

Comparative Study of Neural Language Models

Abhinash Khare

Soham Pal

Sruthi Gorantla

Today we will see...

- A Conditional Random Field Word Segmenter for Sighan Bakeoff 2005 (Tseng et al.)
- Exploring the Limits of Language Modeling (Jozefowicz et al.)
- Character-Based Neural Machine Translation (Ling et al.)
- Preliminary experimental results.

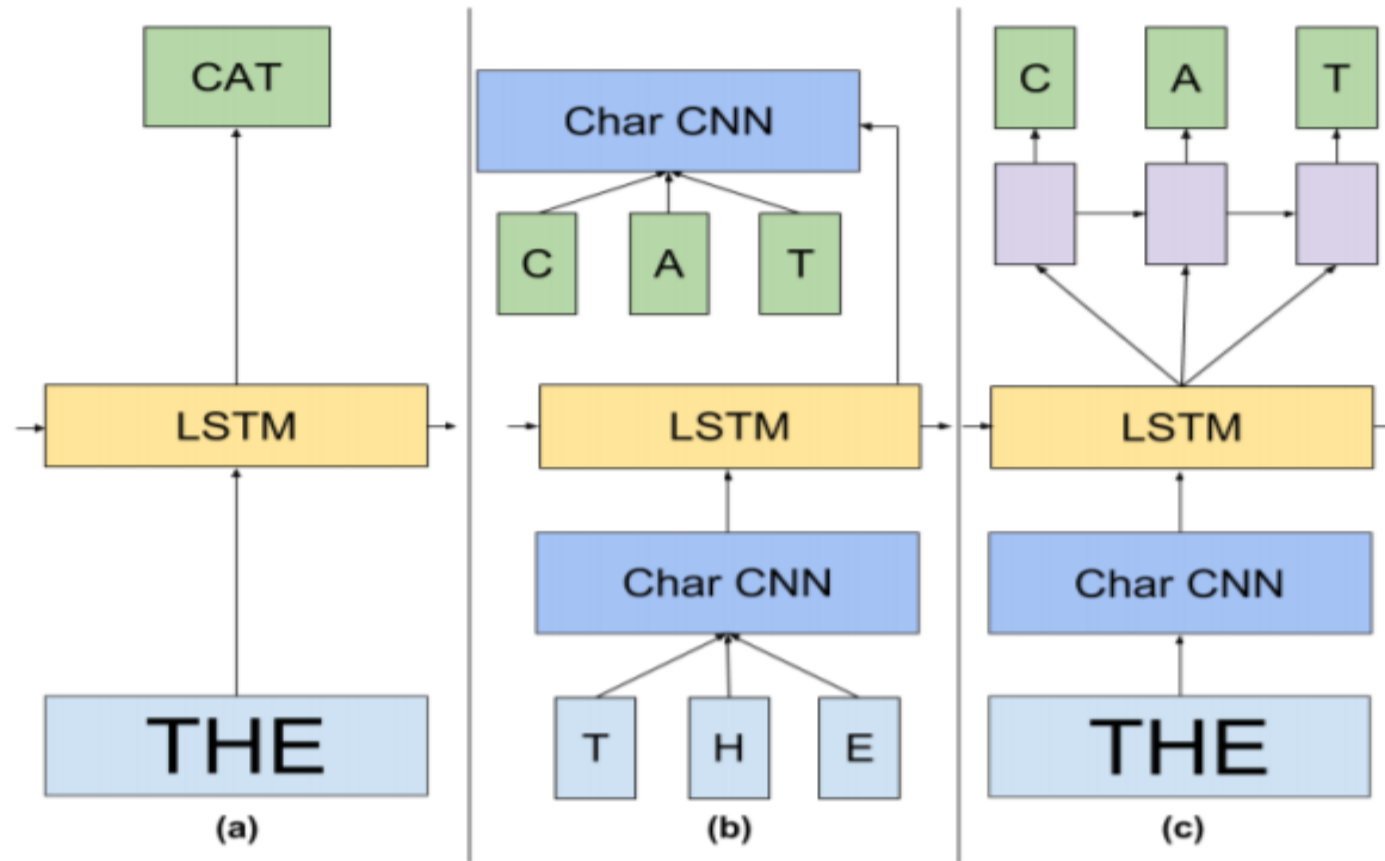
Exploring the Limits of Language Modeling

- Proposed by Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, Yonghui Wu from Google Brain in 2016.
- Performed an empirical study on Language Modeling using CNNs and RNN-LSTMs on One Billion Word Benchmark.
- Proposed models that improve the state-of-the-art perplexity significantly.

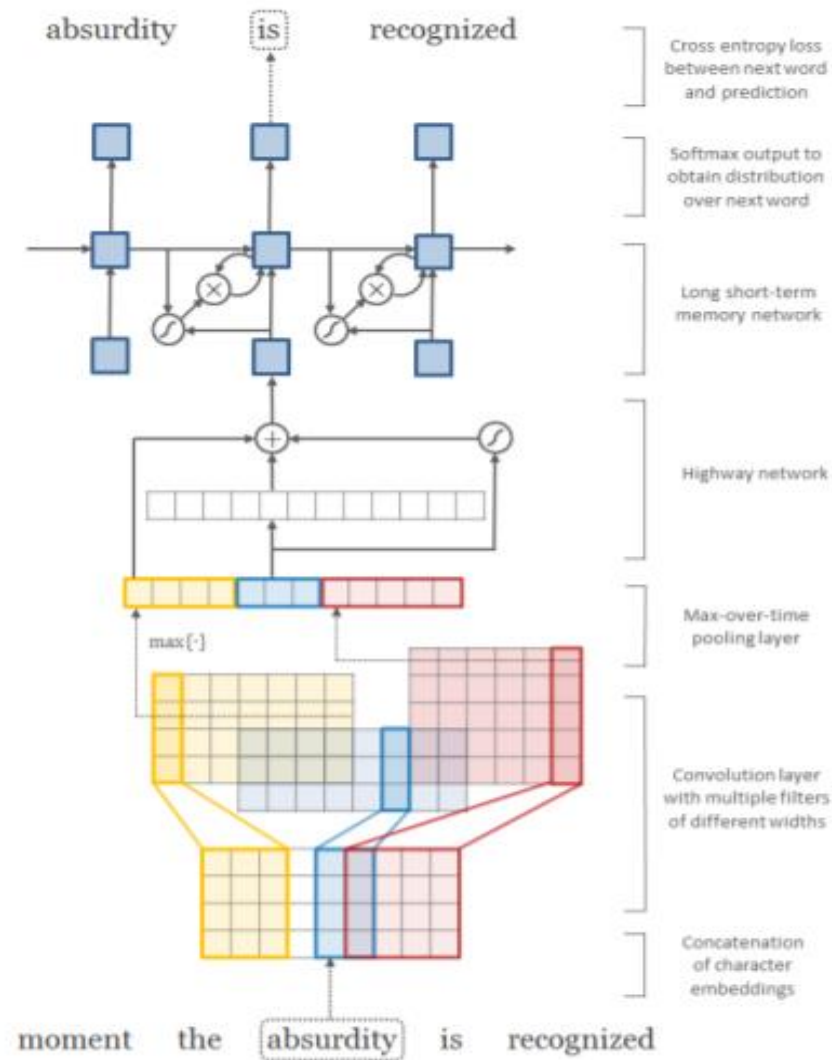
Introduction

- Believe that results on large datasets show much better about what actually matters.
- **Dataset:** One Billion Word Benchmark Dataset (800k word vocabulary).
- **Models:** They have considered three different models using LSTMs and Char CNNs.

Models (Jozefowicz et al.)



Char-CNN (Kim et al.)



Challenges in Language Modeling

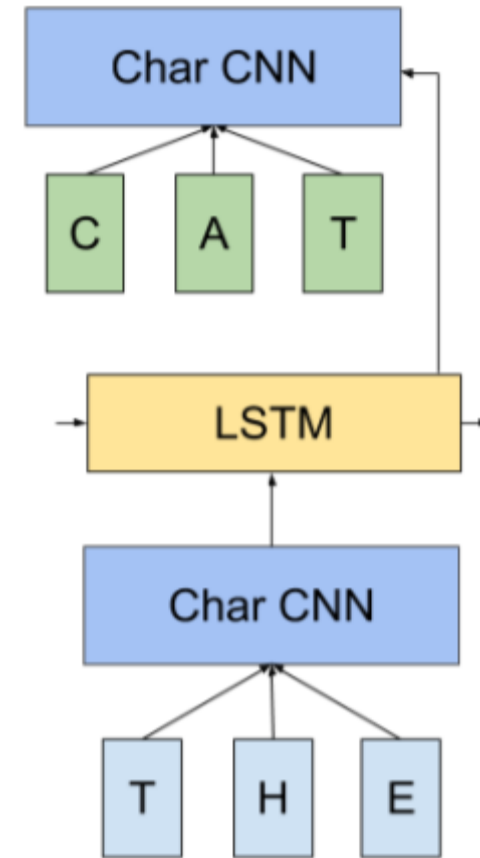
- Learning a very good model for human language which requires large models and thus, large amounts of data.
- Softmax over large vocabularies.

$$p(w) = \frac{\exp(z_w)}{\sum_{w' \in V} \exp(z_{w'})} \quad \text{where} \quad z_w = h^T e_w$$

Improvements in Language Modeling (1)

- **CNN Softmax:**

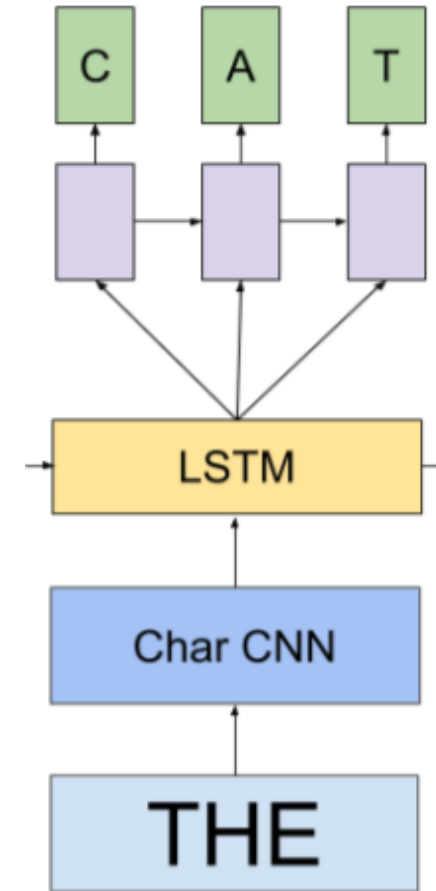
- Instead of computing a huge vector $|V| * |h|$, e_w is calculated by a CNN over the characters of w as $e_w = \text{CNN}(\text{chars}_w)$.
- Though the training time doesn't decrease much, the number of parameters is reduced as the CNN is shared for all the words and the computations can be parallelised.



Improvements in Language Modeling (2)

- **Char LSTM Predictions:**

- Used to predict one character at a time such that the probabilities are calculated over much smaller vocabulary.
- Combination of word-level and char-level LSTMs are used. However, they perform worse than the regular or CNN softmax, even though they scale independently of vocabulary size both during training as well as inference.



Results (Jozefowicz et al.)

MODEL	TEST PERPLEXITY	NUMBER OF PARAMS [BILLIONS]
SIGMOID-RNN-2048 (JI ET AL., 2015A)	68.3	4.1
INTERPOLATED KN 5-GRAM, 1.1B N-GRAMS (CHELBA ET AL., 2013)	67.6	1.76
SPARSE NON-NEGATIVE MATRIX LM (SHAZEER ET AL., 2015)	52.9	33
RNN-1024 + MAXENT 9-GRAM FEATURES (CHELBA ET AL., 2013)	51.3	20
LSTM-512-512	54.1	0.82
LSTM-1024-512	48.2	0.82
LSTM-2048-512	43.7	0.83
LSTM-8192-2048 (NO DROPOUT)	37.9	3.3
LSTM-8192-2048 (50% DROPOUT)	32.2	3.3
2-LAYER LSTM-8192-1024 (BIG LSTM)	30.6	1.8
BIG LSTM+CNN INPUTS	30.0	1.04
BIG LSTM+CNN INPUTS + CNN SOFTMAX	39.8	0.29
BIG LSTM+CNN INPUTS + CNN SOFTMAX + 128-DIM CORRECTION	35.8	0.39
BIG LSTM+CNN INPUTS + CHAR LSTM PREDICTIONS	47.9	0.23

Results for Ensemble of Models

MODEL	TEST PERPLEXITY
LARGE ENSEMBLE (CHELBA ET AL., 2013)	43.8
RNN+KN-5 (WILLIAMS ET AL., 2015)	42.4
RNN+KN-5 (JI ET AL., 2015A)	42.0
RNN+SNM10-SKIP (SHAZEER ET AL., 2015)	41.3
LARGE ENSEMBLE (SHAZEER ET AL., 2015)	41.0
OUR 10 BEST LSTM MODELS (EQUAL WEIGHTS)	26.3
OUR 10 BEST LSTM MODELS (OPTIMAL WEIGHTS)	26.1
10 LSTMS + KN-5 (EQUAL WEIGHTS)	25.3
10 LSTMS + KN-5 (OPTIMAL WEIGHTS)	25.1
10 LSTMS + SNM10-SKIP (SHAZEER ET AL., 2015)	23.7

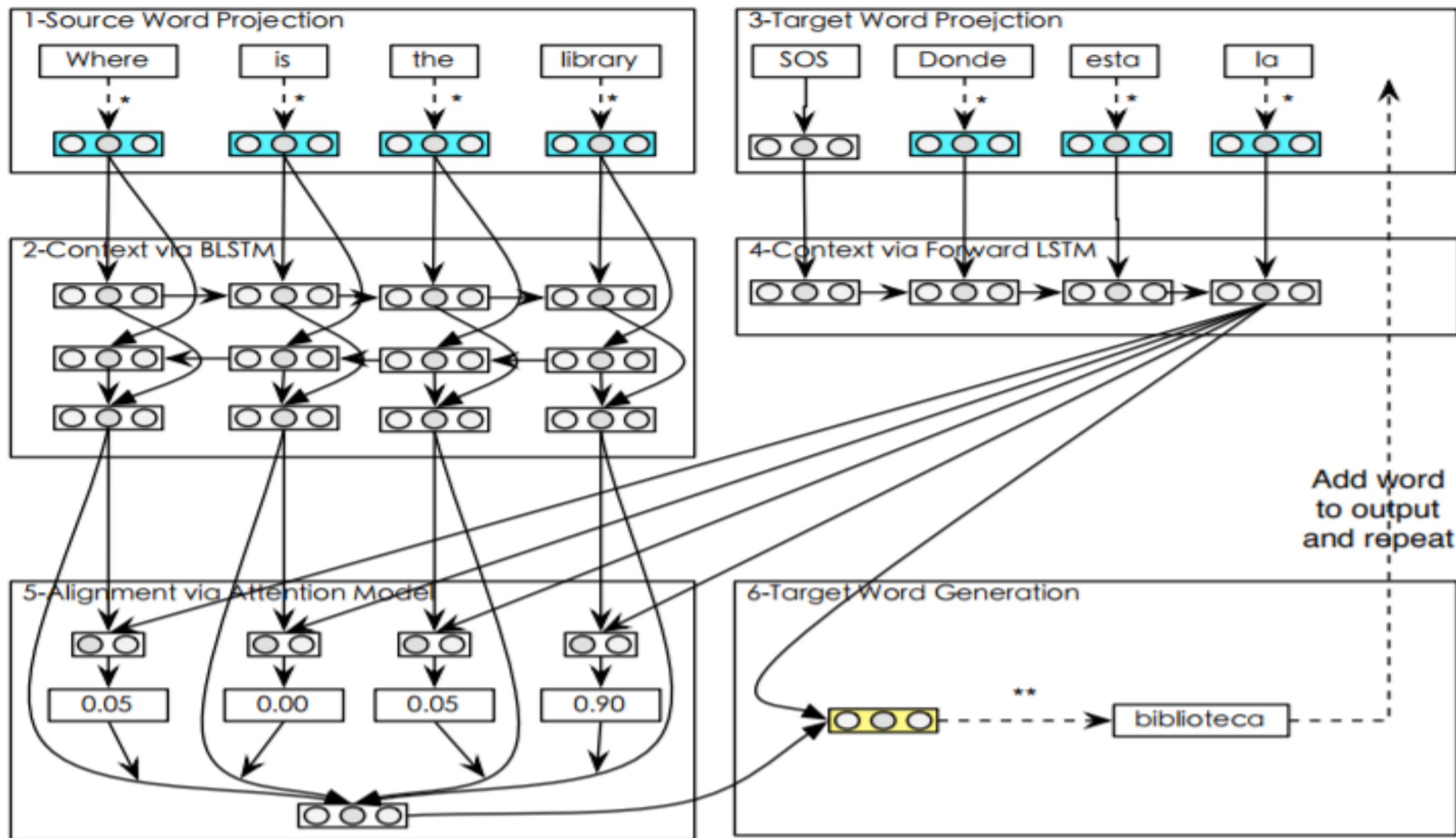
Key Observations

- Size matters.
- Regularization importance.
- Word Embeddings vs Character CNN.
- Smaller models with CNN Softmax.

Character-Based Neural Machine Translation

- Proposed by Wang Ling, Isabel Trancoso, Chris Dyer and Alan Black.
- The input and output sentences are sequences of characters rather than words.
- Proposed models are capable of generating and interpreting unseen word forms.
- Less data preprocessing of the source and target languages.

Joint Alignment and Translation Model



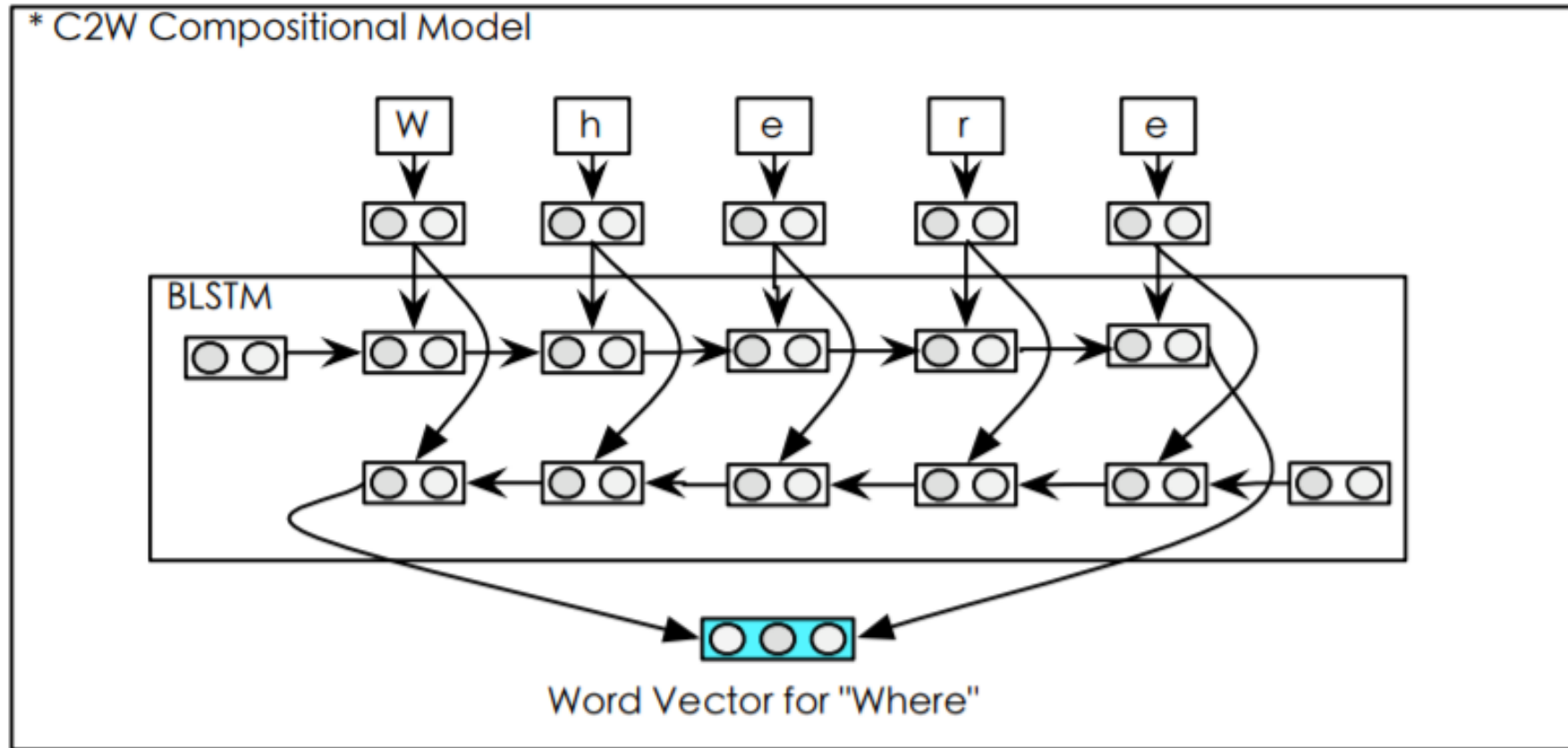
Alignment via Attention

- The model computes a score for each source word and a softmax over all the scores is performed.

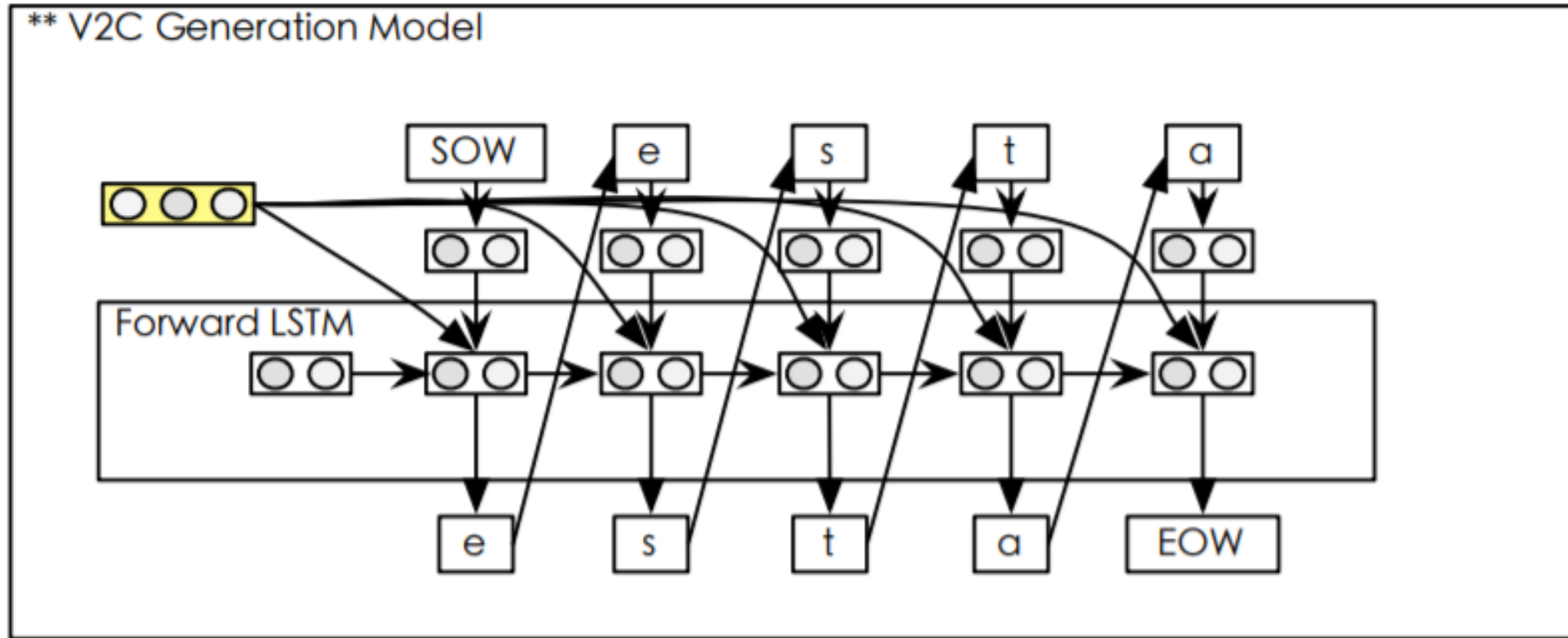
$$z_i = \mathbf{s} \tanh(\mathbf{W}_t \mathbf{l}_{p-1}^f + \mathbf{W}_s \mathbf{b}_i)$$

$$a_i = \frac{\exp(z_i)}{\sum_{j \in [0, n]} \exp(z_j)}$$

Character-Based Word Representation (C2W)



Character-Based Word Generation (V2C)



Layer-wise Training

- Firstly, characters are composed into word vectors using the C2W model.
- Then, the attention model searches for the next source word to translate.
- Finally, the generation of the target word is obtained using the V2C model

Results (1)

- C2W prefers to gather words that are orthographically similar.

C2W Model			Word Lookup Table		
<i>answer</i>	<i>well-founded</i>	<i>responder</i>	<i>answer</i>	<i>well-founded</i>	<i>responder</i>
response	well-balanced	respondeu	reply	described	reagir
answers	much-needed	responderam	response	impressed	aderir
reply	self-employed	responda	answering	bizarre	responda
answered	uncontrolled	responde	join	santer	aceder
replies	inherited	respondem	question	unclear	agradecer

Ref: Character-Based Neural Machine Translation (Ling et al.)

Results (2)

- V2C model generates unseen words.

original	... that does not mean that we want to bring an end to subsidisation .
character translation	... isso não significa a questão de que se trata de um fim à subsidade .
word translation	... Isso não significa que isso , para conseguir reduzir os autores .
original	the budget for the reconstruction of ...
character translation	o orçamento para as reconstruções ...
word translation	o orçamento inerente à reconstrução ...

Ref: Character-Based Neural Machine Translation (Ling et al.)

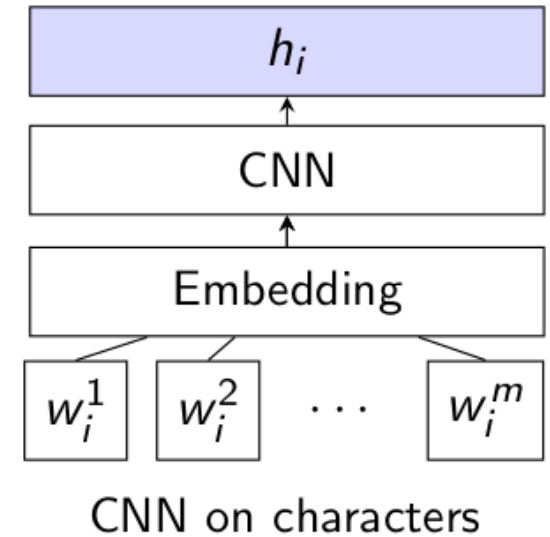
Conclusion

- Using characters as atomic units helps in learning orthographically sensitive word representations.
- Unseen word forms can be generated and thus, showing that character-based neural translation models perform better than word-based models.

Review

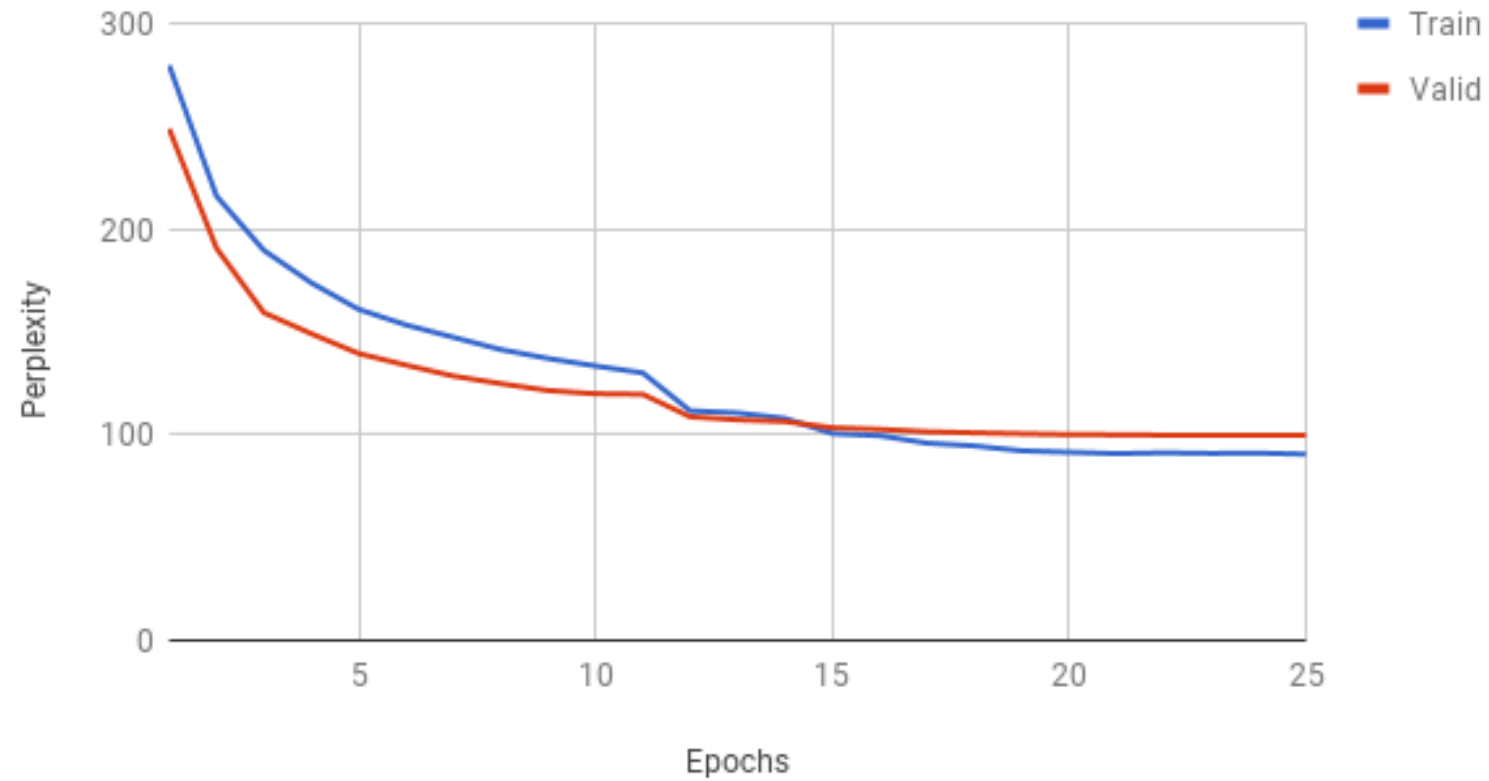
Our 4 proposed models:

- Word embeddings
- Concatenation of character embeddings
- CNN on character embeddings
- RNN on character embeddings



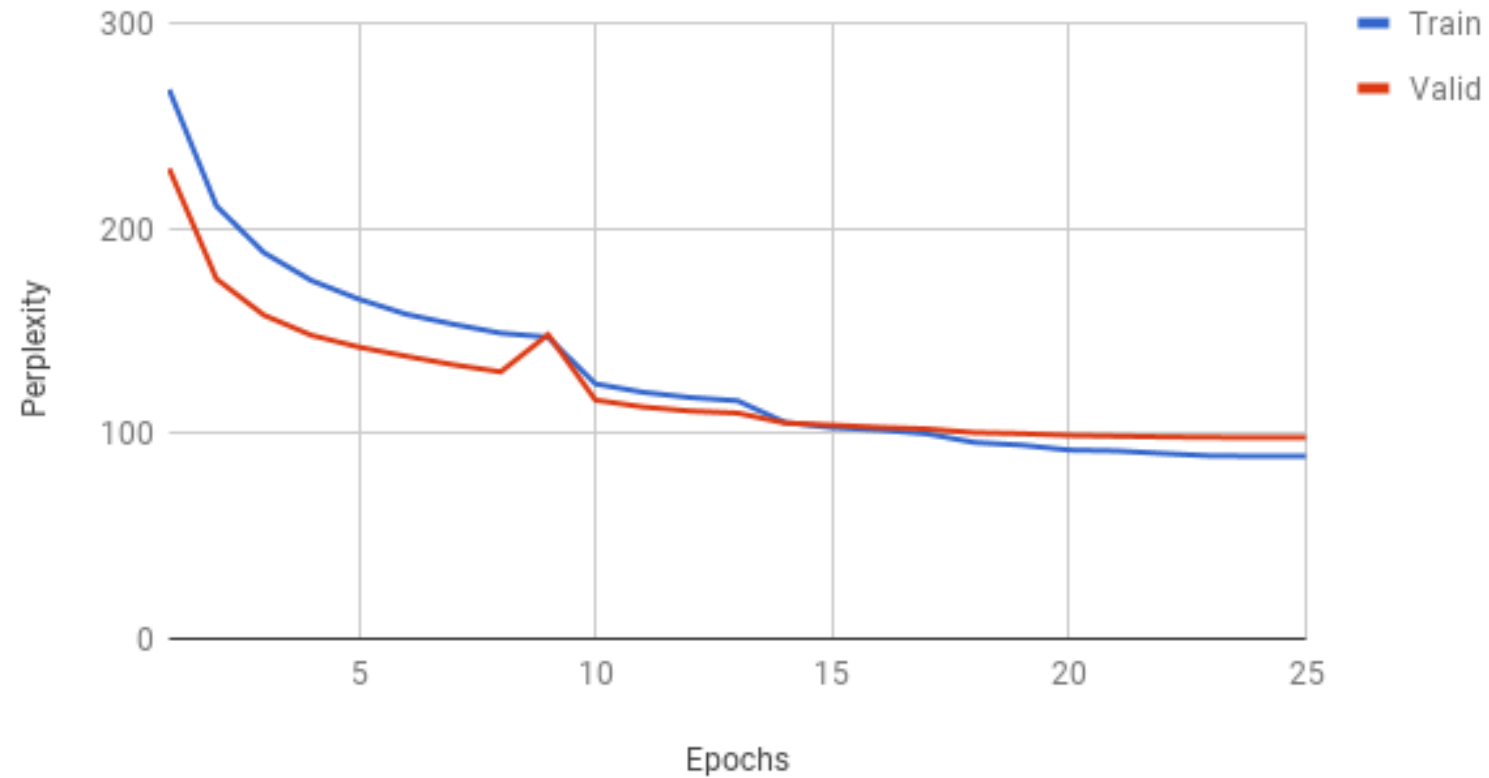
Word embeddings

Word Embeddings



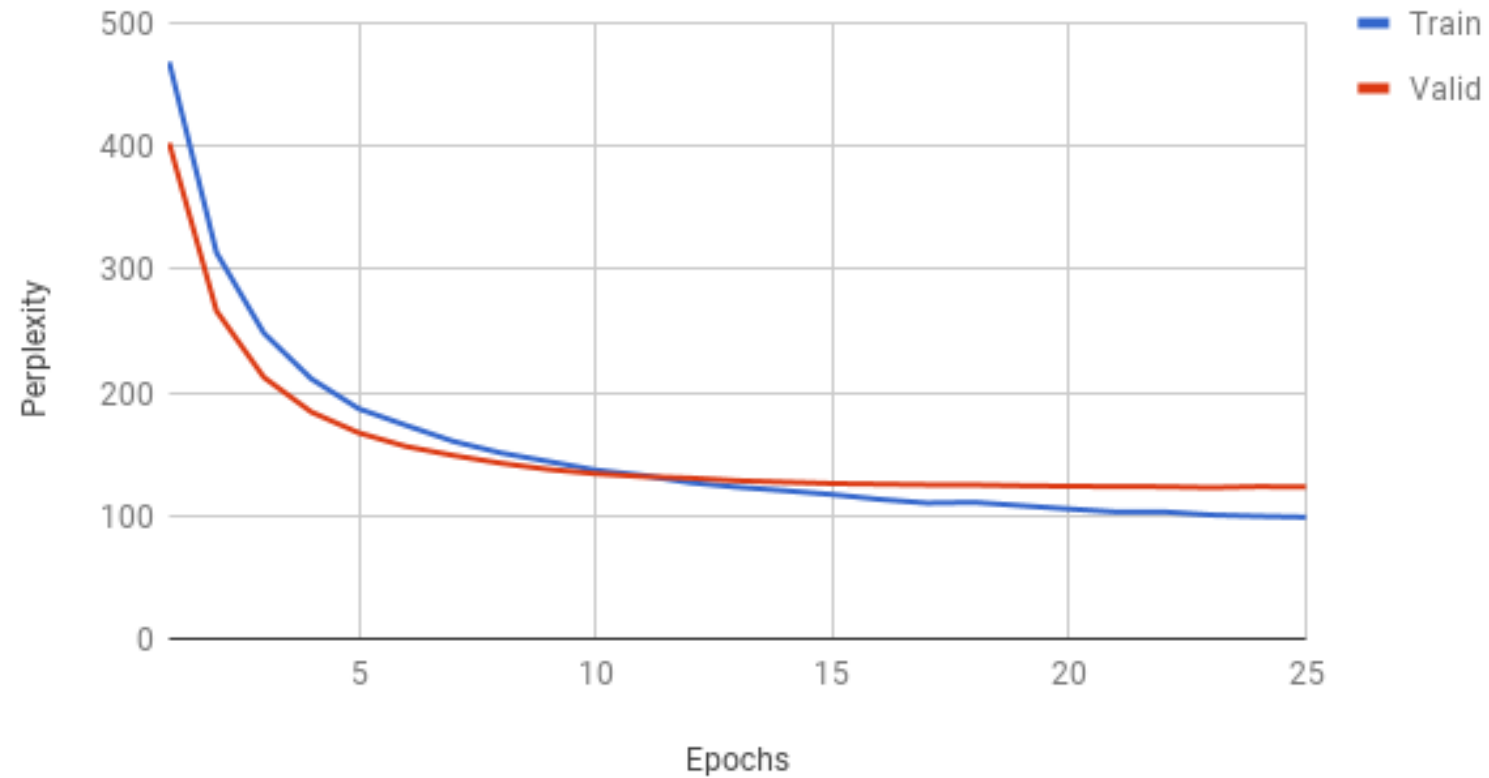
CNN over character embeddings

CNN over character embeddings



Concatenated character embeddings

Character Embeddings



Preliminary Results

- Results reported on the **English Penn Treebank** (PTB), and compared with *Character-Aware Neural Language Models* (Yoon Kim).

Model	Reported	Reproduced	Our Implementation
Word embeddings	97.6	95.814389158859	94.82428529
Concatenation of character embeddings	N/A	N/A	114.72031261
CNN over character embeddings	92.3	92.6425160752	94.28213722