

# Reading Assignment 1

Hamza Rashid  
COMP 550, Fall 2024

The paper "Neural Architectures for Named Entity Recognition" by Lample et al. (2016)., addressed the reliance of state-of-the-art Named Entity Recognition (NER) systems on hand-crafted features and domain-specific knowledge, a common approach at the time due to the small, supervised training corpora that was available. The authors proposed two neural architectures designed to generalize under the limited training corpora: a bidirectional-LSTM encoder pipelining input sequences to a CRF layer, and, a greedy, transition-based chunking algorithm utilizing a Stack-LSTM to manage states. A key component of their methodology was the formation of character-sensitive word embeddings, to capture orthographical and morphological details. While the proposed models achieved state-of-the-art performance, some aspects of the methodology limit the scope for generalization. We continue in more detail.

The Bi-LSTM encoder consists of a forward and backward LSTM, reading the input sequence in those orders, respectively. The resulting left and right representations are then concatenated into a single vector, the word-in-context representation. The bi-directional architecture is critical for making informative encodings, as the left and right encodings may be insufficient, or capture the wrong *context*, if considered in isolation (homographical richness is an intuitive case). Furthermore, the CRF layer utilizes the formulation discussed in class, where the feature sum is taken over the state-transition and observation-emission (given by the Bi-LSTM) scores. The linear-chain CRF's ability to capture the bi-directional relationships between output labels makes it effective in addressing two problems that are closely related to classification accuracy: the strong dependence between output labels (inherent to the NER task), and, the formal grammar induced by the IOBES (Inside, Outside, Beginning, End, Single) tagging scheme utilized by the authors (e.g., I-PER cannot follow B-LOC).

describe transition-based chunking model

In designing the input word embeddings, the authors found character-level treatment crucial for identifying orthographic or morphological evidence that something is a name (or not a name). This was achieved by initializing random embeddings for each character, then feeding each token as a character sequence into a Bi-LSTM – with the final embedding given by the concatenation of the left and right contexts. This embedding is then concatenated with pretrained word embeddings from a large corpus. During testing, words without an embedding in the lookup table are mapped to an UNK embedding, and singletons with the UNK embedding are assigned a probability of 0.5. This approach has been found useful in handling OOV items, as the embedding nonetheless captures character-level details. This differs from techniques shown in class, which deal with OOV items either at the token level (UNK mapping, but no character-level embedding), or at the hyperparameter-tuning level (e.g., the Naive Bayes smoothing parameter). and, embeddings learned from a large corpus that are sensitive to word order.

discuss results (and advantages/disadvantages of using external gazetteers)

## Strengths and Limitations

The researchers achieved the original task, - Strengths: - Language-independent, does not rely on hand-crafted features or gazetteers. - Achieves state-of-the-art results in multiple languages. - Effective use of character-based and word embeddings to handle morphology. - Limitations: - The transition-based chunking model is more dependent on character-based information compared to the LSTM-CRF. - Greedy action selection in the Stack-LSTM model can lead to suboptimal results. The paper includes a detailed outline of the methodologies, and provides strong justifications for its preprocessing decisions, particularly at the input layer (arguing that LSTM's are an a priori better function class for modeling the relationship between words and their characters, as they take into account position-variant features) .

## Bidirectional LSTM-CRF Architecture (Figure 1)

- Describe the key components: - Bidirectional LSTM: Encodes contextual information from both left and right contexts. - Conditional Random Field (CRF): Models dependencies between tags to produce globally optimal sequences. - Explain the importance of these components: - Bidirectional LSTM captures comprehensive context for each word. - CRF layer ensures valid and coherent tag sequences.

## Handling of OOV Items

- Describe how the proposed method addresses out-of-vocabulary (OOV) words: - Uses character-level embeddings generated by a bidirectional LSTM to represent words based on their characters. - Incorporates pre-trained embeddings to handle unseen words by mapping them to a common UNK embedding during training. - Compare to class discussions: - Similar to character-level models we discussed, which also leverage character features to address OOV problems. - Pre-trained embeddings are akin to word2vec embeddings we discussed for capturing distributional semantics.

## Use of Gazetteers (Table 1)

- Discuss methods incorporating gazetteers to improve NER performance: - Gazetteers can provide explicit, domain-specific named entity information, helping models generalize better. - Advantages: Improves recognition accuracy for specific entity types, especially in low-resource settings. - Disadvantages: Dependence on domain-specific resources reduces language independence and increases cost for new domains or languages.

## Conclusion

- Summarize the overall contribution of the paper. - Highlight the effectiveness of neural architectures for NER without relying on language-specific features.