NLP

# Named Entity Recognition

# NER做什么

| 国务院 | 总理 | 李克强 | 调研 | 上海 | 外高桥 | 时 | 提出 | ， | 支持 | 上海 | 积极 | 探索 | 新 | 机制 | 。 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 机构 | | 人名 | | 地名 | | | | | | 地名 | | | | | |

| 中国 | 银行 | 西藏 | 自治区 | 分行 | 行长 | 索朗达吉 |
|---|---|---|---|---|---|---|
| B-ORG | I-ORG | B-GPE | O | O | O | B-PER |

- NER标注一般采用如下标注形式
  - O 非实体
  - S 单独构成实体
  - B 为实体的开始
  - I 为实体的中间
  - E 为实体的结尾

# NER的发展

- 基于词典和规则的方法                          [略]
- 传统机器学习
    - 无监督学习
    - CRF
- Emb[NN/Attention]-NN-CRF/LSTM            [重点]
- 少量标注的训练集的研究                        [重点]
    - 半监督学习
    - 主动学习
    - 对抗生成网络
    - 迁移学习
- 主要的问题[研究热点]
    - 标注集少
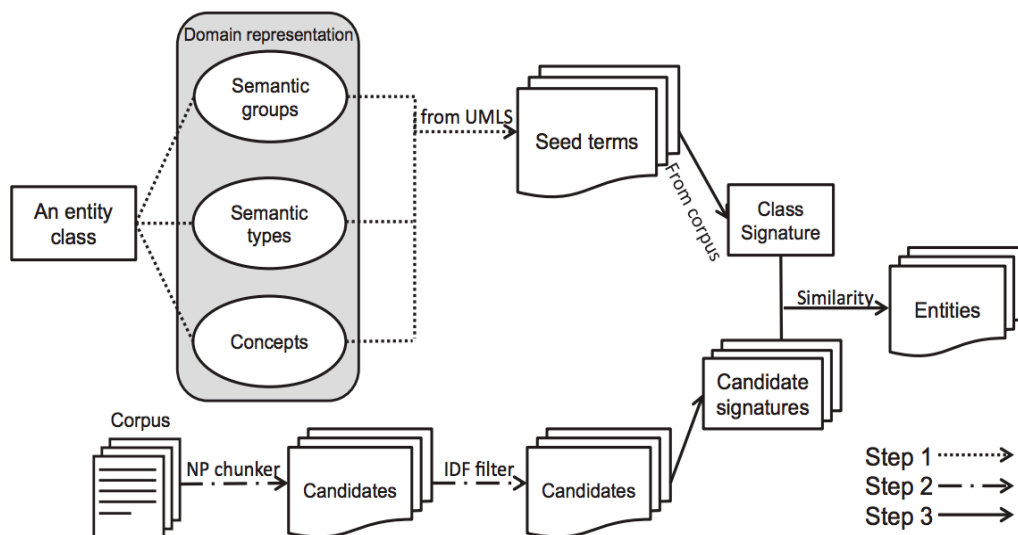    - 标注结果准确率与成本的平衡

# 无监督学习[1][13]

- 不用标注数据，按数据的相似度聚类



**Fig. 1.** Overall approach to unsupervised biomedical named entity recognition.

**Step 1：Seed term collection**

    通过对实体类的语义组，类型，概念从UMLS中获取相关类的术语实体作为种子术语。

**Step 2：Boundary detection**

    对Corpus进行分块，获取名词短语，使用IDF过滤

$$IDF(t, D) = log(|D|/|d \in D : t \in d|)$$

**Step 3：Entity classification**

    构建词表大小V（all possible unigrams）

    t[step1和step2中的结果]用特征维度为2V的$s^t$表示：

$$s^t = \langle s_1^t, s_2^t, \ldots, s_V^t, s_{V+1}^t, \ldots, s_{2V}^t \rangle \quad (2)$$

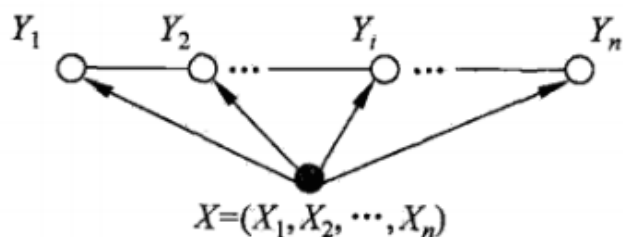Values in the vector are calculated as follows:

$$s_i^t = w_i * f(v_i, t) * IDF(v_i, D), \quad i = 1 \ldots V \quad (3)$$

$$s_i^t = w_o * f(v_i, context_t) * IDF(v_i, D), \quad i = V + i \ldots 2V \quad (4)$$

类别特征按类取平均，候选特征和类别特征用余弦相似度聚类

改进：可以使用K-means继续优化

[Shaodian Zhang et al.2013] Unsupervised biomedical named entity recognition: Experiments with clinical and biological texts

# CRF

- 观测序列X[未标记序列]和状态序列Y[标记序列]

- 转移概率：$Y_i \rightarrow Y_{i+1}$ ， 状态概率：$X \rightarrow Y_i$

- CRF：训练集中寻找实体上下文的模板，测试集中将上下文和模板匹配的中间词作为对应实体标注



$X = (X_1, X_2, \cdots, X_n)$

$$P(y|x) = \frac{1}{Z(x)} \exp\left( \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right)$$
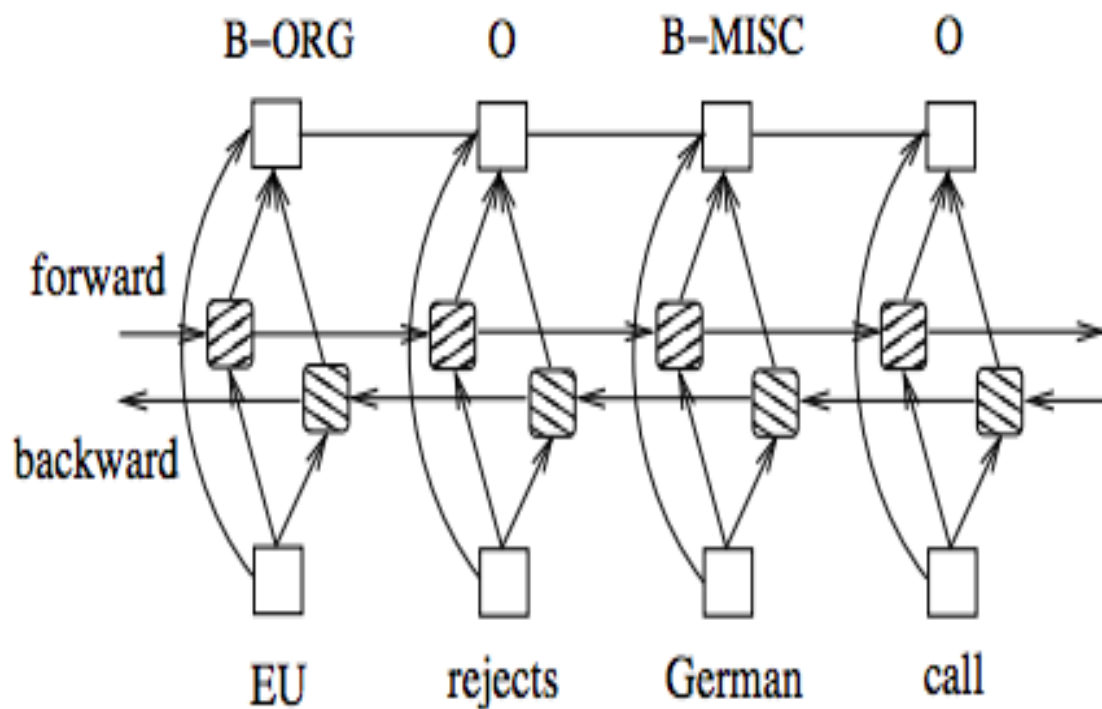
$$Z(x) = \sum_y \exp\left( \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right)$$

```
features = [
    'bias',
    'word='+word,
    # 'word.lower=' + word.lower(),
    # 'word[-3:]=' + word[-3:],
    # 'word[-2:]=' + word[-2:],
    # 'word.isupper=%s' % word.isupper(),
    #'word.istitle=%s' % word.istitle(),
    'word.isdigit=%s' % word.isdigit(),
    'postag=' + postag,
    'cuttag=' + cuttag,
    # 'postag[:2]=' + postag[:2],
]
if i > 0:
    word1 = sent[i - 1][0]
    postag1 = sent[i - 1][2]
    cuttag1 = sent[i - 1][1]
    features.extend([
        '-1:word='+word1,
        '-1:postag=' + postag1,
        '-1:cuttag=' + cuttag1,
        # '-1:postag[:2]=' + postag1[:2],
    ])
else:
    features.append('BOS')

if i < len(sent) - 1:
    word1 = sent[i + 1][0]
    postag1 = sent[i + 1][2]
    cuttag1 = sent[i + 1][1]
    features.extend([
        '+1:word=' + word1,
        '+1:postag=' + postag1,
        '+1:cuttag=' + cuttag1,
```

一般的CRF线性链如左图所示，中间的图是实际中使用CRF进行NER的特征提取代码，比如要预测$Y_i$，word为$X_i$，可以看到除了对$X_i$提取特征外，还涉及到了$X_i$前面的一个词和后面一个词。

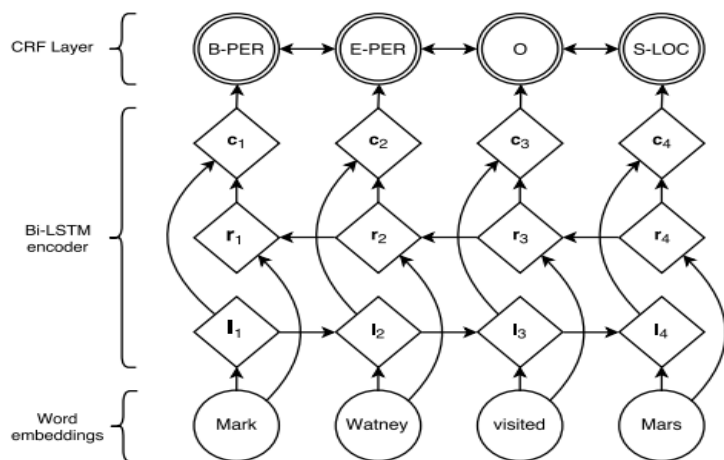但是并没有将全部的上下文考虑进去，这是CRF问题，后面的NN结构解决了该问题

# Emb-BiLSTM-CRF[3][15]



B-ORG    O    B-MISC    O

forward

backward

EU    rejects    German    call

Embedding：使用前置训练，130K Vocabulary，50 维

直接将BiLSTM的输出传给CRF层

CRF层的输入：手工添加了很多特征，比如spelling,context等特征

[Zhiheng Huang et al.2015] Bidirectional LSTM-CRF Models for Sequence Tagging

# Emb[LSTM]-BiLSTM-CRF[5][16]



句子的长度为 n，不同的标记个数为 k，将 Bi-LSTM 输出（要经过 softmax 层，变成 k 维的向量）作为打分矩阵 P，

即 $P_{i,j}$ 表示句子的第 i 个单词对应的是第 j 个标签的分数

$$X = (x_1, x_2, ..., x_n)$$

$$y = (y_1, y_2, ..., y_n)$$

$$s(X, y) = \sum_{i=0}^{n} A_{y_i, y_{i+1}} + \sum_{i=1}^{n} P_{i, y_i}$$
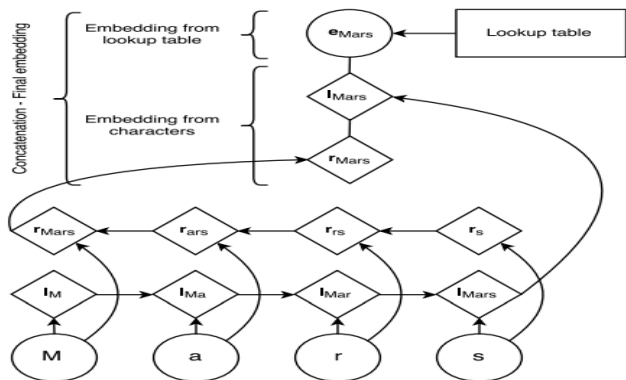
A 表示的是转移分数，$A_{i,j}$ 表示从第 i 个标记转到第 j 个标记的分数

softmax：

$$p(y|X) = \frac{e^{s(X,y)}}{\sum_{y'} e^{s(X,y')}}$$

log 最大似然进行训练：

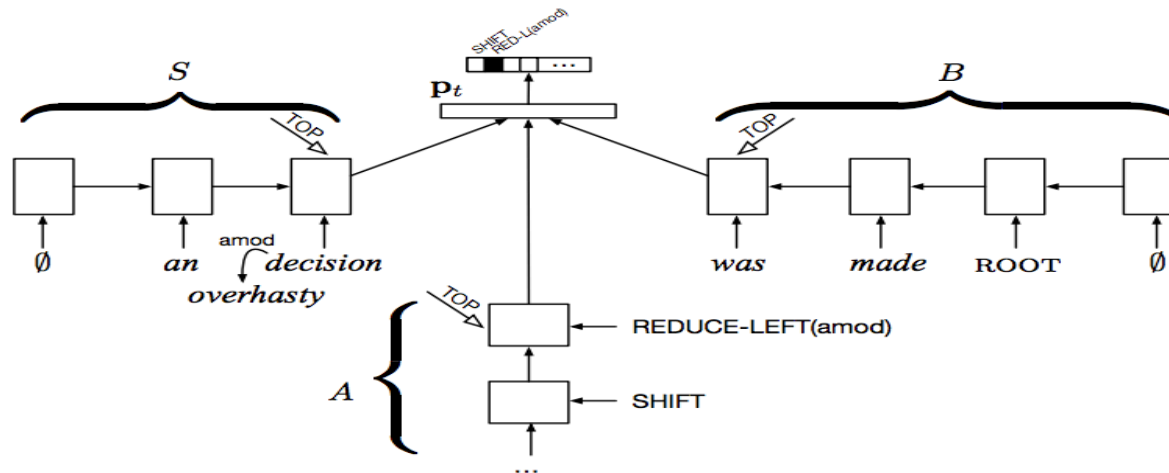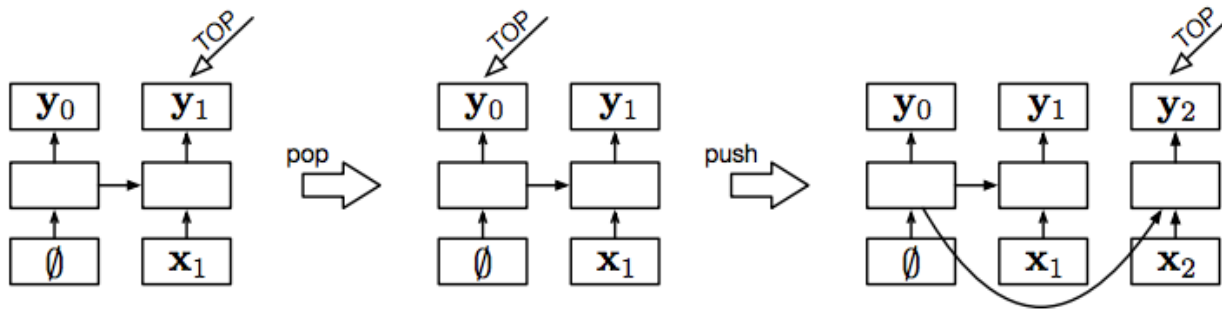$$\log(p(y|X)) = s(X, y) - \log \left( \sum_{y'} e^{s(X,y')} \right)$$

动态规划来预测：

$$y^* = \arg max_{y \in Y_x} s(X, y)$$

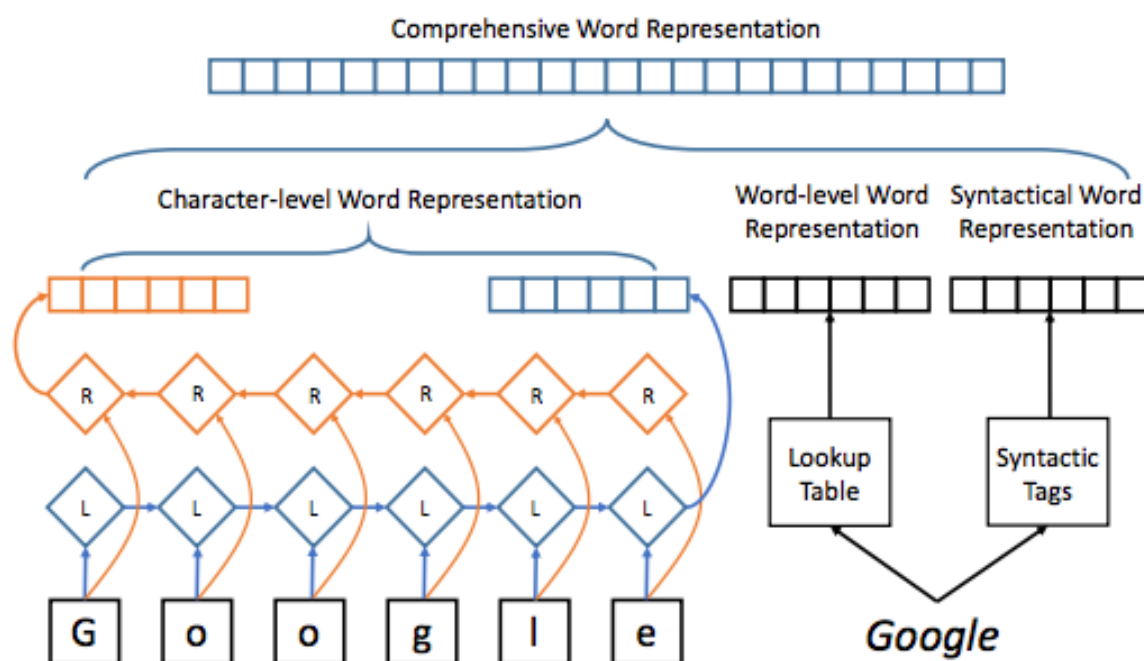[Guillaume Lample et al.2016] Neural Architectures for Named Entity Recognition

# Stack-LSTM[5][16]



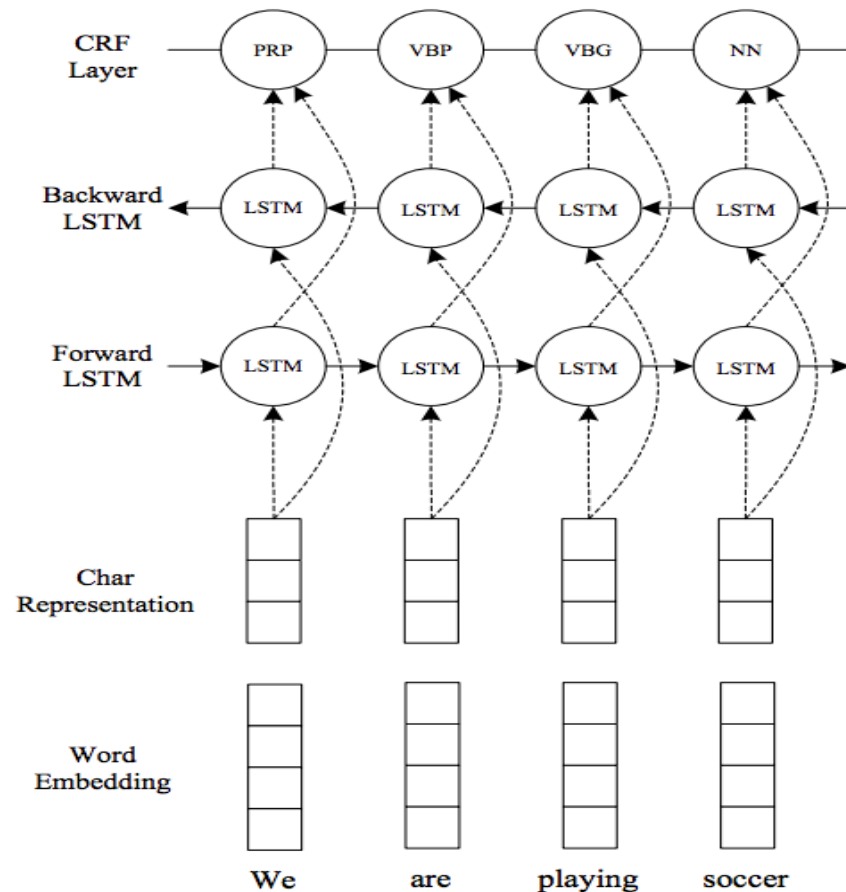[Guillaume Lample et al.2016] Neural Architectures for Named Entity Recognition
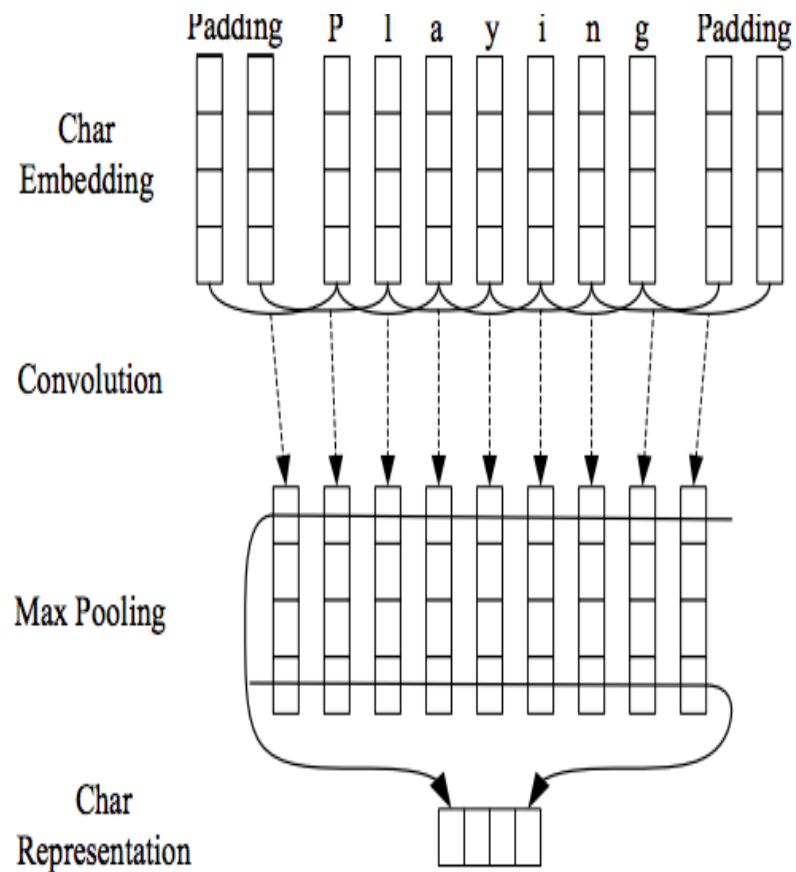[Chris Dyer et al.2015] Transition-Based Dependency Parsing with Stack Long Short-Term Memory

# Multi-channel-BiLSTM-CRF[15][17]



主要是嵌入层的改变，如左图，除了加入字符嵌入，词嵌入外，还加入了词性，句法，语义等信息。

[Bill Y.Lin et al.2017] Multi-channel BiLSTM-CRF Model for Emerging Named Entity Recognition in Social Media

# Emb[CNN]-BiLSTM-CRF[6][16]



[Xuezhe Ma et al.2016] End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF
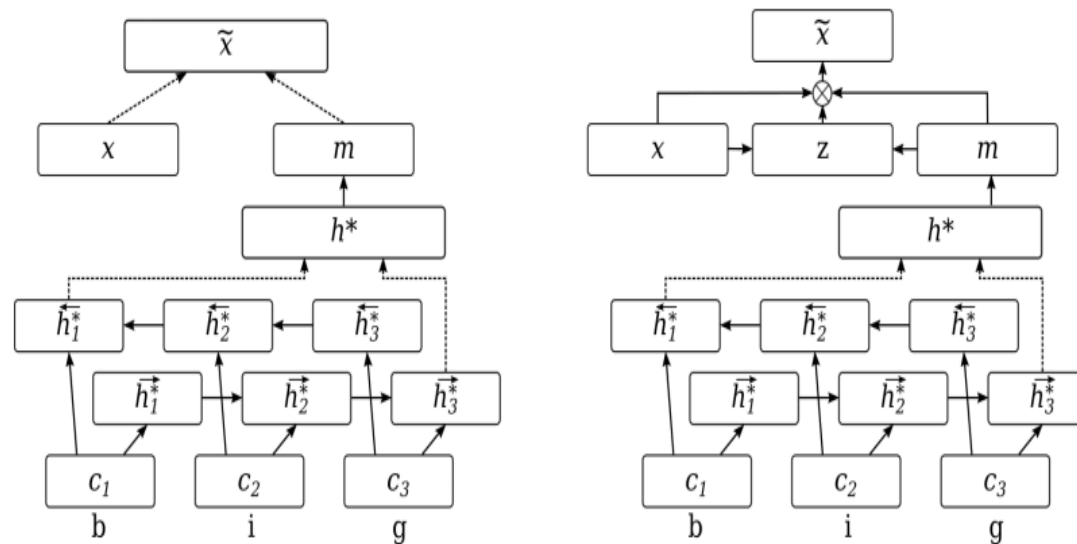
# Emb[Attention]+BiLSTM+CRF[7][16]



Figure 2: Left: concatenation-based character architecture. Right: attention-based character architecture. The dotted lines indicate vector concatenation.

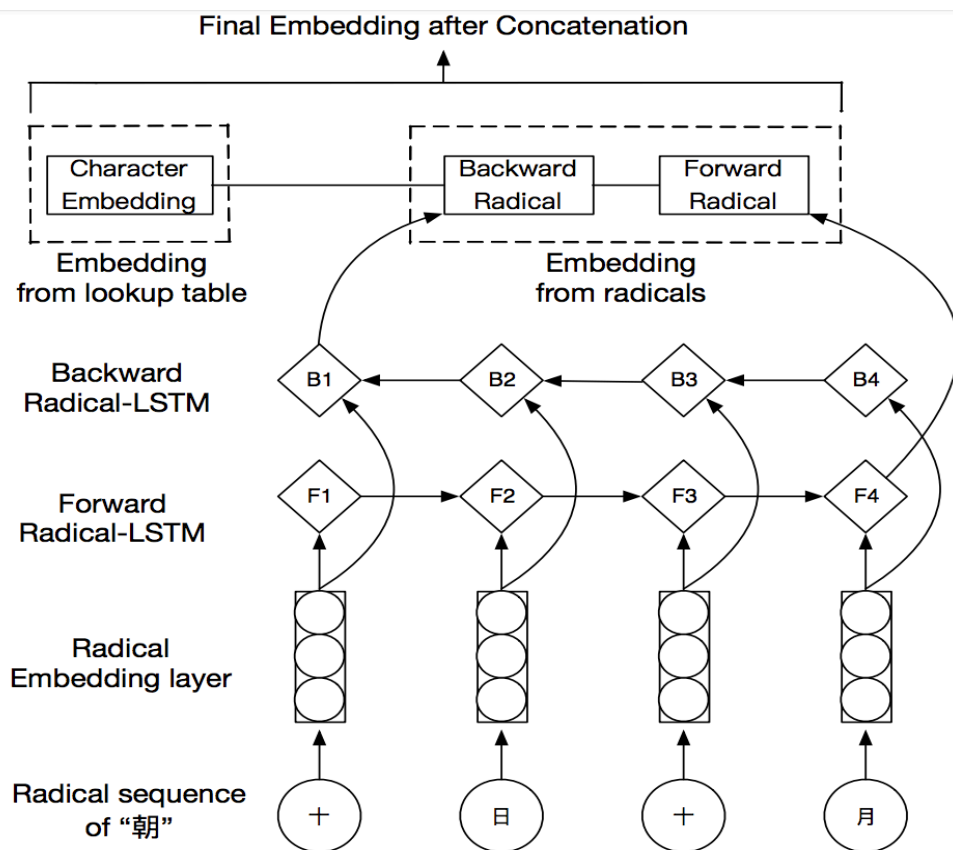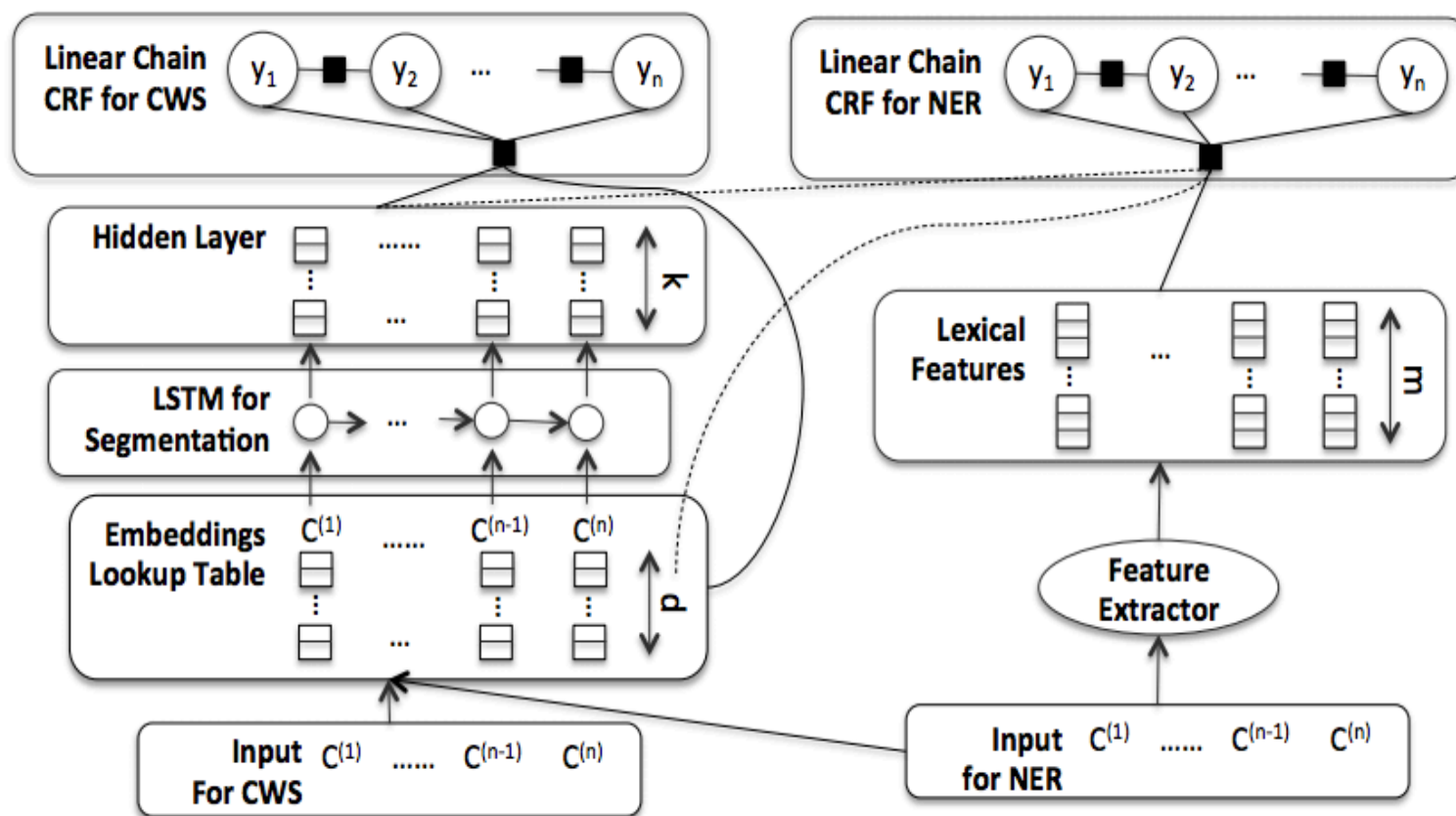$$z = \sigma(W_z^{(3)} tanh(W_z^{(1)}x + W_z^{(2)}m)) \qquad \widetilde{x} = z \cdot x + (1 - z) \cdot m$$

[Marek Rei et al.2016] Attending to Characters in Neural Sequence Labeling Models

# Emb[LSTM]-BiLSTM-CRF[10][16]



Final Embedding after Concatenation

Character Embedding — Embedding from lookup table

Backward Radical — Forward Radical — Embedding from radicals

Backward Radical-LSTM: B1 ← B2 ← B3 ← B4

Forward Radical-LSTM: F1 → F2 → F3 → F4

Radical Embedding layer

Radical sequence of "朝": 十 日 十 月

Table 3: Results with different components.

| Variant | F1 |
|---|---|
| random + dropout | 88.91 |
| random + radical + dropout | 89.44 |
| pretrain + dropout | 90.75 |
| pretrain | 86.87 |

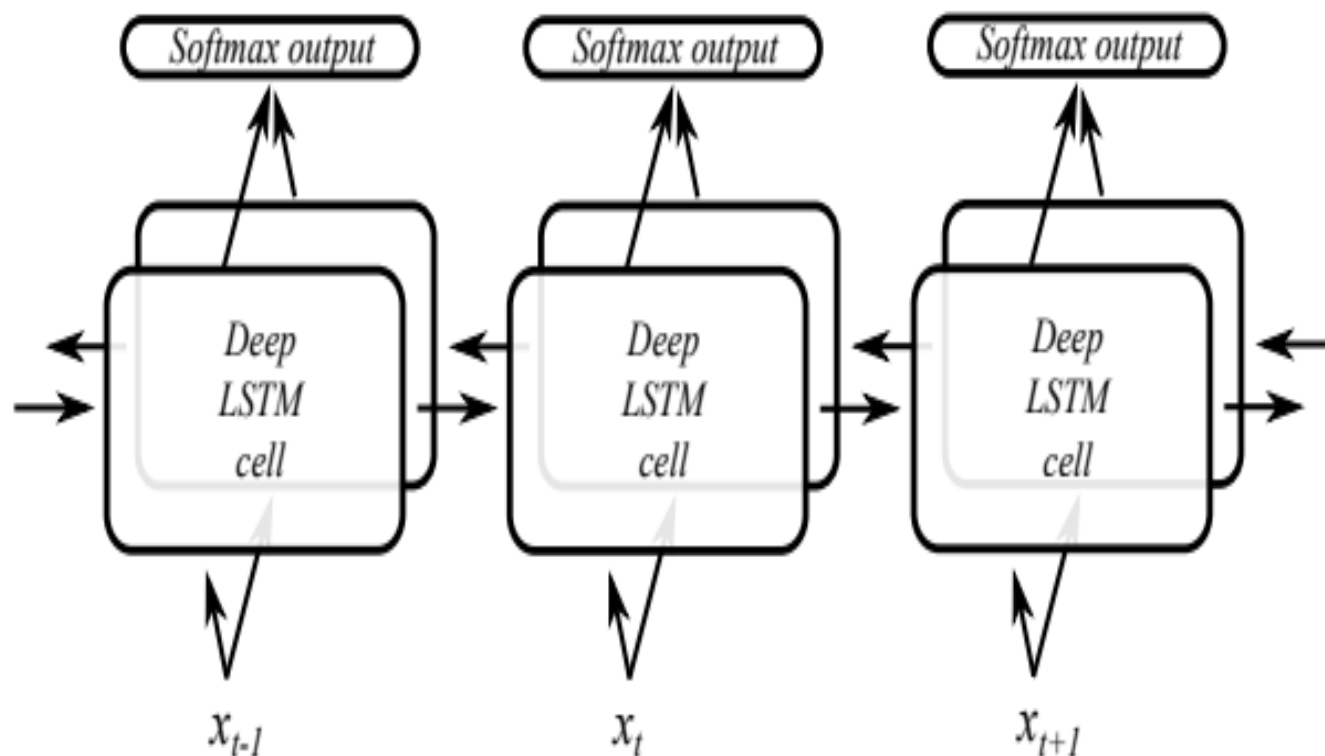| Model | PER-F | LOC-F | ORG-F | P | R | F |
|---|---|---|---|---|---|---|
| Zhou2006 | 90.09 | 85.45 | 83.10 | 88.94 | 84.20 | 86.51 |
| Chen2006 | 82.57 | 90.53 | 81.96 | 91.22 | 81.71 | 86.20 |
| Zhou2013 | 90.69 | 91.90 | 86.19 | 91.86 | 88.75 | 90.28 |
| Zhang2006* | 96.04 | 90.34 | 85.90 | 92.20 | 90.18 | **91.18** |
| BLSTM-CRF + radical | 89.62 | 91.76 | 85.79 | 91.39 | 88.22 | 89.78 |
| BLSTM-CRF + pretrain | 91.77 | **92.10** | **87.30** | 91.28 | **90.62** | **90.95** |

[Chuanhai Dong et al.2016] Character-Based LSTM-CRF with Radical-Level Features for Chinese Named Entity Recognition

# With Word Segmentation[8][16]



[Nanyun Peng et al.2016] Improving Named Entity Recognition for Chinese Social Media with Word Segmentation Representation Learning

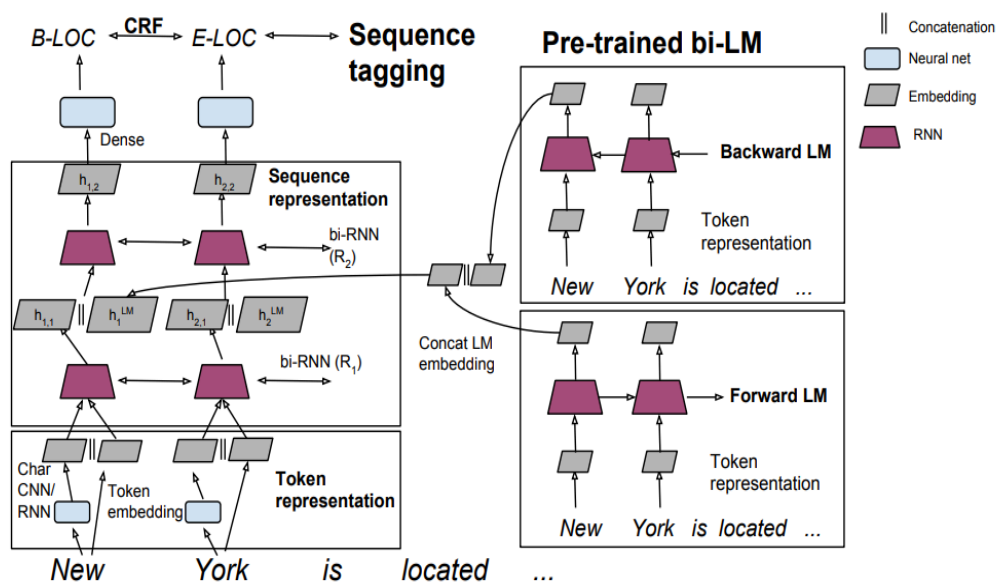# Emb[Character]-DBi-LSTM-Softmax[9][16]



这篇paper中发现以下三种实体：
disorders and findings[疾病 and …]
pharmaceutical drugs[药物]
body structure[身体结果]

[Simon Almgren et al.2016] Named Entity Recognition in Swedish Health Records with Character-Based Deep Bidirectional LSTMs

# Semi-supervised sequence tagging[11][17]

- 利用无标签数据来优化有标签数据训练的模型

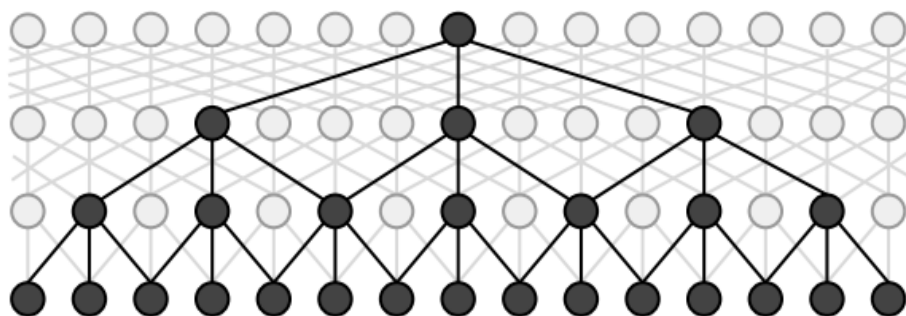| Model | $F_1 \pm$ **std** |
|---|---|
| Chiu and Nichols (2016) | $90.91 \pm 0.20$ |
| Lample et al. (2016) | $90.94$ |
| Ma and Hovy (2016) | $91.37$ |
| Our baseline without LM | $90.87 \pm 0.13$ |
| **TagLM** | $\mathbf{91.93 \pm 0.19}$ |

Table 1: Test set $F_1$ comparison on CoNLL 2003 NER task, using only CoNLL 2003 data and unlabeled text.

左边是NER的三层结构，中间层使用了两层BiLSTM
右边使用语言模型(前面词预测下一个词)训练词的上下文向量[语料是未标记的数据]

[Peters et al.2017] Semi-supervised sequence tagging with bidirectional language models

# Emb-IDCNN-CRF[12][17]

- 从CNN和RNN对句子特征提取方面出发，使用迭代扩张卷积替换BiLSTM层

窗口大小为3，4层扩张CNN
本文将block定义为上面的一个结构，将这些block进行堆叠意为迭代

GPU资源匮乏，CNN在并行上的优势[特征独立]
CNN固定窗口大小无法获取更多上下文信息，不断增加CNN虽然可以扩大窗口，但参数增多，分辨表现差
ID-CNN不会损失分辨率，限定参数

| | |
|---|---|
| Bi-LSTM-CRF (re-impl) | $90.43 \pm 0.12$ |
| ID-CNN-CRF | $\mathbf{90.54 \pm 0.18}$ |

| Model | Speed |
|---|---|
| Bi-LSTM-CRF | $1\times$ |
| Bi-LSTM | $9.92\times$ |
| ID-CNN-CRF | $1.28\times$ |
| 5-layer CNN | $12.38\times$ |
| ID-CNN | $14.10\times$ |

[Emma Strubell et al.2017] Fast and Accurate Entity Recognition with Iterated Dilated Convolutions

# Active Learning - CRF in clinical text[2][15]

- 主动学习分为以下几个部分：
  - 1.初始模型生成：使用少量带标签的样本来构建模型，样本的采样可以使用下面两种方法：一种是随机采样，一种是使用最长句子采样
  - 2.查询：在pool中没有标注的句子使用查询算法被排序，一些算法是使用CRF模型来进行排序，一些不是，选择前N个句子进行注释，然后放到注释过的集合中，每次迭代的batch size按照8，16，32，64的规律选择。
  - 3.训练：CRF模型在更新的注释集中重新训练
  - 4.迭代：重复2，3过程，直到达到某个标准停止。
- 查询算法：
  - 基于不确定的查询算法，最不确定的句是最有信息量的句子[高度依赖于模型的质量，比如通过CRF模型来排序句子] [论文中比较的结果是该类算法整体优于下面的算法]
  - 基于多样性的查询算法，通过单词，语义，句法来构建向量，通过各个句子之间的相似度来排序

[Yukun Chen et al.2015] A study of active learning methods for named entity recognition in clinical text

# Proactive Learning[14][17]

- ## Active learning
  - 假设标记的句子没有错误[乏味/困难导致出错在所难免]
  - 句子只拿给专家标注[成本高]

- ## Proactive learning
  - 对于困难的标注问题出现错误不可避免，允许专家出错
  - 使用两类注释器(者)，reliable expert 和 fallible expert[节约成本]

- ## Algorithm

  1.在未标记数据集中的所有句子使用active learning 准则排序，最有信息的N个句子作为批次采样的输入，这这一步中，批次的句子分布到两个集合分别给reliable 和fallible。分给fallible的句子是fallible有一个很高概率标记正确的句子，同时，只有那些fallible很难标注的句子才会被送给reliable。这样标记的成本就会减少。

  2. 第一步：设定阈值α，当fallible expert判别该句子的概率大于这个阈值时，将这个句子分配给fallible expert，否则进到第二步；

  第二步：计算句子对reliable和fallible的不同：

  $$diff(reliable, fallible, \boldsymbol{x})$$
  $$= |p(CorrectLabels|reliable, \boldsymbol{x})$$
  $$- p(CorrectLabels|fallible, \boldsymbol{x})| \quad (4)$$

  设定一个阈值β，当diff大于该阈值时，将句子交给reliable，否则交给fallible

**Algorithm 1: Proactive Learning for NER**
**Input:** a labelled dataset $L$, an unlabelled dataset $UL$, a test dataset $T$, a budget $B$, a reliable expert $e_r$ with cost $C_r$ for each sentence, a fallible expert $e_f$ with cost $C_f$, the current cost $C$
**Output:** a labelled dataset $L$
1 Estimate the performance of each expert as described in Section 2.1;
2 **while** $C < B$ **do**
3     Train a named entity recognition model $M$ on $L$;
4     Sort all sentences in the unlabelled dataset according to an active learning criterion;
5     Select the top $N$ sentences;
6     $UL_r, UL_f = BatchSampling(M, top N \text{ sentences});$
7     $L_r, L_f \leftarrow e_r$ and $e_f$ annotate $UL_r$ and $UL_f$ respectively;
8     $L = L \cup L_r \cup L_f;$
9     $UL = UL - UL_r - UL_f;$
10     $C = C + C_r * |L_r| + C_f * |L_f|;$
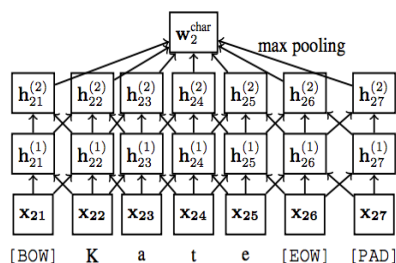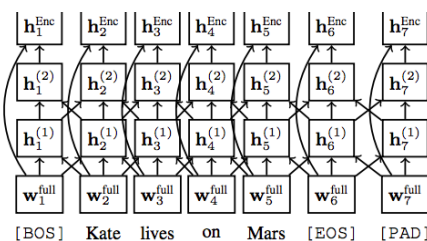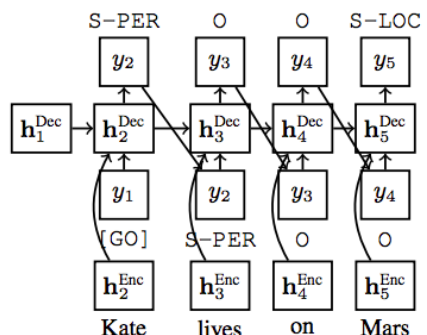11 **end**

**Algorithm 2: Batch Sampling**
**Input:** a named entity recognition model $M$, top-$N$ sentences selected according to an active learning criterion
**Output:** $UL_r, UL_f$
1 $UL_r = \emptyset;$
2 $UL_f = \emptyset;$
3 **while** *Batch Size* **do**
    // Stage 1
4     **foreach** *sentence* $x$ **do**
5        **if** $p(CorrectLabels|fallible, \boldsymbol{x}) > \alpha$ **then**
6           $UL_f = UL_f \cup \{\boldsymbol{x}\};$
7           BatchSize = BatchSize - 1
8        **end**
9     **end**
    // Stage 2
10     **if** *Batch Size* $\neq 0$ **then**
11        Sort the remaining sentences according to a re-ranking criterion;
12        Calculate threshold $\beta$;
13        **foreach** *sentence* $x$ **do**
14           **if** *Batch Size* $\neq 0$ **then**
15              **if** $diff(reliable, fallible, \boldsymbol{x}) < \beta$ **then**
16                 $UL_f = UL_f \cup \{\boldsymbol{x}\};$
17              **else**
18                 $UL_r = UL_r \cup \{\boldsymbol{x}\};$
19              **end**
20           BatchSize = BatchSize - 1;
21           **end**
22        **end**
23     **end**
24 **end**

[Maolin Li et al.2017] Proactive Learning for Named Entity Recognition

# Deep Active Learning-CNN-CNN-LSTM[16][18]



模型实现加速：使用了CNN-CNN-LSTM模型，这个模型在标准数据集上实现了近乎最先进的效果，同时在计算上比最佳性能的模型更有效[时间和表现综合来说模型最优]

模型在小数据集上表现：在训练过程中执行增量的主动学习，仅仅使用原有训练集的25%，模型就能达到原有的最好表现[节约成本]

嵌入层使用character-level和word-level嵌入级联：

$$\mathbf{w}_i^{full} := \left( \mathbf{w}_i^{char}, \mathbf{w}_i^{emb} \right).$$

character-level的嵌入和信息提取层都使用的是CNN层，两层一维卷积+一层最大池化，窗口大小为3

$$\mathbf{h}_i^{Enc} = \left( \mathbf{h}_i^{(l)}, \mathbf{w}_i^{full} \right)$$
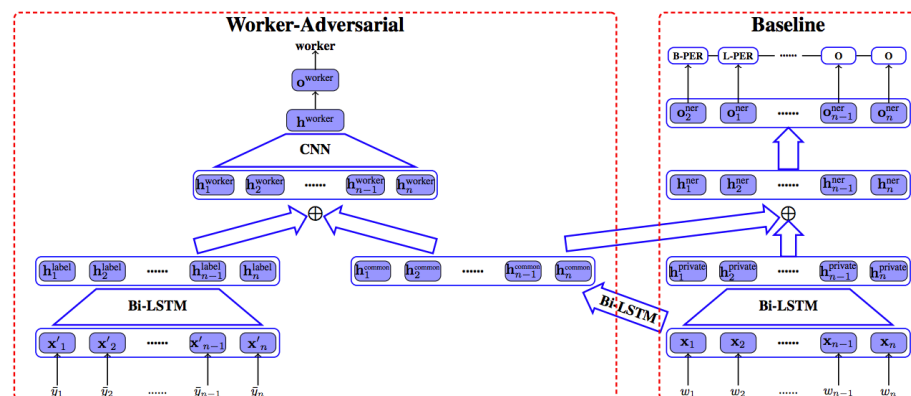
解码层使用CRF模块时间复杂度为O(nT²)，而使用LSTM时间复杂度为O(nT)，且能达到同样的效果。

| Char | Word | Tag | Reference | F1 | Sec/Epoch |
|------|------|-----|-----------|-----|-----------|
| None | CNN | CRF | Collobert et al. (2011) | 88.67 | - |
| None | LSTM | CRF | Huang et al. (2015) | 90.10 | - |
| LSTM | LSTM | CRF | Lample et al. (2016) | 90.94 | - |
| CNN | LSTM | CRF | Chiu & Nichols (2016) | 90.91 ± 0.20 | - |
| GRU | GRU | CRF | Yang et al. (2016) | 90.94 | - |
| None | Dilated CNN | CRF | Strubell et al. (2017) | 90.54 ± 0.18 | - |
| LSTM | LSTM | LSTM | | 90.89 ± 0.19 | 49 |
| CNN | LSTM | LSTM | | 90.58 ± 0.28 | 11 |
| CNN | CNN | LSTM | | 90.69 ± 0.19 | 11 |
| CNN | CNN | CRF | | 90.35 ± 0.24 | 12 |

Table 3: Evaluations on the test set of CoNLL-2003 English

[Yanyao Shen et al.2018] DEEP ACTIVE LEARNING FOR NAMED ENTITY RECOGNITION

# Adversarial Learning for Crowd[17][18]

- 从众包标记数据质量低角度出发，通过对抗学习提取公有特征，减轻注释噪声



对抗学习的目的：优化公有模块的学习质量，使之收敛于真实的数据

原有的结构不变(如右部分)，只是加入了对抗网络(左部分)，使用众包标签作为输入，经过Bi-LSTM，和common级联，一维卷积，窗口大小是5，最大池化，最后经过softmax，维度是标注者的个数：

$$\mathbf{h}_t^{worker} = \mathbf{h}_t^{common} \oplus \mathbf{h}_t^{label}$$

$$\tilde{\mathbf{h}}_t^{worker} = \tanh(\mathbf{W}^{cnn}[\mathbf{h}_{t-2}^{worker}, \mathbf{h}_{t-1}^{worker}, \cdots, \mathbf{h}_{t+2}^{worker}])$$

$$\mathbf{h}^{worker} = \text{max-pooling}(\tilde{\mathbf{h}}_1^{worker}\tilde{\mathbf{h}}_2^{worker}\cdots\tilde{\mathbf{h}}_n^{worker})$$

$$\mathbf{o}^{worker} = \mathbf{W}^{worker}\mathbf{h}^{worker},$$

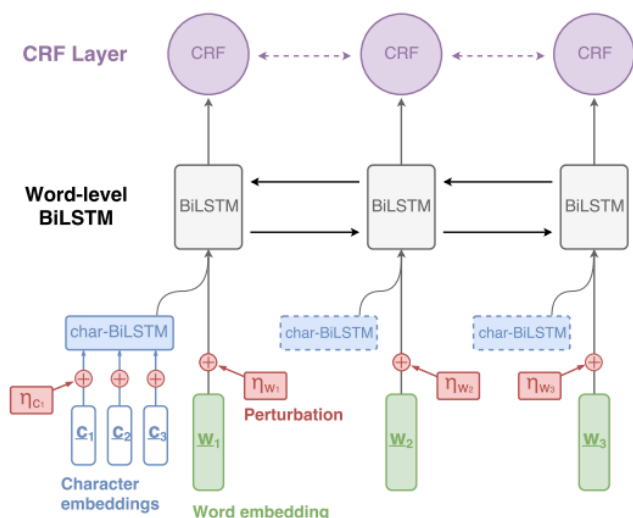$$p(\bar{z}|\mathbf{X}, \bar{\mathbf{y}}) = \frac{\exp(\mathbf{o}_{\bar{z}}^{worker})}{\sum_z \exp(\mathbf{o}_z^{worker})},$$

左右联合：

$$R(\Theta, \Theta', \mathbf{X}, \bar{\mathbf{y}}, \bar{z}) = \text{loss}(\Theta, \mathbf{X}, \bar{\mathbf{y}}) - \text{loss}(\Theta, \Theta', \mathbf{X})$$
$$= -\log p(\bar{\mathbf{y}}|\mathbf{X}) + \log p(\bar{z}|\mathbf{X}, \bar{\mathbf{y}}),$$

其中Θ表示整个模型中和NER有关的参数，Θ′表示仅仅和worker discriminator有关的参数，其中common Bi-LSTM对应的参数属于Θ中，优化下面式子：

$$\hat{\Theta} = \underset{\Theta}{\arg\min} R(\Theta, \Theta', \mathbf{X}, \bar{\mathbf{y}}, \bar{z})$$

$$\hat{\Theta}' = \underset{\Theta'}{\arg\max} R(\hat{\Theta}, \Theta', \mathbf{X}, \bar{\mathbf{y}}, \bar{z})$$

[Yaosheng Yang et al.2018] Adversarial Learning for Chinese NER from Crowd Annotations

# Adversarial Training for Robust[19][18]



在训练词嵌入和字符嵌入的时候加入噪声信息，对于加入噪声的输入，模型仍然需要正确标注出来，提高模型的鲁棒性

输入句子所有的单词/字符嵌入用s表示，模型的参数用θ表示，y为目标的词性标记序列，训练过程中，模型要最小化负对数似然函数：

$$L(\boldsymbol{\theta}; \boldsymbol{s}, \boldsymbol{y}) = -\log p(\boldsymbol{y} \mid \boldsymbol{s}; \boldsymbol{\theta})$$

在s上加一个连续的扰动向量，这个扰动向量为使得我们的模型的损失函数最大：

$$\boldsymbol{\eta} = \arg\max_{\boldsymbol{\eta}': \|\boldsymbol{\eta}'\|_2 \le \epsilon} L(\hat{\boldsymbol{\theta}}; \boldsymbol{s} + \boldsymbol{\eta}', \boldsymbol{y})$$

θ^为模型的参数（在上述求解过程中为常量）对抗训练的样本为：

$$\boldsymbol{s}_{\mathrm{adv}} = \boldsymbol{s} + \boldsymbol{\eta}$$

定义如下对抗训练的损失函数(本文γ = 0.5 )：

$$\tilde{L} = \gamma L(\boldsymbol{\theta}; \boldsymbol{s}, \boldsymbol{y}) + (1 - \gamma) L(\boldsymbol{\theta}; \boldsymbol{s}_{\mathrm{adv}}, \boldsymbol{y})$$

[Michihiro Yasunage et al.2018] Robust Multilingual Part-of-Speech Tagging via Adversarial Training

# Transfer Learning for sequence tagging[13][17]

- 迁移学习通过学习相关领域(source)的知识来提高当前领域(target)模型的性能



(a) Base model: both of Char NN and Word NN can be implemented as CNNs or RNNs.

(b) Transfer model T-A: used for cross-domain transfer where label mapping is possible.

(c) Transfer model T-B: used for cross-domain transfer with disparate label sets, and cross-application transfer.

(d) Transfer model T-C: used for cross-lingual transfer.

该论文中 word-level和character-level部分都使用的是GRU
其中CRF部分添加了max-margin principle，如下：

$$f(\mathbf{h}, \mathbf{y}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{h})} \exp(f(\mathbf{h}, \mathbf{y}') + \text{cost}(\mathbf{y}, \mathbf{y}')),$$

h 为从 word-level 和 character-level 部分学习到的特征表示，h = $(h_1, .., h_T)$，y 为标记序列，y = $(y_1, ..., y_T)$，其中Y(h)表示 h 对应的标记空间，cost 函数表示当y'和 y 的差距越大，惩罚越大。f 函数表示的是 CRF 中的打分函数，前面提到的s(X,y)。

Table 3: Comparison with state-of-the-art results (%).

| Model | CoNLL 2000 | CoNLL 2003 | Spanish | Dutch | PTB 2003 |
|---|---|---|---|---|---|
| Collobert et al. (2011) | 94.32 | 89.59 | – | – | 97.29 |
| Passos et al. (2014) | – | 90.90 | – | – | – |
| Luo et al. (2015) | – | 91.2 | – | – | – |
| Huang et al. (2015) | 94.46 | 90.10 | – | – | 97.55 |
| Gillick et al. (2015) | – | 86.50 | 82.95 | 82.84 | – |
| Ling et al. (2015) | – | – | – | – | **97.78** |
| Lample et al. (2016) | – | 90.94 | 85.75 | 81.74 | – |
| Ma & Hovy (2016) | – | 91.21 | – | – | 97.55 |
| Ours w/o transfer | 94.66 | 91.20 | 84.69 | 85.00 | 97.55 |
| Ours w/ transfer | **95.41** | **91.26** | **85.77** | **85.19** | 97.55 |

[Zhilin Yang et al.2017] TRANSFER LEARNING FOR SEQUENCE TAGGING WITH HIERARCHICAL RECURRENT NETWORKS
Label-aware Double Transfer Learning for Cross-Specialty Medical Named Entity Recognition   28  Apr  2018

# Incorporating dictionaries for NER[20][18]

- 数据驱动的方法典型缺乏处理稀有或没有出现的实体，本文在深层网络中加入词典

CCKS-2017 Task2：临床命名实体识别
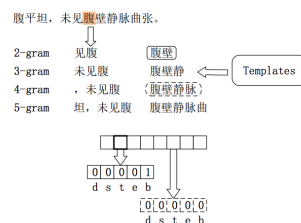本文使用character-level作为嵌入层，同时构造了一个特征向量di
特征构造分三种：
1.N-gram feature
总共有8个temple，每个temple映射到5维上，d，s，t，e，b，分别表示disease，
symptom，treatment，exam，body-part，即每个字符在这个特征类别上有40维度。

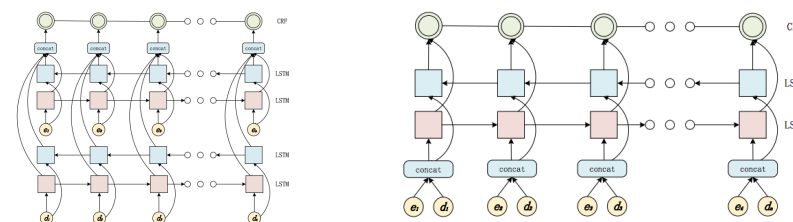| Type | template |
|------|----------|
| 2-gram | $x_{i-1}x_i$, $x_ix_{i+1}$ |
| 3-gram | $x_{i-2}x_{i-1}x_i$, $x_ix_{i+1}x_{i+2}$ |
| 4-gram | $x_{i-3}x_{i-2}x_{i-1}x_i$, $x_ix_{i+1}x_{i+2}x_{i+3}$ |
| 5-gram | $x_{i-4}x_{i-3}x_{i-2}x_{i-1}x_i$, $x_ix_{i+1}x_{i+2}x_{i+3}x_{i+4}$ |

2.Position-Independent Entity Type feature
对句子X基于词典D进行前向最大匹配，并进行归类，然后将每个字符进行映射，
这里其特征表示可以使用one-hot或者特征嵌入矩阵，如下第三行所示
3. Position-Dependent Entity Type feature
这个是有2中切词后的结果来获取的，即该字符是否是实体的开始，中间，结束，
或者是单个字符实体，分别用B，I，E，S，同时可以用one-hot或者特征嵌入矩阵

| Character sequence | 腹 | 平 | 坦 | ， | 未 | 见 | 腹 | 壁 | 静 | 脉 | 曲 | 张 | 。 |
|--------------------|----|----|----|----|----|----|----|----|----|----|----|----|----|
| tag sequence | S-b | O | O | O | O | O | B-b | E-b | B-s | I-s | I-s | E-s | O |
| PIET features | b | None | None | None | None | None | b | b | s | s | s | s | None |
| PDET features | S-b | None | None | None | None | None | B-b | E-b | B-s | I-s | I-s | E-s | None |

这些特征都受词典或上下文的影响，而不受其他句子或统计信息的影响，因此是有别于流行的数据驱动的方法。

提出两种模型如上图所示，3种特征的5种表现形式和2种模型进行组合，
结果如下：

| | | Model-I | | | Model-II | | |
|---|---|---|---|---|---|---|---|
| | | Precision | Recall | $F_1$-Measure | Precision | Recall | $F_1$-Measure |
| N-gram feature | | 88.39 | 88.46 | 88.43 | 88.72 | 88.71 | 88.71 |
| PIET feature | one-hot encoding | 89.53 | 90.58 | 90.05 | 89.38 | 90.49 | 89.93 |
| | feature embedding | 90.11 | 90.01 | 90.56 | 90.00 | 90.60 | 90.30 |
| PDET feature | one-hot encoding | 90.51 | 91.04 | 90.77 | 90.22 | 90.64 | 90.43 |
| | feature embedding | **90.83** | **91.64** | **91.24** | 90.36 | 91.35 | 90.85 |

[Qi Wang et al.2018] Incorporating dictionaries into deep neural networks for the Chinese clinical named entity recognition

# Conclusion

- 经典的三层模型是不会变的
  - Emb representation–info extract–CRF transfor
- 通过各种手段来提高模型的表现
  - 对抗生成网络
  - Attention
  - 加词表，加特征
- 少量标签的学习/加速
  - 迁移学习
  - 主动学习
  - CNN替换LSTM，LSTM替换CRF

# paper：

- [1].[Shaodian Zhang et al.2013] Unsupervised biomedical named entity recognition: Experiments with clinical and biological texts
- [2].[Yukun Chen et al.2015] A study of active learning methods for named entity recognition in clinical text
- [3]. [Zhiheng Huang et al.2015] Bidirectional LSTM-CRF Models for Sequence Tagging
- [4].[Chris Dyer et al.2015] Transition-Based Dependency Parsing with Stack Long Short-Term Memory
- [5].[Guillaume Lample et al.2016] Neural Architectures for Named Entity Recognition
- [6].[Xuezhe Ma et al.2016] End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF
- [7].[Marek Rei et al.2016] Attending to Characters in Neural Sequence Labeling Models
- [8].[Nanyun Peng et al.2016] Improving Named Entity Recognition for Chinese Social Media with Word Segmentation Representation Learning
- [9].[Simon Almgren et al.2016] Named Entity Recognition in Swedish Health Records with Character-Based Deep Bidirectional LSTMs
- [10].[Chuanhai Dong et al.2016] Character-Based LSTM-CRF with Radical-Level Features for Chinese Named Entity Recognition
- [11].[Peters et al.2017] Semi-supervised sequence tagging with bidirectional language models
- [12].[Emma Strubell et al.2017] Fast and Accurate Entity Recognition with Iterated Dilated Convolutions
- [13].[Zhilin Yang et al.2017] TRANSFER LEARNING FOR SEQUENCE TAGGING WITH HIERARCHICAL RECURRENT NETWORKS
- [14].[Maolin Li et al.2017] Proactive Learning for Named Entity Recognition
- [15].[Bill Y.Lin et al.2017] Multi-channel BiLSTM-CRF Model for Emerging Named Entity Recognition in Social Media
- [16].[Yanyao Shen et al.2018] DEEP ACTIVE LEARNING FOR NAMED ENTITY RECOGNITION
- [17].[Yaosheng Yang et al.2018] Adversarial Learning for Chinese NER from Crowd Annotations
- [18].[Zhenghui Wan et al.2018] Label-aware Double Transfer Learning for Cross-Specialty Medical Named Entity Recognition
- [19].[Michihiro Yasunage et al.2018] Robust Multilingual Part-of-Speech Tagging via Adversarial Training
- [20].[Qi Wang et al.2018] Incorporating dictionaries into deep neural networks for the Chinese clinical named entity recognition

NLP