

论文写作的易读性原则

案例分析：基于Seq2Seq的对话数据增广

报告人：刘一佳

合作者：侯宇泰、车万翔、刘挺

<http://yjliu.net/cv/res/2018-08-19-nlpcc-sws.compressed.pdf>

学术报告中的一些设计技巧

报告人：刘一佳
导师：秦兵、车万翔

错误地利用 报告与论文结构的相似性

Challenges and Contribution

- The first challenge is deriving an optimal alignment in ambiguous documents.
- The second challenge is to automatically align the document structure by analyzing the node labels and content of the document.
- The third challenge is to align the document structure with different types of learning.
- We propose an enhanced aligner based on transition-based neural networks.

简介

Overview



Our aligner algorithm

- Enhancing aligner with neural networks
- Producing better alignments

模型

Our oracle parser

模型

Experiments

- We conduct experiments on DocAlign2012.
- We evaluate the alignment score and length of generated alignments.

实验

Conclusion

- We present a novel aligner which is based on a neural network and can automatically generate an alignment between document structures. Our aligner is able to produce better alignments than previous work.
- Such an aligner can be used in other applications such as document summarization, document clustering, etc.
- We also present a neural parser based on our aligner and transition-based and it achieves a performance of 95.4 F-measure F1 score on sentences with entity words and PCFG tags as input.

结论

思考题

- 为什么做学术报告
 - 为了更好地交流
- 做怎样的学术报告
 - “向听众展示我对问题的深入理解”
 - “让听众明白我的论文中的技术”
 - “引起听众的兴趣”

思考题

- 为什么做学术报告
 - 为了更好地交流
- 做怎样的学术报告
 - “向听众展示我对问题的深入理解”
 - “让听众明白我的论文中的技术”
 - “引起听众的兴趣”

听众模型

理想中的听众

- 领域专家
- 已经读过你的论文
- 对于你的工作非常感兴趣

现实中的听众

- 来自其他领域
- 刚刚了解到你的工作
- 这个时段没什么可听的，恰巧发现这屋子网络比较好

类比审稿人模型

审稿

你以为审稿人应该是这样审稿的：

审稿人一定是专家，无所不知。打印出来，仔细研读揣摩数天，对于看不懂的地方反复推敲。即使你的英文写得极其糟糕、即使你的文章组织很混乱、即使你的表述很难看懂，审稿人花费了大量的时间后终于看懂了，他认为你的工作是有意义的，决定给你个border line或以上的分数。

审稿人实际上往往是这样审稿的：

他不一定是专家，一直忙于其他事，在deadline到来之前一天要完成n篇。审稿时他往往先看题目、摘要，扫一下introduction（知道你做什么），然后直接翻到最后找核心实验结果（做得好不好），然后基本确定录还是不录（也许只用5分钟！）。如果决定录，剩下就是写些赞美的话，指出些次要的小毛病。如果决定拒，下面的过程就是细看中间部分找理由拒了。

第一印象定录拒，5分钟内打动审稿人

类比审稿人模型

审稿

你以为审稿人应该是这样审稿的：

审稿人一定是专家，无所不知。打印出来，仔细研读揣摩数天，对于看不懂的地方反复推敲。即使你的英文写得极其糟糕、即使你的文章组织很混乱、即使你的表述很难看懂，审稿人花费了大量的时间后终于看懂了，他认为你的工作是有意义的，决定给你个border line或以上的分数。

审稿人实际上往往是这样审稿的：

他不一定是专家，一直忙于其他事，在deadline到来之前一天要完成n篇。审稿时他往往先看题目、摘要，扫一下introduction（知道你做

“You have two minutes to engage your audience before they start to doze.” -- Simon Peyton Jones in *How to give a great research talk*

简介部分：展示最好的部分

(Zhang and Nirve 2011, Martins et al 2013)



Our Work

- A neural network based dependency parser!

Parsing on English Penn Treebank (§23):

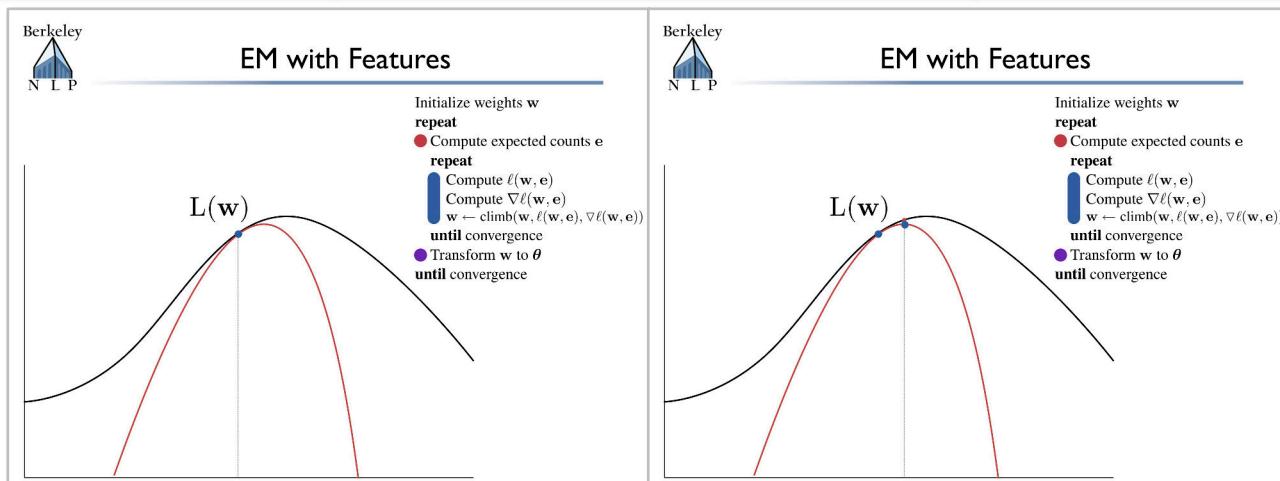
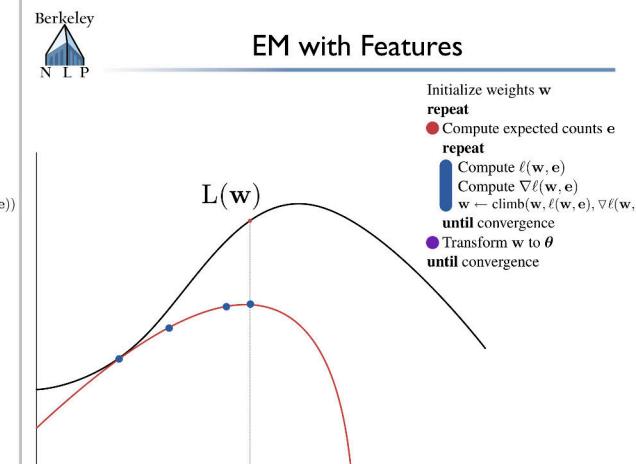
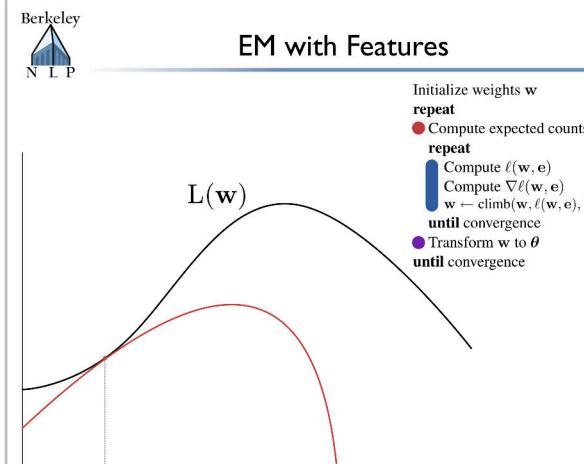
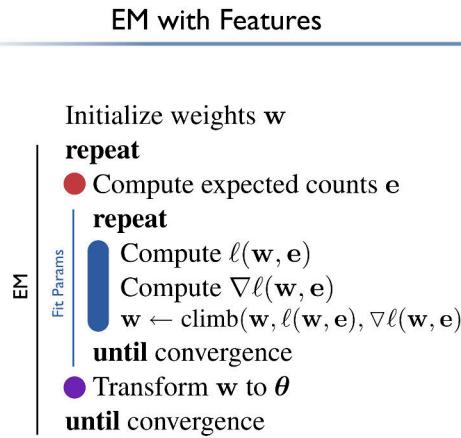
	Unlabeled attachment score (UAS)	sent / s
Transition -based	MaltParser (greedy)	89.9
	Our Parser (greedy)	92.0
Zpar: beam = 64	92.9*	29*
Graph -based	MSTParser	92.0
	TurboParser	93.1*

A Fast and Accurate Dependency Parser using Neural Networks

3

Danqi Chen and Christopher Manning. 2014. A Fast and Accurate Dependency Parser using Neural Networks, 第三页

模型部分：多用例子



模型部分：反例

Transition	Current State	Resulting State	Description
DROP	$[\sigma s_0, \delta, b_0 \beta, A]$	$[\sigma s_0, \delta, \beta, A]$	pops out the word that doesn't convey any semantics (e.g., function words and punctuations).
MERGE	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{b}_1 \bar{\beta}, \bar{A}]$	$[\sigma \bar{s}_0, \bar{\delta}, \bar{b}_0\bar{b}_1 \bar{\beta}, \bar{A}]$	concatenates a sequence of words into a span, which can be derived as a named entity (name) or date-entity.
CONFIRM(\bar{c})	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma \bar{s}_0, \bar{\delta}, \bar{c} \bar{\beta}, \bar{A}]$	derives the first element of the buffer (a word or span) into a concept c .
ENTITY(\bar{c})	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma \bar{s}_0, \bar{\delta}, \bar{c} \bar{\beta}, \bar{A} \cup \text{relations}(c)]$	a special form of CONFIRM that derives the first element into an entity and builds the internal entity AMR fragment.
NEW(\bar{c})	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma \bar{s}_0, \bar{\delta}, \bar{c} \bar{b}_0 \bar{\beta}, \bar{A}]$	generates a new concept c and pushes it to the front of the buffer.
LEFT(r)	$[\sigma s_0, \delta, b_0 \beta, A]$	$[\sigma s_0, \delta, b_0 \beta, A \cup \{s_0 \xleftarrow{r} b_0\}]$	links a relation r between the top concepts on the stack and the buffer.
RIGHT(r)	$[\sigma s_0, \delta, b_0 \beta, A]$	$[\sigma s_0, \delta, b_0 \beta, A \cup \{s_0 \xrightarrow{r} b_0\}]$	
CACHE	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma, s_0 \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	passes the top concept of the stack onto the deque.
SHIFT	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma \bar{s}_0 \bar{\delta} b_0, [\bar{\beta}, \bar{A}]$	shifts the first concept of the buffer onto the stack along with those on the deque.
REDUCE	$[\bar{\sigma} \bar{s}_0, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	$[\sigma, \bar{\delta}, \bar{b}_0 \bar{\beta}, \bar{A}]$	pops the top concept of the stack.

实验部分：图比表格好

LDC2014T12 Experiments

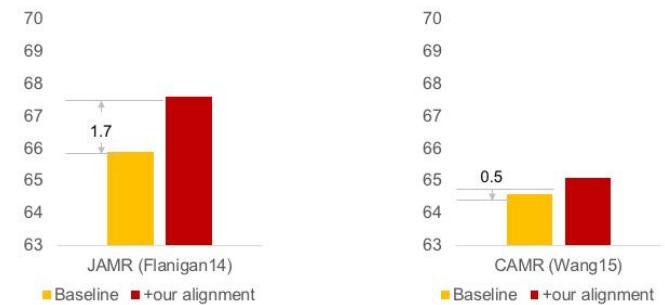
- alignment F-score

Aligner	Alignment F1 (on hand-align)	Oracle's Smatch (on dev. dataset)
JAMR	90.6	91.7
Our	95.2	94.7

- parser improvements

model	newswire	all
JAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	71.3	65.9
+ Our aligner	73.1	67.6
CAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	68.4	64.6
+ Our aligner	68.8	65.1

Aligner Experiments: Two Open-sourced AMR Parsers



实验部分：图比表格好

信息元素的易理解度

step	action	rule	stack	coverage
0				ooooooo
1	S	r_3	[The President will]	oooooo.
2	S	r_1	[The President will] [visit]	oooooo*
3	R_t		[The President will visit]	oooooo*
4	S	r_4	[The President will visit] [London in April]	oooooo*
5	R_e		[The President will visit London in April]	oooooo*

图

*

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \theta_k} &= \sum_{i=1}^I \sum_{y \in \mathcal{Y}(x^{(i)})} P(y|x^{(i)}, \theta) \phi_k(x^{(i)}, y) \\ &\quad - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}(x)} P(x, y|\theta) \phi_k(x, y) \\ &= \sum_{i=1}^I \mathbb{E}_{y|x^{(i)}, \theta} [\phi_k(x^{(i)}, y)] - \mathbb{E}_{x, y|\theta} [\phi_k(x, y)] \end{aligned}$$

公式

System	Setting	English-French	Chinese-English
GIZA++	Model 4+2t	7.7	20.9
	Model 4+2s	9.2	30.3
	Intersection	6.8	21.8
	Union	9.6	28.1
Vigene	Refined method	5.9	18.4
	Cross-EM	5.1	18.9
	HMM, joint	5.6	18.3
	+Model 4+2s	5.5	17.7
Vigene	+cross	5.4	17.6
	+neighbor count	5.2	17.4
	+exact match	5.3	-
	+linked word count	5.2	17.3
	+linked word dictionary	5.1	17.1
	+link co-occurrence count (GIZA++)	5.1	16.3
	+link co-occurrence count (Cross-EM)	4.0	15.7

表格

**

```
Algorithm 1 A beam search algorithm for word alignment
1: procedure ALIGN(f, e)
2:   open ← ∅                                ▷ a list of active alignments
3:   N ← ∅                                     ▷ n-best list
4:   a ← ∅                                     ▷ begin with an empty alignment
5:   ADD(open, a, β, b)                         ▷ initialize the list
6:   while open ≠ ∅ do
7:     closed ← ∅                               ▷ a list of promising alignments
8:     for all a ∈ open do
9:       for all i ∈ s → a do
10:         a' ← a ∪ {i}                           ▷ enumerate all possible new links
11:         g ← GAIN(f, e, a')                   ▷ produce a new alignment
12:         if g > 0 then
13:           ADD(closed, a', β, b)               ▷ compute the link gain
14:           ensure that the score will increase
15:           ADD(N, a', 0, n)                  ▷ update promising alignments
16:         end if
17:       end for
18:     end for
19:     open ← closed                            ▷ update active alignments
20:   end while
21:   return N                                    ▷ return n-best list
22: end procedure
```

算法

Shift-reduce parsing is efficient but suffers from parsing errors caused by syntactic ambiguity. Figure 3 shows two (partial) derivations for a dependency tree. Consider the item on the top, the algorithm can either apply a shift action to move a new item or apply a reduce left action to obtain a bigger structure. This is often referred to as *conflict* in the shift-reduce dependency parsing literature (Huang et al., 2009). In this work, the shift-reduce parser faces four types of conflicts:

正文

Proof of Theorem 1: Let $\bar{\alpha}^k$ be the weights before the k 'th mistake is made. It follows that $\bar{\alpha}^1 = 0$. Suppose the k 'th mistake is made at the i 'th example. Take z to the output proposed at this example, $z = \arg\max_{y \in \text{GEN}(x_i)} \Phi(x_i, y)$. $\bar{\alpha}^k$ follows from the algorithm updates that $\bar{\alpha}^{k+1} = \bar{\alpha}^k + \Phi(x_i, y) - \Phi(x_i, z)$. We take inner products of both sides with the vector U :

$$\begin{aligned} U \cdot \bar{\alpha}^{k+1} &= U \cdot \bar{\alpha}^k + U \cdot \Phi(x_i, y) - U \cdot \Phi(x_i, z) \\ &\geq U \cdot \bar{\alpha}^k + \delta \end{aligned}$$

where the inequality follows because of the property of U assumed in Eq. 3. Because $\bar{\alpha}^1 = 0$, and therefore $U \cdot \bar{\alpha}^1 = 0$, it follows by induction on k that for all k , $U \cdot \bar{\alpha}^{k+1} \geq k\delta$. Because $U \cdot \bar{\alpha}^{k+1} \leq \|U\| \|\bar{\alpha}^{k+1}\|$, it follows that $\|\bar{\alpha}^{k+1}\| \geq k\delta$.

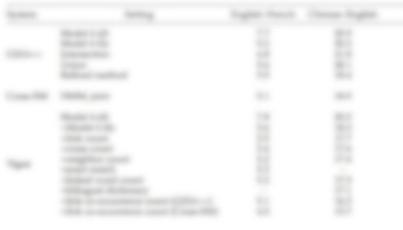
证明

实验部分：图比表格好

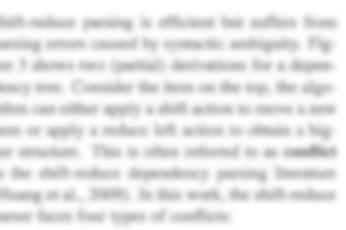
信息元素的易理解度



图



表格



正文

Shift-reduce parsing is efficient but suffers from parsing errors caused by syntactic ambiguity. Figure 3 shows two (partial) derivations for a dependency tree. Consider the item on the top, the algorithm can either apply a shift action to move a new item or apply a reduce left action to obtain a big gap structure. This is often referred to as conflict in the shift-reduce dependency parsing literature (Bisong et al., 2009). In this work, the shift-reduce parser faces four types of conflicts.

Proof of Theorem 3: Let α^2 be the sentence before the i^{th} word is read. It follows that α^2 is a dependency tree. We make a move at i^{th} word. Take β to be the current grammar in this example. $\beta = \text{reduce_bottom}(\mathcal{W}_1, \alpha^2)$. It follows from the expression system that $\beta^{i+1} = \beta^i \cdot \mathcal{W}_{i+1} = \mathcal{W}_{i+1} \cdot \beta^i$. We rightmost generate of both sides with the inverse \mathcal{W} .

用图与例子来描述方法和实验

结论部分：新的展现形式

Conclusion

Problem

The diagram shows a sequence of tokens: "Shareholders took their money". Above the tokens, a purple box contains the text "Shareholders took their money" and "arg max". Below the tokens, another purple box contains the text "Shareholders took their money" and "Downstream task". A blue arrow labeled $\nabla_{\theta} \mathcal{L}$ points from the bottom box to the top box. A red arrow labeled "Loss \mathcal{L} " points from the bottom box to the bottom right. A blue arrow labeled "Intermediate parser θ " points from the top box to the bottom box. A question mark "A layer?" is placed between the two boxes.

Method SPIGOT

A graph with four nodes. Three nodes are red and one is orange. Blue arrows connect the red nodes in a cycle. A black arrow labeled $-\nabla_s \mathcal{L}$ points from the bottom red node to the orange node. A dashed black arrow labeled $\hat{\mathbf{z}} - \nabla_{\hat{\mathbf{z}}} \mathcal{L}$ points from the orange node back to the bottom red node.

Results

The chart compares five methods: Neuro, Pipeline, STE, Structured Att., and SPIGOT. For each method, there are two bars: a purple bar for "in-domain" and a yellow bar for "out-of-domain". The y-axis represents the F1 score, ranging from 78 to 88. SPIGOT shows the highest performance across both domains.

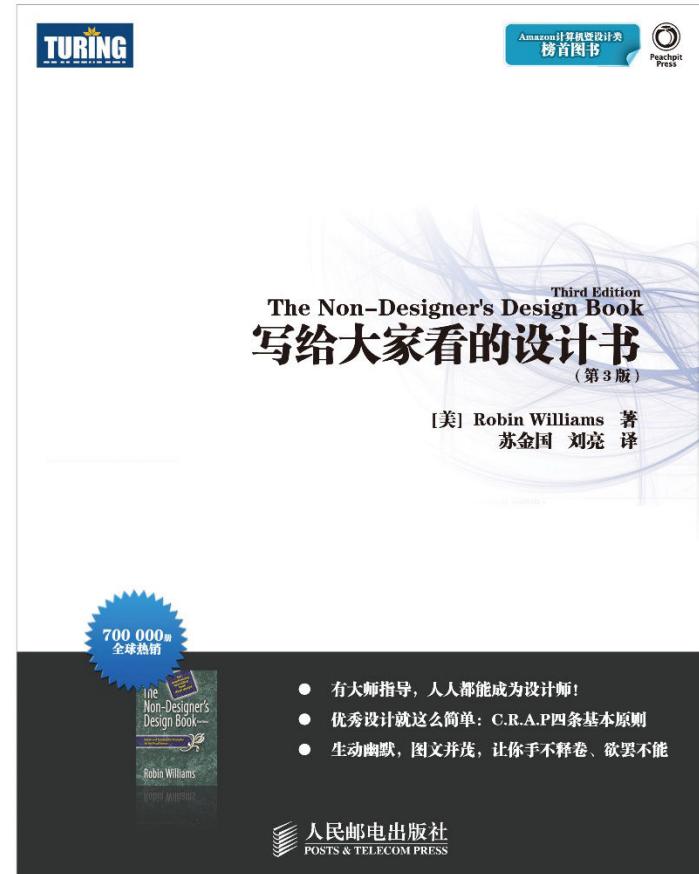
Method	in-domain	out-of-domain
Neuro	~85	~81
Pipeline	~86	~82
STE	~86	~83
Structured Att.	~85	~82
SPIGOT	~87	~83

Below the chart is a table comparing various backpropagation techniques:

	Syntax	Backprop	Well-formedness	Projection
Neuro	N/A	✓	✓	✓
Pipeline	N/A	✓	✓	✓
STE	N/A	✓	✓	✓
Structured Att.	N/A	✓	✓	✓
SPIGOT	N/A	✓	✓	✓

设计原则

- **亲密性**：相关的元素应该组织到一起
- **重复**：相同的内容达到形式的统一
- **对比**：如果两项不完全相同，就应使之截然不同
- **对齐**：使元素之间产生关联，有关联的都应对齐



根据设计原则做幻灯片

Challenges and Contribution

- The first challenge is deriving an optimal alignment in ambiguous situations.
- The second challenge is recalling more semantically matched word-concept pair without harming the alignment precision.
- The final challenge which is faced by both the rule-based and unsupervised aligners is tuning the alignment with downstream parser learning.
- We proposed an enhanced aligner tuned by transition-based oracle parser

加入空行提高相关
元素的亲密性

Challenges and Contribution

- The first challenge is deriving an optimal alignment in ambiguous situations.
- The second challenge is recalling more semantically matched word-concept pair without harming the alignment precision.
- The final challenge which is faced by both the rule-based and unsupervised aligners is tuning the alignment with downstream parser learning.
- We proposed an enhanced aligner tuned by transition-based oracle parser

Challenges and Contribution

- Challenges
 - deriving an optimal alignment in ambiguous situations.
 - recalling more semantically matched word-concept pair without harming the alignment precision.
 - tuning the alignment with downstream parser learning.
- Contribution
 - an enhanced aligner tuned by transition-based oracle parser

相同内容使用相同样式
即提高了一致性又形成
了必要的对比

避免不对齐

Our aligner algorithm

- Enhancing aligner with rich semantic resources
- Producing multiple alignments

```
Input: An AMR graph with a set of graph fragments  $C$ ;  
a sentence  $W$ ; a set of matching rules  $P_M$ ; and  
a set of updating rules  $P_U$ .  
Output: a set of alignments  $\mathcal{A}$ .  
1 for  $c \in C$  do  
2    $A_c \leftarrow \emptyset$ ;  
3 for  $\rho_M \in P_M$  do  
4   for  $w_{s,e} \leftarrow spans(W)$  do  
5     for  $c \in C$  do  
6       if  $\rho_M(c, w_{s,e})$  then  
7          $A_c \leftarrow A_c \cup (s, e, \text{nil})$ ;  
8 updated  $\leftarrow$  true ;  
9 while updated is true do  
10   updated  $\leftarrow$  false;  
11   for  $\rho_U \in P_U$  do  
12     for  $c, c' \in C \times C$  do  
13       for  $(s, e, d) \in A'_c$  do  
14         if  $\rho_U(c, w_{s,e}) \wedge (s, e, c') \notin A_c$  then  
15            $A_c \leftarrow A_c \cup (s, e, c')$ ;  
16           updated  $\leftarrow$  true;  
17  $\mathcal{A} \leftarrow \emptyset$ ;  
18 for  $(a_1, \dots, a_c) \in \text{CartesianProduct}(A_1, \dots, A_{|C|})$  do  
19   legal  $\leftarrow$  true;  
20   for  $a \in (a_1, \dots, a_c)$  do  
21      $(s, e, c') \leftarrow a$ ;  
22      $(s', e', d) \leftarrow a_{c'}$ ;  
23     if  $s \neq s' \wedge e \neq e'$  then  
24       legal  $\leftarrow$  false ;  
25   if legal then  
26      $\mathcal{A} \leftarrow \mathcal{A} \cup (a_1, \dots, a_c)$ ;
```

“乱”的原因：视线跳动过多

Experiments

- We conduct experiments on LDC2014T12
- We evaluate the alignment F-score and Smatch of resulted parsers

Aligner	Alignment F1 (on hand-align)	Oracle's Smatch (on dev. dataset)
JAMR	90.6	91.7
Our	95.2	94.7

model	newswire	all
JAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	71.3	65.9
+ Our aligner	73.1	67.6
CAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	68.4	64.6
+ Our aligner	68.8	65.1

model	newswire	all
Our single parser: Word only		
+ JAMR aligner	68.6	63.9
+ Our aligner	69.3	64.7
Our single parser: Word, POS		
+ JAMR aligner	68.8	64.6
+ Our aligner	69.8	65.2
Our ensemble: Word only + Our aligner		
x3	71.9	67.4
x10	72.5	68.1
Our ensemble: Word, POS + Our aligner		
x3	72.5	67.7
x10	73.3	68.4

“乱”的原因：视线跳动过多

Experiments

- We conduct experiments on LDC2014T12
- We evaluate the alignment F-score and Smatch of resulted parsers

Aligner	Alignment F1 (on hand-align)	Oracle's Smatch (on dev. dataset)
JAMR	90.6	91.7
Our	95.2	94.7

model	newswire	all
JAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	71.3	65.9
+ Our aligner	73.1	67.6
CAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	68.4	64.6
+ Our aligner	68.8	65.1

model	newswire	all
Our single parser: Word only		
+ JAMR aligner	68.6	63.9
+ Our aligner	69.3	64.7

model	newswire	all
Our single parser: Word, POS		
+ JAMR aligner	68.8	64.6
+ Our aligner	69.8	65.2

model	newswire	all
Our ensemble: Word only + Our aligner		
x3	71.9	67.4
x10	72.5	68.1
Our ensemble: Word, POS + Our aligner		
x3	72.5	67.7
x10	73.3	68.4

“乱”的解法：重新组织内容

Experiments

- We conduct experiments on LDC2014T12
- We evaluate the alignment F-score and Smatch of resulted parsers

Aligner	Alignment F1	Oracle's Smatch (on hand-align)	Oracle's Smatch (on dev. dataset)
JAMR	90.6	91.7	
Our	95.2	94.7	

model	newswire	all
Our single parser: Word only		
+ JAMR aligner	68.6	63.9
+ Our aligner	69.3	64.7
Our single parser: Word, POS		
+ JAMR aligner	68.8	64.6
+ Our aligner	69.8	65.2
Our ensemble: Word only + Our aligner		
x3	71.9	67.4
x10	72.5	68.1
Our ensemble: Word, POS + Our aligner		
x3	72.5	67.7
x10	73.3	68.4



LDC2014T12 Experiments

- alignment F-score
- parser improvements

model	newswire	all
JAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	71.3	65.9
+ Our aligner	73.1	67.6
CAMR parser: Word, POS, NER, DEP		
+ JAMR aligner	68.4	64.6
+ Our aligner	68.8	65.1

视线跳动在论文写作中的作用

信息流的变化

Tree-to-String Alignment Template for Statistical Machine Translation

Yi Liu, Qun Lin, and Shouxun Lin
Institute of Computing Technology
Chinese Academy of Sciences
No 6 Kexueyuan South Road, Haidian District
P.O. Box 2704, Beijing, 100080, China
{yliu,linqun,sxlin}@ict.ac.cn

Abstract

We present a novel template model based on *tree-to-string alignment template* (TAT) which describes the relation between a source tree and a target string. A TAT is composed of generating both terminals and non-terminals and performing template matching at both low and high levels. The models are template-based because TATs are extracted automatically from word-aligned, source side parsed parallel texts. To translate a source sentence, we first extract TATs to produce a source parse tree and then apply TATs to transform the tree into a target string. Our experiments show that the TAT model performs significantly outperforms Pharaoh, a state-of-the-art decoder for phrase-based SMT.

1 Introduction

Phrase-based translation models (Marcus and Wong, 2002; Koehn et al., 2003; Och and Ney, 2004), which lie beyond the traditional statistical models (Brown et al., 1990), by modeling translations of phrases rather than individual words, have been suggested to be the state-of-the-art in statistical machine translation by empirical evaluations.

In phrase-based models, phrases are usually strings of adjacent words instead of syntactic constituents, excelling at capturing local reordering and performing translations that are localized to

strings that are common enough to be observed on training data. However, a key limitation of phrase-based models is that they only consider relations at the phrase level robustly. Typically, phrase reordering is modeled in terms of offset positions at the word level (Koehn, 2004; Och and Ney, 2004), resulting little or no direct use of syntactic information.

Recent research on statistical machine translation has led to the development of syntax-based models (Xia et al., 1997) proposed by Translation Grammar. Tree-to-string translation as a process of parallel parsing of the source and target language via a synchronized grammar. Alshabani et al. (2000) propose tree production in a dependency tree as a parser transducer. Melamed (2004) formalizes machine translation problem as synchronous parsing problem on multi-text grammar. Yamada and Knight (2004) describe tree and decoding algorithm for both generalized tree-to-tree and tree-to-string translators. Chiang (2005) presents a hierarchical tree-to-tree model that has two parallel main parts, which are freely produced by a synchronous context-free grammar. Ding and Palmer (2005) propose a syntax-based translation model based on partial parallel synchronous dependency insert grammar, a version of synchronous grammars defined on dependency trees. All these approaches, though different in formalism, make use of a synchronous grammar or tree-based translation rules to model both source and target languages.

Another class of approaches make use of syntactic tree representations. In this approach, treating the translation problem as a parsing problem. Yamada and Knight (2001) use a parser in the target language to train probabilities on a set of the source string.

That is, given a source string $s = s_1, \dots, s_n$, we want to find a target string $t = t_1, \dots, t_m$ to be translated into a target string $t = t_1, \dots, t_m$. This is a local reordering problem of the target string, and β is the length of the source string.

609
Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL, pages 609–616, Sydney, July 2006. ©2006 Association for Computational Linguistics

Joint Tokenization and Translation

Xinyan Xiao[†], Yang Liu[†], Young-Soo Hwang[‡], Qun Lin[†], Shouxun Lin[†]
[†]Key Lab. of Intelligent Info. Processing,
Institute of Computing Technology
Chinese Academy of Sciences
{xianxinyan,yliu,liuqun,sxlin}@ict.ac.cn
[‡]H3Lab Convergence Technology Center
C4I Business
SK Telecom
yhwang@sktelecom.com

Abstract

As tokenization is usually ambiguous for many natural languages such as Chinese and Korean, tokenization errors might potentially introduce translation mistakes for translation systems that rely on word tokenizations. While using simple to offer more alternatives to translation systems have elegantly alleviated this problem, we take a further step to tokenize and translate jointly. That is, the sequence of atoms units that can be combined to form words are taken as input, our joint decoder produces a tokenization on the source side and a translation on the target side simultaneously. By integrating tokenization and translation features in a single pipeline, our joint decoder outperforms the baseline translation systems using beam search to obtain the best tokenizations and latencies significantly on both Chinese-English and Korean-English pairs. Interestingly, as a tokenizer, our joint decoder achieves significant improvements over monolingual Chinese parts.

1 Introduction

Tokenization plays an important role in statistical machine translation (SMT). Given a source-target language pair, always the first step in SMT systems. Based on the type of input, Mi and Huang (2008) distinguish between two categories of SMT systems : *string-based* systems (Koehn et al., 2003; Chiang, 2007; Galley et al., 2006; Shan et al., 2008) that take a string as input and *tree-based* systems (Liu et al., 2006; Mi et al., 2008) that take a tree as input. Note that a tree-based system still needs to first tokenize the source sentence and then obtain the tokens one of the sentence. As shown in Figure 1(a), we refer to this pipeline as *separate tokenization and translation* because they are divided into single steps.

A major problem for most languages is usually ambiguous, SMT systems for separate tokenization and translation suffer from a major drawback: tokenization errors might potentially introduce translation mistakes. As some languages such as Chinese have no spaces in the writing system, how to segment sentences into appropriate words has a dramatic impact on the system performance (Xu et al., 2003; Chiang et al., 2008; Zhang et al., 2009). In addition, although agglutinative languages such as Korean incorporate spaces between "words", which consist of multiple morphemes, the granularity is too coarse and makes the training data

1200
Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010), pages 1200–1208, Beijing, August 2010.

参考文献

- Simon Peyton Jones: How to give a great talk
- 写给大家看的设计书
- 机器翻译学术论文写作方法与技巧
- 知乎专栏：跟我学个P

总结

(Zhang and Nirve 2011, Martins et al 2013)

 Our Work

- A neural network based dependency parser!

Parsing on English Penn Treebank (§23):

	Unlabeled attachment score (UAS)	sent / s
Transition-based	MaltParser (greedy) 89.9 Our Parser (greedy) 92.0	+2.1 $\times 1.8$ 560 1013
	Zpar: beam = 64 92.9*	29*
Graph-based	MSTParser 92.0 TurboParser 93.1*	12 31*

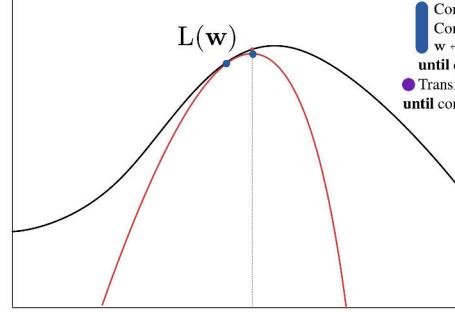
A Fast and Accurate Dependency Parser using Neural Networks 3

为了抓住听众，把最好的部分前置

Berkeley NLP

EM with Features

Initialize weights w
repeat
 ● Compute expected counts e
 repeat
 Compute $\ell(w, e)$
 Compute $\nabla \ell(w, e)$
 $w \leftarrow \text{climb}(w, \ell(w, e), \nabla \ell(w, e))$
 until convergence
 ● Transform w to θ
until convergence



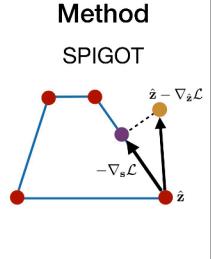
模型部分有取舍，用好图和例子

Conclusion

Problem

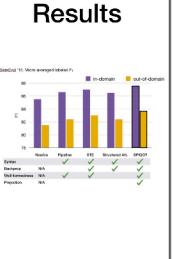
Shareholders took their money
 $\arg \max$
Shareholders took their money
A layer?
Downstream task $\nabla_\theta \mathcal{L}$?
Loss \mathcal{L} ?

Method
SPIGOT



Results

StatCat 11: Mean average precision P_r :
■ in-domain ■ out-of-domain



“结论”也有新思路

亲密性 Challenges and Contributions

- Challenges
 - deriving an optimal alignment in ambiguous situations
 - recalling more semantically matched word-concept pair without harming the alignment precision.
 - tuning the alignment with downstream parser
- Contribution
 - an enhanced aligner tuned by transition-based parser

重复

对比

对齐

四项设计的基本原则

祝大家产出优秀的学术工作