# Deep contextualized word representation
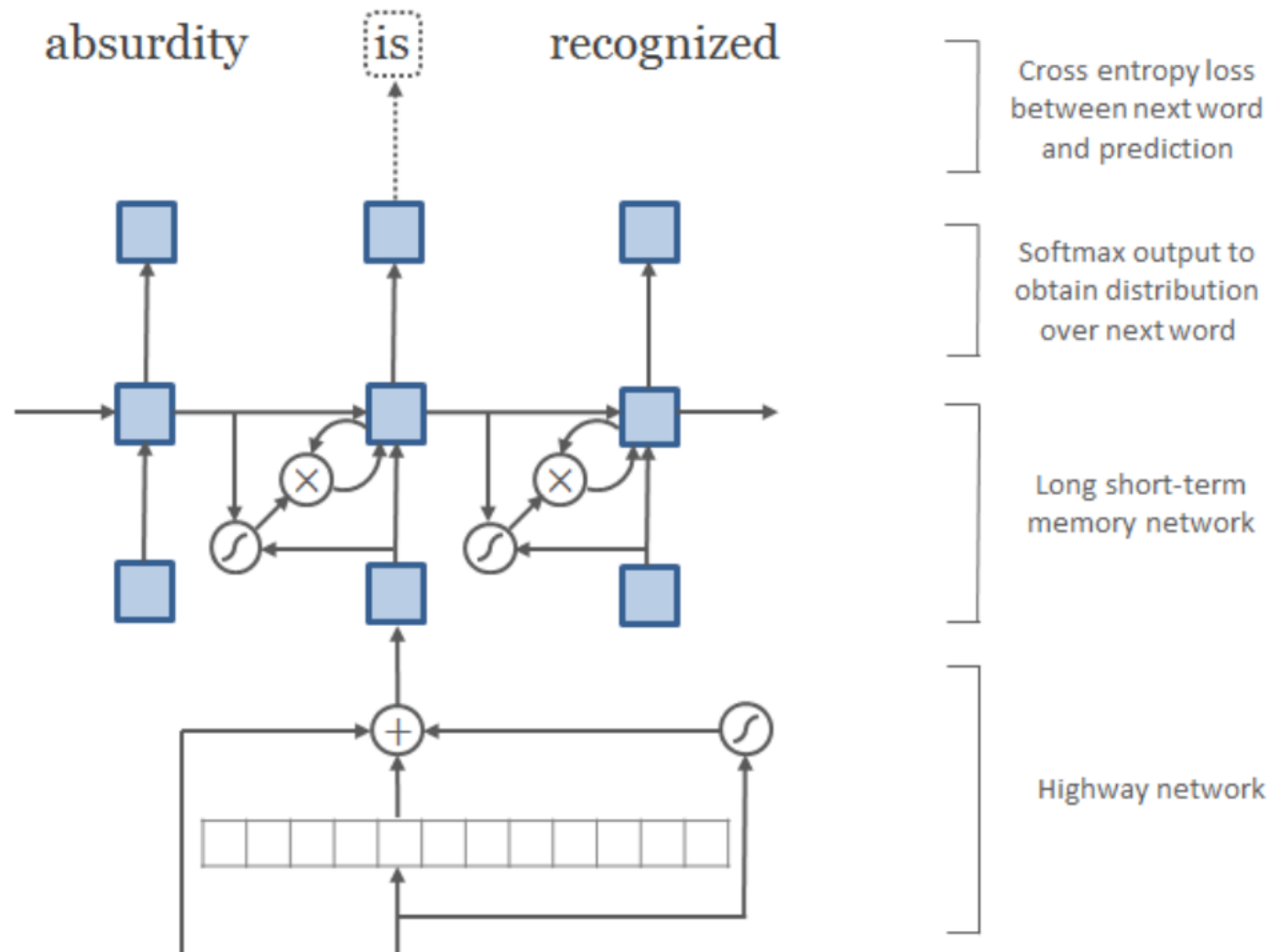
AntNLP

Yupei Du

# Language Models



Max-over-time pooling layer

$\max\{\cdot\}$

Convolution layer with multiple filters of different widths

Concatenation of character embeddings

moment     the     absurdity     is     recognized

# Language Models

absurdity    is    recognized



Cross entropy loss between next word and prediction

Softmax output to obtain distribution over next word

Long short-term memory network

Highway network

# ELMo

## Deep-BiLSTM

$$\sum_{k=1}^{N} (\log p(t_k \mid t_1, \ldots, t_{k-1}; \Theta_x, \overrightarrow{\Theta}_{LSTM}, \Theta_s)$$

$$+ \log p(t_k \mid t_{k+1}, \ldots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s))$$

## ELMo

$$R_k = \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \ldots, L\}$$

$$= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \ldots, L\},$$

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

# Experiments

*SQuAD(Question answering)* : 100K+ crowd sourced question-answer pairs where the answer is a span in a given Wikipedia paragraph

*SNLI(Textual entailment)* : 550K hypothesis/premise pairs.determining whether a "hypothesis" is true, given a "premise"

*SRL(Semantic role labeling)* : models the predicate–argument structure of a sentence, and is often described as answering "Who did what to whom"

*Coref(Coreference resolution)* : clustering mentions in text that re- fer to the same underlying real world entities

# Experiments

*NER(Named entity extraction)* : tag entities with four different entity types

*SST-5(Sentiment analysis)* : selecting one of five labels (from very negative to very positive) to describe a sentence from a movie review

| 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% |

# Analysis

Regularize

Adding $\lambda||w||_2^2$ to loss function

| Task | Baseline | Last Only | All layers | |
|------|----------|-----------|------------|--|
| | | | $\lambda=1$ | $\lambda=0.001$ |
| SQuAD | 80.8 | 84.7 | 85.0 | **85.2** |
| SNLI | 88.1 | 89.1 | 89.3 | **89.5** |
| SRL | 81.6 | 84.1 | 84.6 | **84.8** |

Input or Output?

| Task | Input Only | Input & Output | Output Only |
|------|------------|----------------|-------------|
| SQuAD | 85.1 | **85.6** | 84.8 |
| SNLI | 88.9 | **89.5** | 88.7 |
| SRL | **84.7** | 84.3 | 80.9 |

# Analysis

What information is captured?

| Model | F$_1$ |
|---|---|
| WordNet 1st Sense Baseline | 65.9 |
| Raganato et al. (2017a) | 69.9 |
| Iacobacci et al. (2016) | **70.1** |
| CoVe, First Layer | 59.4 |
| CoVe, Second Layer | 64.7 |
| biLM, First layer | 67.4 |
| biLM, Second layer | 69.0 |

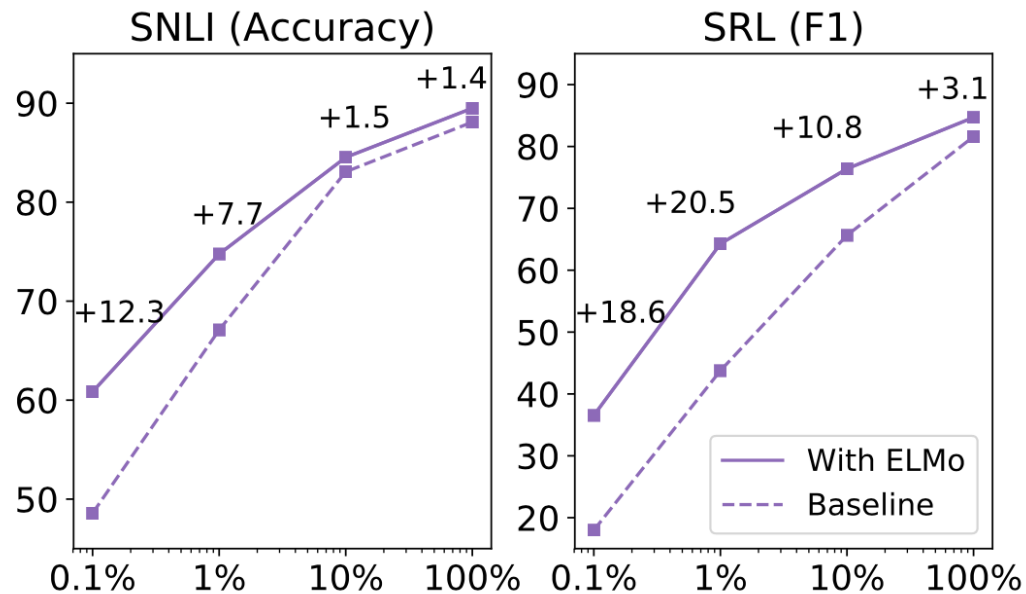| Model | Acc. |
|---|---|
| Collobert et al. (2011) | 97.3 |
| Ma and Hovy (2016) | 97.6 |
| Ling et al. (2015) | **97.8** |
| CoVe, First Layer | 93.3 |
| CoVe, Second Layer | 92.8 |
| biLM, First Layer | 97.3 |
| biLM, Second Layer | 96.8 |

# Analysis

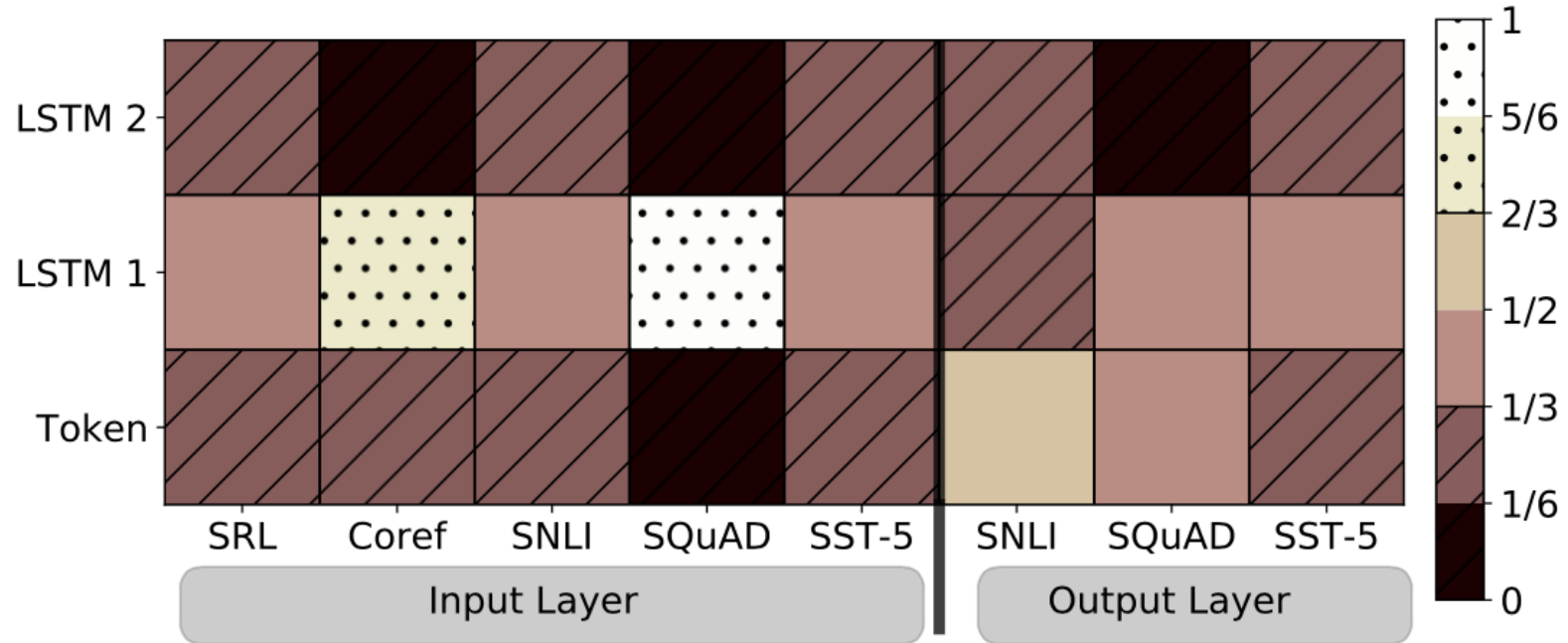Sample efficiency

*Speed*

$$486 \rightarrow 10 \; epochs \; in \; SRL$$

*dataset size*

# Analysis

Visualization of weights

# Thank you!

Q&A