

---

# Efficient Beam Tree Recursion (Appendix)

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Organization

In Section 2, we describe the settings of all the tasks and datasets that we have tested our models on. In Section 3, we provide additional results on logical inference and sentiment classification. Then, in Section 4, we present an extended survey of related works. In Section 5, we detail our architecture setup including the sequence-interaction models. In Section 6, we provide our hyperparameters. Last, further details on BT-RvNN (if necessary) can be found in the document titled “Beam Tree Recursive Cells” provided in the supplementary (however, we described the salient aspects of the model in the main paper).

## 2 Task Details

**ListOps:** ListOps was introduced by Nangia and Bowman [36] and is a task for solving nested lists of mathematical operations. It is a 10-way classification task. Similar to Chowdhury and Caragea [6], we train our models on the original training set with all samples  $\geq 100$  sequence lengths filtered out. We use the original development set for validation. We test on the following sets: the original test set (near-IID split); the length generalization splits from Havrylov et al. [20] that include samples of much higher lengths; the argument generalization splits from Anonymous [1] that involve an unseen number of maximum arguments for each operator; and the LRA split (which has both higher sequence length and higher argument number) from Tay et al. [52].

**Logical Inference:** Logical Inference was introduced by Bowman et al. [3] and is a task that involves classifying fine-grained inferential relations between two given sequences in a form similar to that of formal sentences of propositional logic. Similar to Tran et al. [54], our models were trained on splits with logical connectives  $\leq 6$ . We show the results in OOD test sets with logical connections 10-12. We use the same splits as Shen et al. [44], Tran et al. [54], Chowdhury and Caragea [6].

**SST5:** SST5 is a fine-grained 5-way sentiment classification task introduced by Socher et al. [50]. We use the original splits.

**IMDB:** IMDB is a binary sentiment classification task from Maas et al. [31]. We use the same train, validation, and IID test sets as created in Anonymous [1]. We also use the contrast set Gardner et al. [15] and counterfactual set Kaushik et al. [26] as additional test splits.

**QQP:** QQP<sup>1</sup> [25] is a task of classifying whether two given sequences in a pair are paraphrases of each other or not. As standard Wang et al. [56], we randomly sample 10,000 samples for validation and IID test set such that for each split 5,000 samples are maintained to be paraphrases and the other 5,000 are maintained to be not paraphrases. We also use the adversarial test sets PAWS<sub>QQP</sub> and PAWS<sub>WIKI</sub> from Zhang et al. [62].

**SNLI:** SNLI [2] is a natural language inference (NLI) task. It is a 3-way classification task to classify the inferential relation between two given sequences. We use the same train, development, and IID test set splits as in Chowdhury and Caragea [6]. Any data with a sequence of length  $\geq 150$  is filtered

---

<sup>1</sup><https://data.quora.com/First-Quora-Dataset-Release-QuestionPairs>

Table 1: Mean accuracy and standard deviation on the Logical Inference [3] for  $\geq 10$  number of operations after training on samples with  $\leq 6$  operations, and on SST5 [50] and IMDB [31]. **Count.** represents counterfactual test split from Kaushik et al. [26] and **Cont.** represents contrast test split from Gardner et al. [15] The best results are shown in bold. Our models were run 3 times on different seeds. Subscript represents standard deviation. As an example,  $90_1 = 90 \pm 0.1$

Model	Logical Inference			SST5	IMDB		
	Number of Operations			IID	IID	Cont.	Count.
	10	11	12				
GT-GRC	90.33 <sub>22</sub>	88.43 <sub>18</sub>	85.70 <sub>24</sub>	51.67 <sub>8.8</sub>	85.11 <sub>10</sub>	70.63 <sub>21</sub>	81.97 <sub>5</sub>
EGT-GRC	75.79 <sub>61</sub>	73.38 <sub>68</sub>	69.68 <sub>7.8</sub>	51.63 <sub>14</sub>	86.58 <sub>2.7</sub>	72 <sub>9.2</sub>	81.76 <sub>14</sub>
CRvNN	94.51 <sub>2.9</sub>	<b>94.48</b> <sub>5.6</sub>	92.73 <sub>15</sub>	51.75 <sub>11</sub>	91.47 <sub>1.2</sub>	76.98 <sub>5.8</sub>	83.68 <sub>7.8</sub>
OM	94.95 <sub>2</sub>	93.9 <sub>2.2</sub>	93.36 <sub>6.2</sub>	52.30 <sub>2.7</sub>	<b>91.69</b> <sub>0.5</sub>	<b>77.80</b> <sub>15</sub>	<b>85.38</b> <sub>3.5</sub>
BT-GRC	95.04 <sub>2.3</sub>	94.29 <sub>3.8</sub>	93.36 <sub>2.4</sub>	<b>52.32</b> <sub>4.7</sub>	91.29 <sub>1.2</sub>	75.07 <sub>29</sub>	82.86 <sub>23</sub>
BT-GRC OS	<b>95.43</b> <sub>4.5</sub>	94.21 <sub>6.6</sub>	<b>93.39</b> <sub>1.5</sub>	51.92 <sub>7.2</sub>	90.86 <sub>9.3</sub>	75.68 <sub>21</sub>	84.77 <sub>11</sub>
EBT-GRC	94.95 <sub>1.5</sub>	93.87 <sub>7.4</sub>	93.04 <sub>6.7</sub>	52.22 <sub>1</sub>	91.47 <sub>1.2</sub>	76.16 <sub>17</sub>	84.29 <sub>12</sub>

36 out from the training set for efficiency. We use also additional test set splits for stress tests. We use  
 37 the hard test set split from Gururangan et al. [19], the break test set from Glockner et al. [16], and the  
 38 counterfactual test set from Kaushik et al. [26].

39 **MNLI:** MNLI [57] is another NLI dataset, which is similar to SNLI in format. We use the original  
 40 development sets (match and mismatch) as test sets. We filter out all data with any sequence  
 41 length  $\geq 150$  from the training set. Our actual development set is a random sample of 10,000  
 42 data-points from the filtered training set. As additional testing sets, we use the development set of  
 43 Conjunctive NLI (ConjNLI) [41] and a few of the stress sets from Naik et al. [35]. These stress test  
 44 sets include - Negation Match (NegM), Negation Mismatch (NegMM), Length Match (LenM), and  
 45 Length Mismatch (LenMM). NegM and NegMM add tautologies containing “not” terms - this can  
 46 bias the models to classify contradiction as the inferential relation because the training set contains  
 47 spurious correlations between existence of “not” related terms and the class of contradiction. LenM  
 48 and LenMM add tautologies to artificially increase the lengths of the samples without changing the  
 49 inferential relation class.

### 50 3 Additional Results

51 In Table 1, we show that our EBT-GRC model can keep up fairly well with BT-GRC and BT-GRC  
 52 OS on logical inference [3] and sentiment classification tasks like SST5 [50], and IMDB [15] while  
 53 being much more computationally efficient as demonstrated in the main paper.

### 54 4 Extended Related Works

55 **RvNN History:** Recursive Neural Networks (RvNNs) in the more specified sense of building  
 56 representations through trees and directed acyclic graphs were proposed in [40, 18]. Socher et al.  
 57 [48, 49, 50] extended the use of RvNNs in Natural Language Processing (NLP) by considering  
 58 constituency trees and dependency trees. A few works [63, 51, 28, 64] started adapting Long  
 59 Shot-term Memory Networks [21] as a cell function for recursive processing. Le and Zuidema  
 60 [29], Maillard et al. [33] proposed a chart-based method for simulating bottom-up Recursive Neural  
 61 Networks through dynamic programming. Shi et al. [47], Munkhdalai and Yu [34] explored heuristics-  
 62 based tree-structured RvNNs.

63 RvNNs can also be simulated by stack-augmented recurrent neural networks (RNNs) to an extent  
 64 (similar to how pushdown automata can model context-free grammar [42, 27]). There are multiple  
 65 works on stack-augmented RNNs [4, 60, 32]. Ordered Memory [44] is one of the more modern such  
 66 examples. More recently, DuSell and Chiang [12, 13] explored non-deterministic stack augmented  
 67 RNNs and [9] explored other expressive models. Wu [58] presented a survey of latent structure  
 68 models.

69 Choi et al. [5] proposed a greedy search strategy based on easy-first algorithm [17, 30] for auto-  
70 parsing structures for recursion utilizing STE gumbel softmax for gradient signals. Peng et al. [38]  
71 extended the framework with SPIGOT and Havrylov et al. [20] extended it with reinforcement  
72 learning (RL). Anonymous [1] extended it with beam search and soft top-k. Chowdhury and Caragea  
73 [6], Zhang et al. [61] introduced different forms of soft-recursion.

74 **Top-down Signal:** Similar to us, Teng and Zhang [53] explored bidirectional signal propagation  
75 (bottom-up and top-down). However, they sent top-down signal in a sequential manner which  
76 can be expensive - either it can get slow without parallelization or memory-wise expensive with  
77 parallelization of contextualization of nodes in the same height. Our approach in EBT-GAU also  
78 has some kinship with BP-Transformer [59]. BP-Transformer allows message passing between a  
79 fixed subset of parent nodes and terminal nodes created using a heuristics-based balanced binary tree.  
80 Chart-based models can also create sequence contextualized representations [10, 11] but they can be  
81 quite expensive by default [1] needing their own separate techniques [22, 23].

82 **Transformers + RvNNs:** There have been several approaches to incorporating RvNN-like inductive  
83 biases to Transformers. For instance, Universal Transformer [8] introduced weight-sharing and  
84 dynamic halt to Transformers. Csordás et al. [7] extended on universal transformer with geometric  
85 attention for locality bias and gating. Shen et al. [46] built on weight-shared transformers with high  
86 layer depth and group self-attention. Wang et al. [55], Nguyen et al. [37], Shen et al. [45] added  
87 hierarchical structural biases to self-attention. Fei et al. [14] biased pre-trained Transformers to have  
88 constituent information in intermediate representations. Hu et al. [22] used Transformer as binary  
89 recursive cells in chart-based encoders.

## 90 5 Architecture details

### 91 5.1 Sentence Encoder Models

92 For the sentence encoder models the architectural framework we use is the same siamese dual-encoder  
93 setup as Anonymous [1].

### 94 5.2 Sentence Interaction Models

95 **GAU-Block:** Our specific implementation of a GAU-block [24] is detailed below. Our GAU-  
96 Block can be defined as  $\text{GAUBlock}(x, p, G)$ . The function arguments are of the following forms:  
97  $x \in \mathbb{R}^{n \times d}$ ,  $p \in \mathbb{R}^{l \times d}$  and  $G \in \{0, 1\}^{n \times l}$ .  $x$  accepts the main sequence of vectors that is to serve as  
98 attention queries;  $p$  accepts either the sequence of intermediate node representations created from our  
99 RvNN (for parent attention) or it accepts the same input as  $x$  (for usual cases);  $p$  serves as keys and  
100 values for attention;  $G$  accepts either the adjacency matrix in case of parent attention (where  $G_{ij} = 1$   
101 iff  $p_j$  is a parent of  $x_i$  else  $G_{ij} = 0$ ), otherwise, it accepts just the usual attention mask; either way,  
102  $G$  serves as an attention mask.

$$103 \quad x' = \text{LN}(xW_{init} + b_{init}); \quad p' = \text{LN}(pW_{init} + b_{init}) \quad (1)$$

$$104 \quad u = \text{SiLU}(x'W_u + b_u); \quad v = \text{SiLU}(p'W_v + b_v) \quad (2)$$

$$105 \quad q = z_q \odot \text{SiLU}(x'W_z + b_z) + zb_q; \quad k = z_k \odot \text{SiLU}(p'W_z + b_z) + zb_k \quad (3)$$

$$106 \quad A = \text{Softmax}\left(\frac{qk^T + pos}{\sqrt{2d}}, \text{mask} = G\right) \quad (4)$$

$$107 \quad v' = Av \quad (5)$$

$$108 \quad o = (u \odot v')W_o + b_o \quad (6)$$

$$109 \quad g = \text{Sigmoid}([o; x]W_{gate} + b_{gate}) \quad (7)$$

$$110 \quad \text{out} = g \odot o + (1 - g) \odot x \quad (8)$$

110 Here,  $W_{init} \in \mathbb{R}^{d \times d}$ ;  $W_z \in \mathbb{R}^{d \times d_h}$ ,  $W_u, W_v \in \mathbb{R}^{d \times 2d}$ ,  $b_{init}, b_z, b_o \in \mathbb{R}^d$ ;  $z_q, zb_q, z_k, zb_k \in$   
111  $\mathbb{R}^{d_h}$ ;  $b_u, b_v \in \mathbb{R}^{2d}$ ,  $W_o, W_{gate} \in \mathbb{R}^{2d \times d}$ .  $[\cdot; \cdot]$  represents concatenation.

112  $\text{LN}$  is layer normalization.  $pos$  is calculated using the technique of Shaw et al. [43] using relative  
113 tree height distance for parent attention, or relative positional distance for usual cases.

114 **GAU Sequence Interaction Setup:** Let GAUStack represent some arbitrary number of compositions  
 115 of GAUBlocks (multilayered GAU block). GAUStack has the same function arguments as GAUBlock.  
 116 Given two sequences  $(x_1, x_2)$  and their corresponding attention masks  $(M_1, M_2)$  as inputs where  
 117  $x_1 \in \mathbb{R}^{n_1 \times d}, x_2 \in \mathbb{R}^{n_2 \times d}, M_1 \in \{0, 1\}^{n_1 \times n_1}, M_2 \in \{0, 1\}^{n_2 \times n_2}$ , the GAU setup can be expressed  
 118 as:

$$inp = [CLS + seg_1; x_1 + seg_1; SEP; CLS + seg_2, x_2 + seg_2] \quad (9)$$

$$r = \text{GAUStack}(x = inp, p = inp, G = [M_1; M_2]) \quad (10)$$

$$\alpha = \text{Softmax}(\text{GELU}(rW_1 + b_1)W_2 + b_2) \quad (11)$$

$$cls' = \sum_i \alpha_i r \quad (12)$$

$$logits = \text{GELU}(cls'W_1^{logits} + b_1^{logits})W_2^{logits} + b_2^{logits} \quad (13)$$

123 Here,  $CLS, SEP, seg_1, seg_2 \in \mathbb{R}^{1 \times d}$  are randomly initialized trainable vectors;  $seg_1, seg_2$  are  
 124 segment embeddings.  $W_1 \in \mathbb{R}^{d \times d}, W_2 \in \mathbb{R}^{d \times 1}; b_1, b_2, b_1^{logits} \in \mathbb{R}^d; b_2^{logits} \in \mathbb{R}^c; W_1^{logits} \in$   
 125  $\mathbb{R}^{d \times d}, W_2^{logits} \in \mathbb{R}^{d \times c}$ .  $c$  is the number of classes for the task.

126 **EGT-GAU Sequence Interaction Setup:** EGT-GAU starts from the same input as above. Let us  
 127 also assume we have the EGT-GRC( $x$ ) module which takes a sequence of vectors  $x \in \mathbb{R}^{n \times d}$  as  
 128 the input to recursively process and outputs  $(cls, p, G)$  where  $cls \in \mathbb{R}^{1 \times d}$  is the root representation,  
 129  $p \in \mathbb{R}^{l \times d}$  is the sequence of non-terminal representations from the tree, and  $G \in \{0, 1\}^{n \times l}$  is the  
 130 adjacency matrix for parent attention (i.e.,  $G_{ij} = 1$  iff  $p_j$  is a parent of  $x_i$ , else  $G_{ij} = 0$ ). Technically,  
 131 tree height information is also extracted for relative position but we do not express that explicitly for  
 132 the sake of brevity. With these elements, EGT-GAU can be expressed as below:

$$cls_1, p_1, G_1 = \text{EGT-GRC}(x = x_1); \quad cls_2, p_2, G_2 = \text{EGT-GRC}(x = x_2) \quad (14)$$

$$x'_1 = \text{GAUStack}_1(x = x_1, p = p_1, G = G_1); \quad x'_2 = \text{GAUStack}_1(x = x_2, p = p_2, G = G_2) \quad (15)$$

$$cls'_1 = \text{GELU}(cls_1W_1^{cls} + b_1^{cls})W_2^{cls} + b_2^{cls}; \quad cls'_2 = \text{GELU}(cls_2W_1^{cls} + b_1^{cls})W_2^{cls} + b_2^{cls} \quad (16)$$

$$inp = [cls'_1 + seg_1; x'_1 + seg_1; SEP; cls'_2 + seg_2, x'_2 + seg_2] \quad (17)$$

$$r = \text{GAUStack}_2(x = inp, p = inp, G = [M_1; M_2]) \quad (18)$$

137 Everything else after eqn. 18 is the same as eqn. 11 to 13.  $SEP, seg_1, seg_2 \in \mathbb{R}^{1 \times d}; seg_1, seg_2$  are  
 138 segment embeddings as before.  $W_1^{cls}, W_2^{cls} \in \mathbb{R}^{d \times d}; b_1^{cls}, b_2^{cls} \in \mathbb{R}^d$ .

139 **EBT-GAU Sequence Interaction Setup:** This setup is similar to that of EGT-GAU but with a few  
 140 changes. EBT-GAU uses EBT-GRC as a module instead of EGT-GRC. EBT-GAU returns outputs of  
 141 the form  $(cls, bp, bG, s)$  where  $cls \in \mathbb{R}^{1 \times d}$  is the beam-score-weighted-averaged root representation,  
 142  $bp \in \mathbb{R}^{b \times l \times d}$  are the beams (beam size  $b$ ) of sequences of non-terminal representations from the  
 143 tree,  $bG \in \{0, 1\}^{b \times n \times l}$  are the beams of adjacency matrices for parent attention, and  $s \in \mathbb{R}^b$  are the  
 144 softmax-normalized beam scores. Let NGAUStack represent the same function as GAUStack but  
 145 formalized for batched processing of multiple beams of sequences. With these elements, EBT-GAU  
 146 can be expressed as:

$$cls_1, bp_1, bG_1, s_1 = \text{EBT-GRC}(x = x_1); \quad cls_2, bp_2, bG_2, s_2 = \text{EBT-GRC}(x = x_2) \quad (19)$$

$$bx_1 = \text{repeat}(x_1, b); \quad bx_2 = \text{repeat}(x_2, b) \quad (20)$$

$$bx'_1 = \text{NGAUStack}_1(bx_1, bp_1, bG_1); \quad bx'_2 = \text{NGAUStack}_1(bx_2, bp_2, bG_2) \quad (21)$$

$$x'_1 = \sum_i s[i] \cdot bx'_1[i]; \quad x'_2 = \sum_i s[i] \cdot bx'_2[i] \quad (22)$$

150 Everything else after eqn. 22 is the same as the equations 16-18 followed by the equations 11 to 13.  
 151  $\text{repeat}(x, b)$  changes  $x \in \mathbb{R}^{n \times d}$  to  $bx \in \mathbb{R}^{b \times n \times d}$  by batching the same  $x$  for  $b$  times.

## 152 6 Hyperparameter details

153 For sentence encoder models, we use the same hyperparameters as [1] (the preprint of the paper  
154 is available in the supplementary in anonymized form) for all the datasets. The only new hyper-  
155 parameter for EBT-GRC is  $d_s$  which we set as 64; otherwise the hyperparameters are the same  
156 as that of BT-GRC or BT-GRC OS. We discuss the hyperparameters of the sequence interaction  
157 models next. For EBT-GAU/EGT-GAU, we used a two-layered weight-shared GAU-Blocks for  
158 NGAUStack<sub>1</sub>/GAUStack<sub>1</sub> and a three-layered weight-shared GAU-Blocks for GAUStack<sub>2</sub> (for pa-  
159 rameter efficiency and regularization). GAU uses a five-layered GAU-Blocks (weights unshared) for  
160 GAUStack so that the parameters are similar to that of EBT-GAU or EGT-GAU. We use a dropout  
161 of 0.1 after the multiplication with  $W_o$  in each GAUBlock layer and a head size  $d_h$  of 128 (similar  
162 to Hua et al. [24]). For relative position, we set  $k = 5$  ( $k$  here corresponds the receptive field for  
163 relative attention in Shaw et al. [43]) for normal GAUBlocks and  $k = 10$  for parent attention (since  
164 parent attention is only applied to higher heights, we do not need to initialize weights for negative  
165 relative distances). Other hyperparameters are kept same as the sentence encoder models. The  
166 hyperparameters of MNLI, SNLI, and QQP are shared. Note that all the natural language tasks are  
167 trained with fixed 840B Glove Embeddings [39] as in Anonymous [1]. All models were trained in a  
168 single Nvidia RTX A6000. The code is available in the supplementary.

## 169 References

- 170 [1] Anonymous. Beam tree recursive cells. In *Proceedings of the International Conference on*  
171 *Machine Learning*, 2023.
- 172 [2] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large  
173 annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference*  
174 *on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal,  
175 September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1075. URL  
176 <https://aclanthology.org/D15-1075>.
- 177 [3] Samuel R. Bowman, Christopher D. Manning, and Christopher Potts. Tree-structured compo-  
178 sition in neural networks without tree-structured architectures. In *Proceedings of the 2015th*  
179 *International Conference on Cognitive Computation: Integrating Neural and Symbolic Ap-*  
180 *proaches - Volume 1583, COCO’15*, page 37–42, Aachen, DEU, 2015. CEUR-WS.org.
- 181 [4] Samuel R. Bowman, Jon Gauthier, Abhinav Rastogi, Raghav Gupta, Christopher D. Manning,  
182 and Christopher Potts. A fast unified model for parsing and sentence understanding. In  
183 *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*  
184 *(Volume 1: Long Papers)*, pages 1466–1477, Berlin, Germany, August 2016. Association for  
185 Computational Linguistics. doi: 10.18653/v1/P16-1139. URL [https://aclanthology.org/](https://aclanthology.org/P16-1139)  
186 [P16-1139](https://aclanthology.org/P16-1139).
- 187 [5] Jihun Choi, Kang Min Yoo, and Sang-goo Lee. Learning to compose task-specific tree structures.  
188 In Sheila A. McIlraith and Kilian Q. Weinberger, editors, *Proceedings of the Thirty-Second AAAI*  
189 *Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial*  
190 *Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial*  
191 *Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5094–5101.  
192 AAAI Press, 2018. URL [https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/](https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16682)  
193 [view/16682](https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16682).
- 194 [6] Jishnu Ray Chowdhury and Cornelia Caragea. Modeling hierarchical structures with continuous  
195 recursive neural networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the*  
196 *38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*  
197 *Learning Research*, pages 1975–1988. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/chowdhury21a.html)  
198 [mlr.press/v139/chowdhury21a.html](https://proceedings.mlr.press/v139/chowdhury21a.html).
- 199 [7] Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. The neural data router: Adaptive  
200 control flow in transformers improves systematic generalization. In *International Conference on*  
201 *Learning Representations*, 2022. URL [https://openreview.net/forum?id=KBQP4A\\_J1K](https://openreview.net/forum?id=KBQP4A_J1K).

- 202 [8] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Uni-  
203 versal transformers. In *International Conference on Learning Representations*, 2019. URL  
204 <https://openreview.net/forum?id=HyzdRiR9Y7>.
- 205 [9] Gr'egoire Del'etang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot  
206 Catt, Marcus Hutter, Shane Legg, and Pedro A. Ortega. Neural networks and the chomsky  
207 hierarchy. *ArXiv*, abs/2207.02098, 2022.
- 208 [10] Andrew Drozdov, Patrick Verga, Mohit Yadav, Mohit Iyyer, and Andrew McCallum. Unsuper-  
209 vised latent tree induction with deep inside-outside recursive auto-encoders. In *Proceedings*  
210 *of the 2019 Conference of the North American Chapter of the Association for Computational*  
211 *Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1129–  
212 1141, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi:  
213 10.18653/v1/N19-1116. URL <https://aclanthology.org/N19-1116>.
- 214 [11] Andrew Drozdov, Subendhu Rongali, Yi-Pei Chen, Tim O’Gorman, Mohit Iyyer, and Andrew  
215 McCallum. Unsupervised parsing with S-DIORA: Single tree encoding for deep inside-outside  
216 recursive autoencoders. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*  
217 *Language Processing (EMNLP)*, pages 4832–4845, Online, November 2020. Association  
218 for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.392. URL [https://](https://aclanthology.org/2020.emnlp-main.392)  
219 [aclanthology.org/2020.emnlp-main.392](https://aclanthology.org/2020.emnlp-main.392).
- 220 [12] Brian DuSell and David Chiang. Learning context-free languages with nondeterministic stack  
221 RNNs. In *Proceedings of the 24th Conference on Computational Natural Language Learning*,  
222 pages 507–519, Online, November 2020. Association for Computational Linguistics. doi:  
223 10.18653/v1/2020.conll-1.41. URL <https://aclanthology.org/2020.conll-1.41>.
- 224 [13] Brian DuSell and David Chiang. Learning hierarchical structures with differentiable nonde-  
225 terministic stacks. In *International Conference on Learning Representations*, 2022. URL  
226 [https://openreview.net/forum?id=5LXw\\_Qp1BiF](https://openreview.net/forum?id=5LXw_Qp1BiF).
- 227 [14] Hao Fei, Yafeng Ren, and Donghong Ji. Retrofitting structure-aware transformer language  
228 model for end tasks. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*  
229 *Language Processing (EMNLP)*, pages 2151–2161, Online, November 2020. Association  
230 for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.168. URL [https://](https://aclanthology.org/2020.emnlp-main.168)  
231 [aclanthology.org/2020.emnlp-main.168](https://aclanthology.org/2020.emnlp-main.168).
- 232 [15] Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep  
233 Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi,  
234 Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire,  
235 Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace,  
236 Ally Zhang, and Ben Zhou. Evaluating models’ local decision boundaries via contrast sets. In  
237 *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323,  
238 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
239 findings-emnlp.117. URL <https://aclanthology.org/2020.findings-emnlp.117>.
- 240 [16] Max Glockner, Vered Shwartz, and Yoav Goldberg. Breaking NLI systems with sentences that  
241 require simple lexical inferences. In *Proceedings of the 56th Annual Meeting of the Association*  
242 *for Computational Linguistics (Volume 2: Short Papers)*, pages 650–655, Melbourne, Australia,  
243 July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-2103. URL  
244 <https://aclanthology.org/P18-2103>.
- 245 [17] Yoav Goldberg and Michael Elhadad. An efficient algorithm for easy-first non-directional  
246 dependency parsing. In *Human Language Technologies: The 2010 Annual Conference of*  
247 *the North American Chapter of the Association for Computational Linguistics*, pages 742–  
248 750, Los Angeles, California, June 2010. Association for Computational Linguistics. URL  
249 <https://aclanthology.org/N10-1115>.
- 250 [18] C. Goller and A. Kuchler. Learning task-dependent distributed representations by backprop-  
251 agation through structure. In *Proceedings of International Conference on Neural Networks*  
252 *(ICNN’96)*, volume 1, pages 347–352 vol.1, 1996. doi: 10.1109/ICNN.1996.548916.

- 253 [19] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and  
254 Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of*  
255 *the 2018 Conference of the North American Chapter of the Association for Computational*  
256 *Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New  
257 Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/  
258 N18-2017. URL <https://aclanthology.org/N18-2017>.
- 259 [20] Serhii Havrylov, Germán Kruszewski, and Armand Joulin. Cooperative learning of disjoint  
260 syntax and semantics. In *Proceedings of the 2019 Conference of the North American Chapter*  
261 *of the Association for Computational Linguistics: Human Language Technologies, Volume 1*  
262 *(Long and Short Papers)*, pages 1118–1128, Minneapolis, Minnesota, June 2019. Association  
263 for Computational Linguistics. doi: 10.18653/v1/N19-1115. URL <https://aclanthology.org/N19-1115>.
- 265 [21] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9  
266 (8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL  
267 <https://doi.org/10.1162/neco.1997.9.8.1735>.
- 268 [22] Xiang Hu, Haitao Mi, Zujie Wen, Yafang Wang, Yi Su, Jing Zheng, and Gerard de Melo.  
269 R2D2: Recursive transformer based on differentiable tree for interpretable hierarchical lan-  
270 guage modeling. In *Proceedings of the 59th Annual Meeting of the Association for Com-*  
271 *putational Linguistics and the 11th International Joint Conference on Natural Language*  
272 *Processing (Volume 1: Long Papers)*, pages 4897–4908, Online, August 2021. Associa-  
273 tion for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.379. URL <https://aclanthology.org/2021.acl-long.379>.
- 275 [23] Xiang Hu, Haitao Mi, Liang Li, and Gerard de Melo. Fast-R2D2: A pretrained recursive neural  
276 network based on pruned CKY for grammar induction and text representation. In *Proceedings of*  
277 *the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2809–2821,  
278 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.  
279 URL <https://aclanthology.org/2022.emnlp-main.181>.
- 280 [24] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. Transformer quality in linear time. In  
281 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato,  
282 editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
283 *Proceedings of Machine Learning Research*, pages 9099–9117. PMLR, 17–23 Jul 2022. URL  
284 <https://proceedings.mlr.press/v162/hua22a.html>.
- 285 [25] Shankar Iyer, Nikhil Dandekar, and Kornél Csernai. First quora dataset release: Question pairs.  
286 In *Quora*, 2017.
- 287 [26] Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. Learning the difference that makes a  
288 difference with counterfactually-augmented data. In *International Conference on Learning*  
289 *Representations*, 2020. URL <https://openreview.net/forum?id=Sk1gs0NFvr>.
- 290 [27] Donald E. Knuth. On the translation of languages from left to right. *Information and Control*, 8  
291 (6):607 – 639, 1965. ISSN 0019-9958.
- 292 [28] Phong Le and Willem Zuidema. Compositional distributional semantics with long short term  
293 memory. In *Proceedings of the Fourth Joint Conference on Lexical and Computational Seman-*  
294 *tics*, pages 10–19, Denver, Colorado, June 2015. Association for Computational Linguistics.  
295 doi: 10.18653/v1/S15-1002. URL <https://aclanthology.org/S15-1002>.
- 296 [29] Phong Le and Willem Zuidema. The forest convolutional network: Compositional distributional  
297 semantics with a neural chart and without binarization. In *Proceedings of the 2015 Conference*  
298 *on Empirical Methods in Natural Language Processing*, pages 1155–1164, Lisbon, Portugal,  
299 September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1137. URL  
300 <https://aclanthology.org/D15-1137>.
- 301 [30] Ji Ma, Jingbo Zhu, Tong Xiao, and Nan Yang. Easy-first POS tagging and dependency  
302 parsing with beam search. In *Proceedings of the 51st Annual Meeting of the Association*  
303 *for Computational Linguistics (Volume 2: Short Papers)*, pages 110–114, Sofia, Bulgaria,  
304 August 2013. Association for Computational Linguistics. URL <https://aclanthology.org/P13-2020>.
- 305

- 306 [31] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher  
307 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting*  
308 *of the Association for Computational Linguistics: Human Language Technologies*, pages 142–  
309 150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL  
310 <https://aclanthology.org/P11-1015>.
- 311 [32] Jean Maillard and Stephen Clark. Latent tree learning with differentiable parsers: Shift-reduce  
312 parsing and chart parsing. In *Proceedings of the Workshop on the Relevance of Linguistic*  
313 *Structure in Neural Architectures for NLP*, pages 13–18, Melbourne, Australia, July 2018.  
314 Association for Computational Linguistics. doi: 10.18653/v1/W18-2903. URL [https://](https://aclanthology.org/W18-2903)  
315 [aclanthology.org/W18-2903](https://aclanthology.org/W18-2903).
- 316 [33] Jean Maillard, Stephen Clark, and Dani Yogatama. Jointly learning sentence embeddings and  
317 syntax with unsupervised tree-lstms. *Natural Language Engineering*, 25(4):433–449, 2019. doi:  
318 10.1017/S1351324919000184.
- 319 [34] Tsendsuren Munkhdalai and Hong Yu. Neural tree indexers for text understanding. In *Pro-*  
320 *ceedings of the 15th Conference of the European Chapter of the Association for Computational*  
321 *Linguistics: Volume 1, Long Papers*, pages 11–21, Valencia, Spain, April 2017. Association for  
322 Computational Linguistics. URL <https://aclanthology.org/E17-1002>.
- 323 [35] Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig.  
324 Stress test evaluation for natural language inference. In *Proceedings of the 27th International*  
325 *Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA,  
326 August 2018. Association for Computational Linguistics. URL [https://www.aclweb.org/](https://www.aclweb.org/anthology/C18-1198)  
327 [anthology/C18-1198](https://www.aclweb.org/anthology/C18-1198).
- 328 [36] Nikita Nangia and Samuel Bowman. ListOps: A diagnostic dataset for latent tree learning.  
329 In *Proceedings of the 2018 Conference of the North American Chapter of the Association for*  
330 *Computational Linguistics: Student Research Workshop*, pages 92–99, New Orleans, Louisiana,  
331 USA, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-4013. URL  
332 <https://aclanthology.org/N18-4013>.
- 333 [37] Xuan-Phi Nguyen, Shafiq Joty, Steven Hoi, and Richard Socher. Tree-structured attention with  
334 hierarchical accumulation. In *International Conference on Learning Representations*, 2020.  
335 URL <https://openreview.net/forum?id=HJxK5pEYvr>.
- 336 [38] Hao Peng, Sam Thomson, and Noah A. Smith. Backpropagating through structured argmax  
337 using a SPIGOT. In *Proceedings of the 56th Annual Meeting of the Association for Com-*  
338 *putational Linguistics (Volume 1: Long Papers)*, pages 1863–1873, Melbourne, Australia,  
339 July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1173. URL  
340 <https://aclanthology.org/P18-1173>.
- 341 [39] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for  
342 word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural*  
343 *Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for  
344 Computational Linguistics. doi: 10.3115/v1/D14-1162. URL [https://aclanthology.org/](https://aclanthology.org/D14-1162)  
345 [D14-1162](https://aclanthology.org/D14-1162).
- 346 [40] Jordan B. Pollack. Recursive distributed representations. *Artificial Intelligence*, 46(1):77 –  
347 105, 1990. ISSN 0004-3702. doi: [https://doi.org/10.1016/0004-3702\(90\)90005-K](https://doi.org/10.1016/0004-3702(90)90005-K). URL  
348 <http://www.sciencedirect.com/science/article/pii/000437029090005K>.
- 349 [41] Swarnadeep Saha, Yixin Nie, and Mohit Bansal. ConjNLI: Natural language inference over  
350 conjunctive sentences. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*  
351 *Language Processing (EMNLP)*, pages 8240–8252, Online, November 2020. Association  
352 for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.661. URL [https://](https://aclanthology.org/2020.emnlp-main.661)  
353 [aclanthology.org/2020.emnlp-main.661](https://aclanthology.org/2020.emnlp-main.661).
- 354 [42] M.P. Schützenberger. On context-free languages and push-down automata. *Information and*  
355 *Control*, 6(3):246 – 264, 1963. ISSN 0019-9958.

- 356 [43] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position rep-  
357 resentations. In *Proceedings of the 2018 Conference of the North American Chapter of the*  
358 *Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short*  
359 *Papers)*, pages 464–468, New Orleans, Louisiana, June 2018. Association for Computational  
360 Linguistics. doi: 10.18653/v1/N18-2074. URL <https://aclanthology.org/N18-2074>.
- 361 [44] Yikang Shen, Shawn Tan, Arian Hosseini, Zhouhan Lin, Alessandro Sordoni, and Aaron C  
362 Courville. Ordered memory. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-  
363 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,  
364 pages 5037–5048. Curran Associates, Inc., 2019. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8748-ordered-memory.pdf)  
365 [8748-ordered-memory.pdf](http://papers.nips.cc/paper/8748-ordered-memory.pdf).
- 366 [45] Yikang Shen, Yi Tay, Che Zheng, Dara Bahri, Donald Metzler, and Aaron Courville. Struct-  
367 Former: Joint unsupervised induction of dependency and constituency structure from masked  
368 language modeling. In *Proceedings of the 59th Annual Meeting of the Association for Com-*  
369 *putational Linguistics and the 11th International Joint Conference on Natural Language*  
370 *Processing (Volume 1: Long Papers)*, pages 7196–7209, Online, August 2021. Associa-  
371 tion for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.559. URL [https:](https://aclanthology.org/2021.acl-long.559)  
372 <https://aclanthology.org/2021.acl-long.559>.
- 373 [46] Zhiqiang Shen, Zechun Liu, and Eric Xing. Sliced recursive transformer. In *Computer Vision–*  
374 *ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings,*  
375 *Part XXIV*, pages 727–744. Springer, 2022.
- 376 [47] Haoyue Shi, Hao Zhou, Jiase Chen, and Lei Li. On tree-based neural sentence modeling. In  
377 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*,  
378 pages 4631–4641, Brussels, Belgium, October–November 2018. Association for Computational  
379 Linguistics. doi: 10.18653/v1/D18-1492. URL <https://aclanthology.org/D18-1492>.
- 380 [48] Richard Socher, Christopher D. Manning, and Andrew Y. Ng. Learning continuous phrase  
381 representations and syntactic parsing with recursive neural networks. In *In Proceedings of the*  
382 *NIPS-2010 Deep Learning and Unsupervised Feature Learning Workshop*, 2010.
- 383 [49] Richard Socher, Jeffrey Pennington, Eric H. Huang, Andrew Y. Ng, and Christopher D. Manning.  
384 Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proceedings*  
385 *of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 151–  
386 161, Edinburgh, Scotland, UK., July 2011. Association for Computational Linguistics. URL  
387 <https://aclanthology.org/D11-1014>.
- 388 [50] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew  
389 Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a  
390 sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural*  
391 *Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association  
392 for Computational Linguistics. URL <https://aclanthology.org/D13-1170>.
- 393 [51] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representa-  
394 tions from tree-structured long short-term memory networks. In *Proceedings of the 53rd Annual*  
395 *Meeting of the Association for Computational Linguistics and the 7th International Joint Con-*  
396 *ference on Natural Language Processing (Volume 1: Long Papers)*, pages 1556–1566, Beijing,  
397 China, July 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1150. URL  
398 <https://aclanthology.org/P15-1150>.
- 399 [52] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng  
400 Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for  
401 efficient transformers. In *International Conference on Learning Representations*, 2021. URL  
402 <https://openreview.net/forum?id=qVyeW-grC2k>.
- 403 [53] Zhiyang Teng and Yue Zhang. Head-lexicalized bidirectional tree LSTMs. *Transactions of the*  
404 *Association for Computational Linguistics*, 5:163–177, 2017. doi: 10.1162/tacl\_a\_00053. URL  
405 <https://aclanthology.org/Q17-1012>.

- 406 [54] Ke Tran, Arianna Bisazza, and Christof Monz. The importance of being recurrent for  
407 modeling hierarchical structure. In *Proceedings of the 2018 Conference on Empirical  
408 Methods in Natural Language Processing*, pages 4731–4736, Brussels, Belgium, October-  
409 November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1503. URL  
410 <https://aclanthology.org/D18-1503>.
- 411 [55] Yaushian Wang, Hung-Yi Lee, and Yun-Nung Chen. Tree transformer: Integrating tree  
412 structures into self-attention. In *Proceedings of the 2019 Conference on Empirical Meth-  
413 ods in Natural Language Processing and the 9th International Joint Conference on Natural  
414 Language Processing (EMNLP-IJCNLP)*, pages 1061–1070, Hong Kong, China, Novem-  
415 ber 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1098. URL  
416 <https://aclanthology.org/D19-1098>.
- 417 [56] Zhiguo Wang, Wael Hamza, and Radu Florian. Bilateral multi-perspective matching for natural  
418 language sentences. In *Proceedings of the 26th International Joint Conference on Artificial  
419 Intelligence, IJCAI’17*, page 4144–4150. AAAI Press, 2017. ISBN 9780999241103.
- 420 [57] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus  
421 for sentence understanding through inference. In *Proceedings of the 2018 Conference of the  
422 North American Chapter of the Association for Computational Linguistics: Human Language  
423 Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana, June  
424 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL <https://aclanthology.org/N18-1101>.
- 425 [58] Zhaofeng Wu. Learning with latent structures in natural language processing: A survey. *ArXiv*,  
426 abs/2201.00490, 2022.
- 428 [59] Zihao Ye, Qipeng Guo, Quan Gan, Xipeng Qiu, and Zheng Zhang. Bp-transformer: Modelling  
429 long-range context via binary partitioning. *arXiv preprint arXiv:1911.04070*, 2019.
- 430 [60] Dani Yogatama, Phil Blunsom, Chris Dyer, Edward Grefenstette, and Wang Ling. Learning to  
431 compose words into sentences with reinforcement learning. In *5th International Conference on  
432 Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track  
433 Proceedings*. OpenReview.net, 2017.
- 434 [61] Aston Zhang, Yi Tay, Yikang Shen, Alvin Chan, and SHUAI ZHANG. Self-instantiated  
435 recurrent units with dynamic soft recursion. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S.  
436 Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*,  
437 volume 34, pages 6503–6514. Curran Associates, Inc., 2021. URL [https://proceedings.  
438 neurips.cc/paper/2021/file/3341f6f048384ec73a7ba2e77d2db48b-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/3341f6f048384ec73a7ba2e77d2db48b-Paper.pdf).
- 439 [62] Yuan Zhang, Jason Baldridge, and Luheng He. PAWS: Paraphrase adversaries from word  
440 scrambling. In *Proceedings of the 2019 Conference of the North American Chapter of the  
441 Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long  
442 and Short Papers)*, pages 1298–1308, Minneapolis, Minnesota, June 2019. Association for  
443 Computational Linguistics. doi: 10.18653/v1/N19-1131. URL [https://aclanthology.  
444 org/N19-1131](https://aclanthology.org/N19-1131).
- 445 [63] Xiaodan Zhu, Parinaz Sobihani, and Hongyu Guo. Long short-term memory over recursive  
446 structures. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International  
447 Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*,  
448 pages 1604–1612, Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.  
449 press/v37/zhub15.html](https://proceedings.mlr.press/v37/zhub15.html).
- 450 [64] Xiaodan Zhu, Parinaz Sobhani, and Hongyu Guo. DAG-structured long short-term memory  
451 for semantic compositionality. In *Proceedings of the 2016 Conference of the North American  
452 Chapter of the Association for Computational Linguistics: Human Language Technologies*,  
453 pages 917–926, San Diego, California, June 2016. Association for Computational Linguistics.  
454 doi: 10.18653/v1/N16-1106. URL <https://aclanthology.org/N16-1106>.