

| Project Proposals | This Lecture |
|--|---|
| ~1 page Define a problem, give context of related work (at least 3-4 relevant papers) Propose a direction that you think is feasible and outline steps to get there, including what dataset you'll use | Neural CRFs Tagging / NER Parsing |
| Okay to change directions after the proposal is submitted, but run it by me if it's a big change | |













| ۵ "N | LP (<i>f</i> | Almos | t) Fr | om So | cratch" | | How do w | ve use a tag | ger for | SRL? |
|--|-------------------------------|--------------------------------------|--------------------------|-----------------|--|---------|----------------------|-----------------------|--------------|-------------------------|
| Approach | POS (PWA) | CHUNK (F1) | NER (F1) | SRL (F1) | $\begin{array}{c c} \textbf{Input Window} & & & & \\ \hline \text{Text} & \text{cat sat on the mat} \\ \text{Feature 1} & w_1^1 & w_2^1 & \dots & w_N^1 \\ \vdots \end{array}$ | Gold | ARG1 | VA | RG2 | ARG3 |
| Benchmark Systems | 97.24 | 94.29 | 89.31 | 77.92 | Feature K $w_1^K w_2^K \dots w_N^K$ | | | | | |
| NN+WLL | 96.31 | 89.13 | 79.53 | 55.40 | Lookup Table | | Housing starts are e | expected to quicken a | ι bit from A | ugust's pace |
| NN+SLL | 96.37 | 90.33 | 81.47 | 70.99 | $LT_{W^1} \rightsquigarrow$ | | | | | |
| NN+WLL+LM1 | 97.05 | 91.91 | 85.68 | 58.18 | d | | | | | |
| NN+SLL+LM1 | 97.10 | 93.65 | 87.58 | 73.84 | $LT_{W^K} \rightsquigarrow$ | Taggir | ng problem with re | espect to a particul | ar verb | |
| NN+WLL+LM2 | 97.14 | 92.04 | 86.96 | 58.34 | concat | | 01 | | | |
| NN+SLL+LM2 | 97.20 | 93.63 | 88.67 | 74.15 | Linear | Con't | do this with foodf | anward nativaries a | fficiently | raumonto ara ta a |
| WLL: independen LM1/LM2: pretrai a language model | t classi ned wo over la | fication; S ord embe arge corp | SLL: ne ddings ora | ural CRF | HardTanh Linear $M^2 \times 0 \longrightarrow \frac{n_{hu}}{n_{hu}^2} = \theta \tan \theta$ | far fro | om the verb to use | e fixed context wind | dow sizes | inguments are too |
| | | | | Collober | , Weston, et al. 2008, 2011 | | | | Figure | e from He et al. (2017) |



| 🛞 Ho |)w " | from s | scrat | tch" | was this? | Neural CRFs with LSTMs |
|---------------------|----------------|----------------|--------------------|--------------------|--------------------------------|--|
| Approach | POS (PWA) | CHUNK (F1) | NER (F1) | SRL (F1) | ▶ NN+SLL isn't great | Neural CRF using character LSTMs to compute word representations |
| Benchmark Systems | 97.24 | 94.29 | 89.31 | 77.92 | LM2: trained for 7 weeks on | |
| NN+WLL NN+SLL | 96.31 96.37 | 89.13 90.33 | 81.47 | 55.40 70.99 | Wikipedia+Reuters — very | |
| NN+WLL+LM1 | 97.05 | 91.91 | 85.68 | 58.18 | expensive! | $\begin{array}{ c c c c c }\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_3 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_4 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_4 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_4 & c_4 & c_4 & \text{the bedding from } \\\hline c_1 & c_2 & c_4 & $ |
| NN+SLL+LM1 | 97.10 | 93.65 | 87.58 | 73.84 | | |
| NN+WLL+LM2 | 97.14 | 92.04 | 86.96 | 58.34 | Sparse features needed | $BiLSTM = \left(\begin{array}{c} r_1 \\ r_2 \end{array} \right) + \left(\begin{array}{c} r_3 \\ r_4 \end{array} \right) + \left(\begin{array}{c} r_4 \\ r_4 \end{array} \right) + \left(\begin{array}{c} r_1 \\ r_4 \end{array} \right) + \left(\begin{array}{c} r_$ |
| NN+SLL+LM2 | 97.20 | 93.63 | 88.67 | 74.15 | to get best performance | |
| NN+SLL+LM2+Suffix2 | 97.29 | - | - | - | | |
| NN+SLL+LM2+Gazettee | r _ | - | 89.59 | - | OII NER+SRL dilywdy | |
| NN+SLL+LM2+POS | - | 94.32 | 88.67 | - | No use of sub word | |
| NN+SLL+LM2+CHUNK | - | - | - | 74.72 | No use of sub-word | |
| | | | | | teatures | embeddings [Mark Watney visited Mars M a r s |
| | | | | C | ollobert and Weston 2008, 2011 | Chiu and Nichols (2015), Lample et al. (2016 |

| Neural CRFs | s with LSTMs | | |
|--|---|--|-------------------------|
| Chiu+Nichols: character CNNs instead of LSTMs Lin/Passos/Luo: use external resources like Wikipedia | Model Collobert et al. (2011)* Lin and Wu (2009) Lin and Wu (2009)* Huang et al. (2015)* Passos et al. (2014) | F ₁ 89.59 83.78 90.90 90.10 90.05 | |
| LSTM-CRF captures the important aspects of NER: word context (LSTM), sub-word features | Passos et al. $(2014)^*$ Luo et al. $(2015)^* + \text{gaz}$ Luo et al. $(2015)^* + \text{gaz} + \text{linking}$ Chiu and Nichols (2015) Chiu and Nichols $(2015)^*$ | 90.90 89.9 91.2 90.69 90.77 | Neural CRFs for Parsing |
| (character LSTMs), outside knowledge (word embeddings) | LSTM-CRF (no char) LSTM-CRF Chiu and Nichols (2015), Lample | 90.20 90.94 et al. (2016) | |















| Туре | Model | English UAS | PTB-SD 3.3.0 LAS | Chinese UAS | e PTB 5. LAS |
|-------------|----------------------------------|----------------|---------------------|----------------|-----------------|
| Transition | Ballesteros et al. (2016) | 93.56 | 91.42 | 87.65 | 86.21 |
| | Andor et al. (2016) | 94.61 | 92.79 | | - |
| | Kuncoro et al. (2016) | 95.8 | 94.6 | | - |
| Graph | Kiperwasser & Goldberg (2016) | 93.9 | 91.9 | 87.6 | 86.1 |
| | Cheng et al. (2016) | 94.10 | 91.49 | 88.1 | 85.7 |
| | Hashimoto et al. (2016) | 94.67 | 92.90 | - | - |
| | Deep Biaffine | 95.74 | 94.08 | 89.30 | 88.23 |
| Riaffine ar | - Annoach works well (other r | oural C | REs are also | strong | |

| | Νει |
|--|-----|
| | |

Neural CRFs

State-of-the-art for:

POS

۲

- NER without extra data (Lample et al.)
- Dependency parsing (Dozat and Manning)
- Semantic Role Labeling (He et al.)
- Why do they work so well?
- Word-level LSTMs compute features based on the word + context
- Character LSTMs/CNNs extract features per word
- Pretrained embeddings capture external semantic information
- CRF handles structural aspects of the problem

| | Takeaways |
|------------------------|---|
| Any struc compute | tured model / dynamic program + any neural network to potentials = neural CRF |
| Can incor like gram | porate transition potentials or other scores over the structure mar rules |
| State-of-t | the-art for many text analysis tasks |