

Order-sensitive Neural Constituency Parsing

Zhicheng Wang, Tianyu Shi, Liyin Xiao, Cong Liu*

School of Computer Science and Engineering

Sun Yat-sen University

Guangzhou, China

{wangzhch23, shity3, xiaoly28}@mail2.sysu.edu.cn, liucong3@mail.sysu.edu.cn

Abstract—We propose a novel algorithm that improves on the previous neural span-based CKY decoder for constituency parsing. In contrast to the traditional span-based decoding, where spans are combined only based on the sum of their scores, we introduce an order-sensitive strategy, where the span combination scores are more carefully derived from an order-sensitive basis. Our decoder can be regarded as a generalization over existing span-based decoder in determining a finer-grain scoring scheme for the combination of lower-level spans into higher-level spans, where we emphasize on the order of the lower-level spans and use order-sensitive span scores as well as order-sensitive combination grammar rule scores to enhance prediction accuracy. We implement the proposed decoding strategy harnessing GPU parallelism and achieve a decoding speed on par with state-of-the-art span-based parsers. Using the previous state-of-the-art model without additional data as our baseline, we outperform it and improve the F1 score on the Penn Treebank Dataset by 0.26% and on the Chinese Treebank Dataset by 0.35%.

Index Terms—Machine Learning, Natural Language Processing, Constituency Parsing

I. INTRODUCTION

The application of the neural span-based decoder [14] significantly increases the constituency parsing accuracy. The majority of span-based constituency parsers used today are based on the encoder-decoder architectures [22], in which the encoder converts input sentences to vector arrays representing words and then uses a scorer to calculate a score for each span, and finally, the decoder uses the span scores to construct a constituency tree. After replacing the original LSTM model with a transformer-based [23] encoder, the accuracy of constituency parsing is increased further.

The introduction of neural networks simplifies the decoder of the algorithm. While the traditional CKY algorithm matches each non-terminal node in the parse tree with the grammar rule formed by the left child node, right child node, and their parent node, neural CKY utilizes a span scorer to perform chart-parsing, which applies a dynamic programming algorithm to determine the maximum value of the subtree score from bottom to top in order to construct the globally optimal parse trees. This is an efficient approach that results in high parsing performance.

However, this span-based neural decoder relies on the neural scorer to ensure the correctness of the compositional spans and the order between the left and right child spans in these compositions, which is not explicitly modeled in the

neural network. As illustrated in Fig. 1, the neural scorer will give a child span the same score without considering its order; therefore, correctness is not guaranteed. But the traditional CKY algorithm elicits spans based on grammar, which includes information about the label and the order of its child spans. For example, the grammar rule $NP \rightarrow DT\ NN$ implies that for parent node NP, DT is the left node, and NN is the right node. Because there is no rule for DT as a right child, the traditional method implicitly contains order information.

To attack this problem, we propose an order-sensitive scorer in our decoder. In our model, we carefully consider how the order among the children can affect the score of a subtree. Specifically, we determine the per root label score of each subtree based on four values: the per-label span score of the left split, the per-label span score of the right split, the per-label score of the parent span, and the score of composition grammar rule which consists of the labels of the left and right child spans and that of the parent span. To model the difference between the order of spans, we define two scores for each of the scores above.

After incorporating this strategy, as depicted in Fig. 1, the tree score of the word *The* as the left node should be greater than the tree score of it as the right node, so that the order information is factored in when constructing the parse tree.

Our algorithm can be computed in parallel easily. We modify the span-based neural constituency parsing model [14] and use the previous state-of-the-art model which is proposed by [9] as the baseline for our experiment. Our model’s F1 score is 0.26% higher than the baseline on the PTB [15] dataset and 0.35% higher on the CTB [17] dataset.

II. RELATED WORK

A. CKY Algorithm

The CKY algorithm [1]–[3] is the first polynomial-time parsing algorithm applicable to ambiguous CFGs that permit multiple derivations. The conventional CKY parsing algorithm can enumerate all possible parse trees for a given sentence, but it usually does not properly disambiguate between them. A neural application of the CKY algorithm (chart parsing) is proposed to handle the disambiguation. The idea behind neural CKY parsing algorithms [14] is to use a modified version of CKY to combine constituent scores and find the optimal parse tree. Our approach improves the chart parsing algorithm and introduces the order-sensitive scores which enhance the information about the span order.

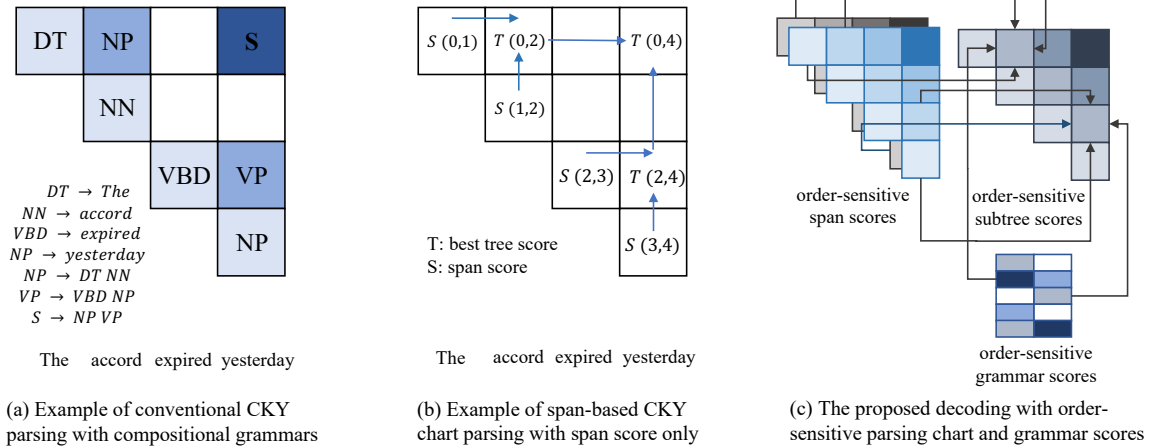


Fig. 1. Comparisons of Classic, Span-based, and Order-sensitive CKY

B. Transition-Based Parsing

Transition-based statistical parsing assigns a score to each shift or reduce actions made during the parsing process, and then chooses the parsing tree constructed from the highest scoring sequence of decisions. It is challenging for transition-based parsing to handle natural languages that are difficult to parse due to their complex semantics and long dependencies. Early attempts use classifier [5] and best-first search [6], while recent work utilizes specific tags [4] and dynamic oracles [7] to improve accuracy.

C. Span-Based Neural Constituency Parsing

Several constituency parsers proposed in the literature in recent years can be categorized as span-based neural constituency parsers. They train neural networks to assign a score to each span and then use a modified version of CKY to find an optimal parse tree, which contains spans whose scores sum up to a maximum value. In [22], the parser combines a chart decoder with a sentence encoder based on self-attention. Label attention layer [9] improves sentence embeddings by appending extra embeddings to the sentences that provide task-specific information in a way similar to prompting [30]. Head-driven phrase structure grammar [10] benefits from a uniform formalization that allows for the representation of rich contextual syntactic and semantic meanings. Recursive semi-Markov model [8] utilizes the 1-order semi-Markov model to predict the immediate children sequence of a constituent candidate. Our proposed framework can be computed in parallel by GPU without introducing higher complexity.

III. METHODOLOGY

A. Order-sensitive Decoding Algorithm

We count the number of labels after splitting to left child and right child after binarization following the same method in [22] for the training set and the test set of PTB, as shown in Table I, which demonstrates a significant difference in the occurrence probability of label's order.

TABLE I
NUMBERS OF LABELS AS LEFT AND RIGHT CHILDREN IN PTB DATASET

Label ^a	Train		Test	
	L	R	L	R
NP	82466	114615	5016	6692
VP	253	68520	16	4173
PP	333	51629	21	3042
S	169	31025	4	1874
SBAR	46	15443	5	952
WHNP	6777	394	395	25
ADJP	1408	5175	89	348
QP	2442	185	136	5
WHADVP	1934	0	124	0
ADVP	1103	2185	64	148

^a The ten most frequently occurring labels are selected in the table.

Table I show that the occurrence probability between left and right child nodes is similar for NP, ADVP in PTB. However, there are more labels with a greater difference of probability for its order, such as VP, PP, SBAR, whose number of occurrences as right child node is at least 10 times that of left, whereas the number of occurrences as left child node for QP, WHNP, WHADVP is much greater than that of the right child node. In this case, it is illogical to assign the same score to each span label when it represents in a different order. For example, WHADVP, whose span score as a right child node should be as lower as possible. As a result, we develop an order-sensitive algorithm to address this issue.

After the encoder section, we obtain the score for each span from the encoder. The purpose of the Order-sensitive Decoding Algorithm is to build a parse tree from the scores. The parse tree's score is calculated by summing the subtree scores of the constituent tree:

$$t(T) = \sum_{(i,j,l) \in T} t(i,j,l) \quad (1)$$

where T represents the parsing tree, and t is the subtree score.

The above equation can be used to represent any type of parse tree obtained through any partition, and the most

reasonable parse tree is the final parse tree with the highest tree score:

$$T^* = \arg \max_T t(T) \quad (2)$$

When span size is 1,

$$t_{best}(i, j) = \max_l s(i, j, l) \quad (3)$$

When span size is larger than 1,

$$t_{best}(i, j) = \max_l s(i, j, l) + \max_k [t_{best}(i, k) + t_{best}(k, j)] \quad (4)$$

where l is the label, $s(i, j, l)$ is the span score of span (i, j) .

However, the preceding strategy disregards the difference between the order of labels and assigns identical scores to each span as the left or right subtree.

To address this problem, we use neural networks to compute distinct scores $s(i, j, l, o)$ for the same span regarding their order where o is the order of a span. In a binarized tree, o is L for the left child span and R for the right child span.

Secondly, we elicit the compositional grammar from the training set, based on which we create a list of grammar rule scores $g(G, o)$, where o represents the order of the parent node in the grammar rule. The grammar rule score chart is randomly initialized and is the parameter in the parsing model.

With the order sensitive span scores and the grammar rule scores introduced above, we can define our order-sensitive decoding as follows:

When span size is 1:

$$t(i, j, l, o) = s(i, j, l, o) \quad (5)$$

When span size is larger than 1,

$$t(i, j, l, o) = s(i, j, l, o) + \max_{i < k < j, G: l \rightarrow l_1 l_2} [t(i, k, l_1, L) + t(k, j, l_2, R) + g(G, o)] \quad (6)$$

The optimal tree is given by:

$$t_{best}(0, n) = \max_l t(0, n, l, L) \quad (7)$$

since the order as a child is meaningless for the root node.

Our order-sensitive decoding algorithm for binarized tree constituency parsing is shown in the Algorithm 1. In the algorithm, each subtree is the sum of four scores. Compared with the conventional chart parsing, our span scores $s(i, j, l, o)$ are order-sensitive. For binarized trees, they are given by $s(i, j, l, L)$ and $s(i, j, l, R)$ for left spans and right spans respectively. For the subtree over the left child span (i, k) , we use the order sensitive score $t(i, k, l_1, L)$, while for the subtree over the right child span, we use $t(k, j, l_2, R)$. Finally, we have a grammar rule score $g(G, o)$ for each grammar rule $G: l \rightarrow l_1 l_2$, which is also order-sensitive.

In our implementation, all computations within the inner loop of the same span are implemented using batchified GPU

Algorithm 1 Order-sensitive Decoding Algorithm for Binarized Trees

Input: Sentence length n , ordered span scores $s(i, j, l, o), 0 \leq i < j \leq n$, order rule scores $g(G, o)$, for grammar rule $G: l \rightarrow l_1 l_2$.

Output: Tree T^* with maximum value $t(T^*)$.

```

1: for  $w \leftarrow 1$  to  $n$  do ▷ span size
2:   for  $i \leftarrow 0$  to  $n - w + 1$  do
3:      $j \leftarrow i + w$ 
4:     for all labels  $l$  do
5:       if  $w = 1$  then ▷ span size 1
6:          $t(i, j, l, o) \leftarrow s(i, j, l, o)$  ▷  $o$  in  $L, R$ 
7:       else ▷ span size larger than 1
8:          $t(i, j, l, o) \leftarrow s(i, j, l, o) +$ 
           $\max_{i < k < j, G: l \rightarrow l_1 l_2} [t(i, k, l_1, L) + t(k, j, l_2, R) + g(G, o)]$ 
9:       end if
10:    end for
11:  end for
12: end for
13:  $t_{best}(0, n) \leftarrow \max_l t(0, n, l, L)$  ▷  $t(T^*)$ 
14: Trace back the tree  $T^*$ 

```

operations. The key concept is to compress the computation of the same span size into a large tensor calculation. Because all parameters involved in the computation for subtree scores with the same span size are independent of each other, the calculation of one sentence can be accomplished in n steps, where n is the length of the sentence.

B. Model

Our parser, which is shown in the Fig. 2, is based on the encoder-decoder model [22], with some improvements resulting from subsequent research. Our model uses XLNet [11], incorporates a label attention layer [9] to improve self-attention [12], and augments the scoring layer and decoder [14] with our order-sensitive decoder.

The model passes through three modules in the encoder: Token Representation shown in Section III-B1, which converts the word to a vector; Attention Layers shown in Section III-B2, which augments the vector with content and position information; and Scoring Layer shown in Section III-B3, which calculates the vector corresponding to the span and scores it.

Following the encoder section, we obtain the score for each span from the scorers. We construct the parse tree from the span score using our order-sensitive decoding algorithm in the decoder shown in Section III-B4.

1) *Token Representation:* Our model’s token representation is bipartite. The initial step is to vectorize the words in the sentence s using pretrained language models.

$$s = \{w_1, w_2, w_3, \dots, w_n\} \quad (8)$$

This section makes use of two pre-trained language models: XLNet [11] for the Penn Treebank Dataset [15] and BERT-Chinese [16] for the Chinese Treebank Dataset [17]. This step

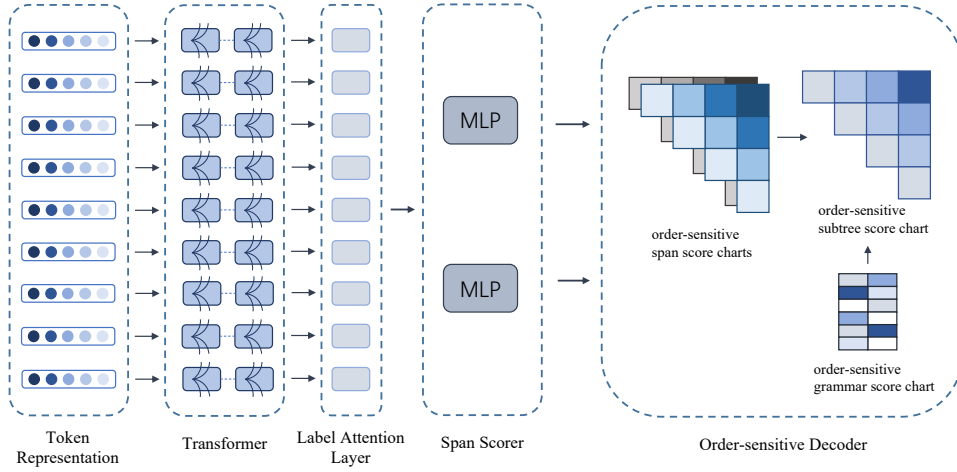


Fig. 2. Structure of Our Model

extracts the content of each word using pre-trained language models.

2) *Attention Layers*: Following the initial embedding of sentences, it enters the attention layers, which are divided into three sections. First, the word vector enters the multilevel embedding module, where the sentences' content and position features are bound to the vector:

$$E_n = [v_n; p_n] \quad (9)$$

where v_n means the token representation passed from the previous section, p_n is the position embedding.

The processed vector is then passed through 12 layers of the normal self-attention layer, which is initialized following the method from [12]. After that, the processed vector is given to a label attention layer [9]. This layer of attention assigns a query vector to each semantic label, thereby substituting the label for the attention heads.

3) *Scoring layer*: In the scoring layer, span representations adhere to the definitions from [22] and [13]. The span vector $v(i, j)$ is calculated based on the span containing the words from i to j :

$$v(i, j) = [\vec{h}_j - \vec{h}_i; \overset{\leftarrow}{h}_{j+1} - \overset{\leftarrow}{h}_{i+1}] \quad (10)$$

where \vec{h}_i and $\overset{\leftarrow}{h}_i$ is the backward and forward expression of the i -th word embedding E_n . The backward and forward representation can be obtained by splitting the representation in half. The span vector $v(i, j)$ in the previous work is then passed through a multilayer perceptron [18] span classifier with two fully connected layers and a single ReLU activation function [19], the output dimensionality of which is equal to the number of possible non-terminal labels:

$$s(i, j, \cdot) = W_2 \text{ReLU}(\text{LN}(W_1 v(i, j))) \quad (11)$$

where W_1 and W_2 are the fully connected layers, $s(i, j)$ is the score of span for the specific label, LN is the layer

normalization. Because the order needs to be considered in this section, each label of each span will have two scores, one for the left subtree and one for the right subtree, resulting in a different scoring equation:

$$s(i, j, o, \cdot) = W_{o,2} \text{ReLU}(\text{LN}(W_{o,1} v(i, j))) \quad (12)$$

where o is the order of spans, $W_{o,2}$ and $W_{o,1}$ are two fully connected layers for the specific order.

In our model, we use two MLPs to generate span score of different order separately.

4) *Decoder*: In this section, we use our improved order-sensitive neural CKY algorithm to form the decoder, which is parallelized in the computation. The model converts the input score into a parse tree and extracts the labels of tree nodes. We use the ordered span score chart from the span scorer and the order rule score generated by the grammar rules. The rule score in the rule score chart, in particular, can be learned during the training process.

5) *Loss*: The model is trained to predict the correct parse tree T^* . In the model, we can predict different tree T with different scores. It can be inferred that the following constraint must be satisfied:

$$t(T^*) \geq t(T) + \Delta(T, T^*) \quad (13)$$

where Δ is the Hamming loss for the label of spans, $t(T^*)$ is the score of the correct tree. The loss function in the model is the hinge loss:

$$L = \max(\max_T (t(T) + \Delta(T, T^*) - t(T^*)), 0) \quad (14)$$

In the training step, we minimize the loss to achieve the optimal model.

IV. EXPERIMENT

We analyze our model on the English Penn Tree-bank (PTB) [15] and Chinese Treebank 5.1 (CTB) [17] which has the data splitting from [21] and left-branching binarization. Following

normal procedure, we apply the EVALB [20] method for assessing the F1 score.

A. Setup

In the English experiment (PTB), we use large-cased pre-trained XLNet, followed by a 12 self-attention layer. For the label attention layer, we follow the setting of [9], using 112 heads (one per syntactic category), and for the MLP, we utilize two two-layer MLPs with the hidden layer of dimension 250.

In our Chinese experiments (CTB), we use BERT as our pre-training model. Except for the label attention layer, which uses 64 heads [9], the rest of parameters are consistent with the English experiments.

Our specific parameters are depicted in Table II.

TABLE II
HYPER-PARAMETERS IN OUR EXPERIMENT.

Parameter	value
batch size	32
decay factor	0.5
max decay	3
attention layer	12
LAL head	112 / 64 ^a
learning rate	3e-5
decay patience	5
dropout	0.2
MLP layer	2
MLP hidden layer	250
rule chart min	-1e6

^a 112 heads for PTB and 64 heads for CTB.

B. Performance

Table III indicates the results of our proposed parser on the PTB dataset. The baseline model in the table refers to the result after removing the order-sensitive component, that is, using the XLNet and the label attention layer in the encoder, and using ordinary neural CKY parsing in the decoder. Our model has 0.26% improvements on PTB compared with the baseline. In comparison to other models that do not incorporate data from dependency parsing, our model’s accuracy is comparable to the state-of-the-art model.

TABLE III
CONSTITUENCY PARSING PERFORMANCE ON PTB TEST SET.

Model	Precision	Recall	F1
Syntactic Distance [29]	92.00	91.70	91.80
Local Models [28]	92.50	92.20	92.40
Sequence Tagging [27]	-	-	90.60
HPSG [26] ^a	96.21	96.46	96.33
Multilingual [23]	95.73	95.46	95.59
CRF [25]	95.85	95.53	95.69
LAL [9] ^a	96.53	96.24	96.38
Linearization [24]	95.50	96.10	95.80
Recursive semi-Markov [8]	96.29	95.55	95.92
Baseline	95.83	95.42	95.61
Our model	95.97	95.77	95.87

^a Use the data from Dependency Parsing.

Table IV indicates the results of CTB dataset. The correctness of our model is 0.35% higher than the baseline. Among

models without using dependency parsing data, our model establishes new state-of-the-art results.

TABLE IV
CONSTITUENCY PARSING PERFORMANCE ON CTB 5.1 TEST SET.

Model	Precision	Recall	F1
Syntactic Distance [29]	86.60	86.40	86.50
Local Models [28]	87.50	87.10	87.30
Sequence Tagging [27]	-	-	85.61
HPSG [26] ^a	92.33	92.03	92.18
Multilingual [23]	91.96	91.55	91.75
CRF [25]	92.51	92.04	92.27
LAL [9] ^a	93.45	91.85	92.64
Linearization [24]	92.70	92.20	92.40
Recursive semi-Markov [8]	92.94	92.06	92.50
Baseline	91.99	92.33	92.16
Our model	92.55	92.46	92.51

^a Use the data from Dependency Parsing.

C. Ablation Study

TABLE V
ABLATION STUDY ON PTB TEST SET.

Model	Precision	Recall	F1
Baseline	95.83	95.42	95.61
Ordered span score	95.92	95.65	95.79
Order & Grammar rule score	95.97	95.77	95.87

We next conduct an ablation study to discover more about the contributions of various components in our proposed model.

a) *Impact of grammar rule scores.*: In our parsing model, we use the ordered span scores and the grammar rule scores to calculate the parsing tree, so we remove the grammar rule scores to evaluate its impact. That is, we simplify our decoder using the following subtree score:

$$t(i, j, o) = \max_l [s(i, j, l, o) + \max_{i < k < j} (t(i, k, L) + t(k, j, R))] \quad (15)$$

The results depicted in the Table V clearly show that the introduction of the grammar rule chart improves the F1 score of the constituency parsing.

b) *Impact of order-sensitive decoding.*: We can evaluate the impact of the order-sensitive decoding by comparing the results of the baseline with the simplified model with ordered span score. The simplified model can have 0.18% improvement of F1 score.

D. Speed Comparison

We average the results of parsing speed from 20 trials which are measured on the PTB test set. With a single Tesla V100, the average processing rate is 65 sentences per second. Table VI shows the speed comparison between our model and the baseline. The speed is slightly slower in live testing, but considering the increase in accuracy, the processing speed is acceptable.

TABLE VI
SPEED COMPARISON ON PTB TEST SET.

Model	Sents/sec
Baseline	103
Ordered span score	89
Ordered & grammar rule score	65

V. CONCLUSION

In this paper, we introduce the order-sensitive decoding algorithm to the neural CKY constituency parser which combines order-sensitive span scores and ordered grammar rule scores, with the benefit of enhancing the prediction accuracy for the parse trees. We discuss the advantages and necessity of the order-sensitive decoding in the label statistics and ablation study sections. Our model is not time-consuming compared with previous model computing on GPU. Experiments demonstrate that our model outperforms the baseline model, which incorporates elements of the previous state-of-the-art model on the PTB and CTB datasets.

REFERENCES

- [1] J. Cocke, "Programming languages and their compilers: Preliminary notes," 1969.
- [2] T. Kasami, "An efficient recognition and syntax-analysis algorithm for context-free languages," 1965.
- [3] D. H. Younger, "Recognition and parsing of context-free languages in time n^3 ," *Information & Computation*, vol. 10, pp. 189–208, 1967.
- [4] N. Kitaev and D. Klein, "Tetra-tagging: Word-synchronous parsing with linear-time inference," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 6255–6261. [Online]. Available: <https://aclanthology.org/2020.acl-main.557>
- [5] K. Sagae and A. Lavie, "A classifier-based parser with linear run-time complexity," in *Proceedings of the Ninth International Workshop on Parsing Technology*. Vancouver, British Columbia: Association for Computational Linguistics, Oct. 2005, pp. 125–132. [Online]. Available: <https://aclanthology.org/W05-1513>
- [6] —, "A best-first probabilistic shift-reduce parser," in *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*. Sydney, Australia: Association for Computational Linguistics, Jul. 2006, pp. 691–698. [Online]. Available: <https://aclanthology.org/P06-2089>
- [7] J. Cross and L. Huang, "Span-based constituency parsing with a structure-label system and provably optimal dynamic oracles," *arXiv: Computation and Language*, 2016.
- [8] X. Xin, J. Li, and Z. Tan, "N-ary constituent tree parsing with recursive semi-markov model," in *Meeting of the Association for Computational Linguistics*, 2021.
- [9] K. Mrini, F. Démoncourt, Q. H. Tran, T. Bui, W. Chang, and N. Nakashole, "Rethinking self-attention: Towards interpretability in neural parsing," in *Empirical Methods in Natural Language Processing*, 2020.
- [10] I. A. Sag and C. J. Pollard, "Head-driven phrase structure grammar," 1994.
- [11] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in *Neural Information Processing Systems*, 2019.
- [12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2017.
- [13] D. Gaddy, M. Stern, and D. Klein, "What's going on in neural constituency parsers? an analysis," *arXiv: Computation and Language*, 2018.
- [14] M. Stern, J. Andreas, and D. Klein, "A minimal span-based neural constituency parser," in *Meeting of the Association for Computational Linguistics*, 2017.
- [15] M. Marcus, M. A. Marcinkiewicz, and B. Santorini, "Building a large annotated corpus of english: the penn treebank," *Computational Linguistics*, vol. 19, pp. 313–330, 1993.
- [16] Y. Cui, W. Che, T. Liu, B. Qin, and Z. Yang, "Pre-training with whole word masking for chinese bert," *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, vol. 29, p. 3504–3514, jan 2021. [Online]. Available: <https://doi.org/10.1109/TASLP.2021.3124365>
- [17] N. Xue, F. Xia, F.-D. Chiou, and M. Palmer, "The penn chinese treebank: Phrase structure annotation of a large corpus," *Natural Language Engineering*, vol. 11, pp. 207–238, 2005.
- [18] F. Rosenblatt, "Principles of neurodynamics. perceptrons and the theory of brain mechanisms," Cornell Aeronautical Lab Inc Buffalo NY, Tech. Rep., 1961.
- [19] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, G. Gordon, D. Dunson, and M. Dudík, Eds., vol. 15. Fort Lauderdale, FL, USA: PMLR, 11–13 Apr 2011, pp. 315–323. [Online]. Available: <https://proceedings.mlr.press/v15/glorot11a.html>
- [20] S. Sekine and M. Collins, "Evalb bracket scoring program," *URL: http://www.cs.nyu.edu/cs/projects/teus/proteus/evalb*, 1997.
- [21] J. Liu and Y. Zhang, "Shift-reduce constituent parsing with neural lookahead features," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 45–58, 2017. [Online]. Available: <https://aclanthology.org/Q17-1004>
- [22] N. Kitaev and D. Klein, "Constituency Parsing with a Self-Attentive Encoder," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018. [Online]. Available: <http://aclweb.org/anthology/P18-1249>
- [23] N. Kitaev, S. Cao, and D. Klein, "Multilingual constituency parsing with self-attention and pre-training," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 3499–3505. [Online]. Available: <https://aclanthology.org/P19-1340>
- [24] Y. Wei, Y. Wu, and M. Lan, "A span-based linearization for constituent trees," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 3267–3277. [Online]. Available: <https://aclanthology.org/2020.acl-main.299>
- [25] Y. Zhang, H. Zhou, and Z. Li, "Fast and accurate neural crf constituency parsing," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, C. Bessiere, Ed. International Joint Conferences on Artificial Intelligence Organization, 7 2020, pp. 4046–4053, main track. [Online]. Available: <https://doi.org/10.24963/ijcai.2020/560>
- [26] J. Zhou and H. Zhao, "Head-Driven Phrase Structure Grammar parsing on Penn Treebank," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 2396–2408. [Online]. Available: <https://aclanthology.org/P19-1230>
- [27] D. Vilares, M. Abdou, and A. Søgaard, "Better, faster, stronger sequence tagging constituent parsers," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 3372–3383. [Online]. Available: <https://aclanthology.org/N19-1341>
- [28] Z. Teng and Y. Zhang, "Two local models for neural constituent parsing," in *Proceedings of the 27th International Conference on Computational Linguistics*. Santa Fe, New Mexico, USA: Association for Computational Linguistics, Aug. 2018, pp. 119–132. [Online]. Available: <https://aclanthology.org/C18-1011>
- [29] Y. Shen, Z. Lin, A. P. Jacob, A. Sordani, A. Courville, and Y. Bengio, "Straight to the tree: Constituency parsing with neural syntactic distance," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, Jul. 2018, pp. 1171–1180. [Online]. Available: <https://aclanthology.org/P18-1108>
- [30] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *CoRR*, vol. abs/2107.13586, 2021. [Online]. Available: <https://arxiv.org/abs/2107.13586>