardless of their usage, leading to inadequate handling of polysemous words (words with multiple meanings).
- The representations are static, meaning the context of surrounding words isn't considered, which can be crucial for understanding the true meaning of a word in a particular sentence.

# Introduction: Need for Enhanced Representations

The challenges mentioned before really point out that we need better and more aware word representations that understand the context. We need a model that can get the subtle details of how words are used, grasp the different meanings a word can have depending on the situation, and adjust to the complicated ways we use language.

# Method: Bidirectional language model

Embeddings from Language Models (ELMo) model uses bidirectional language model as its base.

A forward language model computes the probability of the sequence by modeling the probability of token given tokens before it:

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1})$$

A backward LM is similar to a forward LM, except it runs over the sequence in reverse, predicting the previous token given the future context:

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N)$$

## Method: Bidirectional language model

A biLM combines forward and backward LM. Overall loss is sum of log likelihood of the forward and backward directions:

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

The final model uses  $L = 2$  biLSTM layers. In order to obtain embeddings for LSTMs character convolutions are used.

## Method: General architecture

For each token  $t_k$  ELMo computes all intermediate representations of both LM:

$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

For downstream model, ELMo collapses all layers into a single vector via weighting of all biLM layers:

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}.$$

# Method: Using ELMo for downstream tasks

Process of using ELMo to improve the task model:

- Freeze the weight of biLM.
- Concatenate vector **ELMo**<sub>*k*</sub> with embedding  $x_k$  of supervised architecture and pass this enhanced representation into target model.
- It was noticed that adding another linear combination to the output of the target model can also be beneficial.

# Evaluation

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	$88.7 \pm 0.17$	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)	$91.93 \pm 0.19$	90.15	$92.22 \pm 0.10$	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	$54.7 \pm 0.5$	3.3 / 6.8%

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model baselines across six benchmark NLP tasks. The performance metric varies across tasks – accuracy for SNLI and SST-5;  $F_1$  for SQuAD, SRL and NER; average  $F_1$  for Coref. Due to the small test sizes for NER and SST-5, we report the mean and standard deviation across five runs with different random seeds. The “increase” column lists both the absolute and relative improvements over our baseline.



# Analysis: regularization and utilisation

Task	Baseline	Last Only	All layers	
			$\lambda=1$	$\lambda=0.001$
SQuAD	80.8	84.7	85.0	<b>85.2</b>
SNLI	88.1	89.1	89.3	<b>89.5</b>
SRL	81.6	84.1	84.6	<b>84.8</b>

Table 2: Development set performance for SQuAD, SNLI and SRL comparing using all layers of the biLM (with different choices of regularization strength  $\lambda$ ) to just the top layer.

Task	Input Only	Input & Output	Output Only
SQuAD	85.1	<b>85.6</b>	84.8
SNLI	88.9	<b>89.5</b>	88.7
SRL	<b>84.7</b>	84.3	80.9

Table 3: Development set performance for SQuAD, SNLI and SRL when including ELMo at different locations in the supervised model.

# Analysis: layer contribution

Model	F <sub>1</sub>
WordNet 1st Sense Baseline	65.9
<a href="#">Raganato et al. (2017a)</a>	69.9
<a href="#">Iacobacci et al. (2016)</a>	<b>70.1</b>
CoVe, First Layer	59.4
CoVe, Second Layer	64.7
biLM, First layer	67.4
biLM, Second layer	69.0

Table 5: All-words fine grained WSD F<sub>1</sub>. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.

Model	Acc.
<a href="#">Collobert et al. (2011)</a>	97.3
<a href="#">Ma and Hovy (2016)</a>	97.6
<a href="#">Ling et al. (2015)</a>	<b>97.8</b>
CoVe, First Layer	93.3
CoVe, Second Layer	92.8
biLM, First Layer	97.3
biLM, Second Layer	96.8

Table 6: Test set POS tagging accuracies for PTB. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.

# Analysis: interesting to look at

Source		Nearest Neighbors
GloVe	play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM	Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent <u>play</u> .
	Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}	{...} they were actors who had been handed fat roles in a successful <u>play</u> , and had talent enough to fill the roles competently , with nice understatement .

Table 4: Nearest neighbors to “play” using GloVe and the context embeddings from a biLM.

