

# A Semi-Autoregressive Graph Generative Model for Dependency Parsing

Anonymous ACL submission

## Abstract

Recent years have witnessed impressive progress in Neural Dependency Parsing. According to the different factorization approaches to the graph joint probabilities, existing parsers can be roughly divided into autoregressive and non-autoregressive patterns. The former means that the graph should be factorized into multiple sequentially dependent components, then it can be built up component by component. And the latter assumes these components to be independent so that they can be outputted at once. However, when treating the directed edge in the dependency graph as an explicit dependency, we discover that there is a mixture of independent and interdependent components in the dependency graph, signifying that both fail to precisely capture the explicit dependencies among nodes and edges. Based on this property, we design a Semi-Autoregressive Dependency Parser to generate dependency graphs via adding node groups and edge groups autoregressively while pouring out all group elements in parallel. The model meanwhile deals with two problems in graph generation with respect to the uncertainty of generation orders and edge sparsity, via introducing a novel concept of Topological Hierarchy and a Graph Transformer as the decoder. The experiments show the proposed parser outperforms strong baselines on Enhanced Universal Dependencies of 14 languages. Also, the performances of model variations show the importance of specific parts.

## 1 Introduction

Dependency graph in neural parsing is a directed graph representing semantic dependencies between words, with a transitive relation traveling from the rooted node to all words in the sentence phase by phase. As such, transition-based parsing seems to be a natural choice, as it builds up the parsing graph sequentially so that the dependency relationships can be captured. However, graph-based parsing

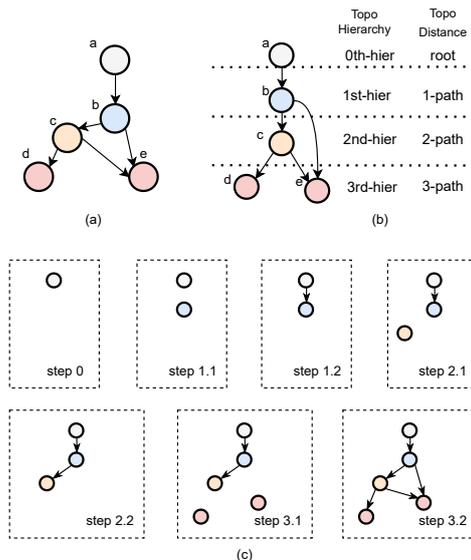


Figure 1: (a) An example graph (b) Divide nodes into different topological hierarchies based on their furthest distances from the root node. (c) Semi-autoregressive graph generation process.

dominates recent competitions on parsing technologies including IWPT 2020 and 2021 (Bouma et al., 2020, 2021), even if using a simple biaffine attention (Dozat and Manning, 2016) only to predict the whole graph at once. To explore a more effective parsing method that can represent these dependency relationships in a rigorous manner, we define and construct Topological Hierarchies for dependency graphs based on the explicit dependencies carried by them. According to the characteristics of topological hierarchies, we propose a Semi-Autoregressive Dependency Parser (SADP) – a novel graph-based parsing fashion via the semi-autoregressive graph generation.

Generally, autoregressive graph generation indicates that the model dynamically adds nodes and edges based on the generated sub-graph structure until reaching the complete graph. Its first challenge is to determine an generation order so that the joint probability of the graph can be factorized

063 into the product of conditional probabilities. Al- 114  
064 though the dependency graphs stipulate the strict 115  
065 sequential dependencies by directed edges, there is 116  
066 a lack of such topological orders between sibling 117  
067 nodes. For instance, node *b* in Figure 1.a depends 118  
068 on the node *a* because there is an explicit edge 119  
069 pointing from *a* to *b*. However, it is hard to decide 120  
070 the dependency relationship between node *d* and 121  
071 node *e* as they are not linked directly or indirectly. 122  
072 Previous works on directed graph generation solve 123  
073 the problem surfacely. Cai and Lam (2019, 2020) 124  
074 sort these sibling nodes randomly at the early stage 125  
075 of training and then change them to a determinis- 126  
076 tic order (e.g., relation-frequency) at later training 127  
077 steps. Some other works do the sorting by referring 128  
078 to the orders of the known sequences like word or- 129  
079 der or alphanumerical order (Zhang et al., 2019a,b; 130  
080 Bevilacqua et al., 2021). Since the dependency 131  
081 graph does not assign an explicit sequential rela- 132  
082 tionship between sibling nodes, such imposed or- 133  
083 ders would lead to exposure bias (Ranzato et al., 134  
084 2016) between training and inference. Once the 135  
085 sibling nodes are not generated in the same order 136  
086 as in the training, the learned knowledge would be 137  
087 invalid and even mislead subsequent predictions. 138  
088 The aforementioned random ordering seems to alle- 139  
089 viate the problem to some extent, but it destabilizes 140  
090 and complicates the training process and generally 141  
091 results in inferior models. 142

092 Instead of imposing orders on these sibling 134  
093 nodes, we assume them (including their incoming 135  
094 edges) to be conditionally independent to construct 136  
095 Topological Hierarchies (TH) as the generation or- 137  
096 ders. As shown in Figure 1.b, we divide nodes 138  
097 into several hierarchies according to their furthest 139  
098 distances from the root node. We can see that there 140  
099 are no explicit dependency relationships between 141  
100 nodes in the same hierarchy. Besides, nodes in 142  
101 the later hierarchies only depends on those in the 143  
102 previous hierarchies, forming a natural generation 144  
103 sequence. For a directed acyclic graph (DAG), it 145  
104 *at least* has one topological ordering but *only* has 146  
105 one topological hierarchy. At each generation step, 147  
106 we firstly predict all new nodes in parallel and then 148  
107 calculate their incoming edges by the biaffine at- 149  
108 tention (Dozat and Manning, 2016). In a word, 150  
109 our model autoregressively adds node groups and 151  
110 edge groups but non-autoregressively generates el- 152  
111 ements in these groups. See Figure 1.c for our 153  
112 semi-autoregressive generation process. 154

113 Another challenge is that incorrect sub-graph

114 structures may be predicted during inference. Tradi- 115  
116 tional graph representation models like GCN (Kipf 117  
118 and Welling, 2016) heavily rely on the given adja- 119  
120 cency to capture context information. That means 121  
122 it may fail to represent historical information com- 123  
124 pletely and efficiently when predicted edges make 124  
125 mistakes. An extreme situation of edge sparsity is 125  
126 that the new nodes have no incoming edges pre- 126  
127 dicted so that the model can only represent its node 127  
128 features rather than the sub-graph structure. To 128  
129 enhance the robustness of the generator, we design 129  
130 a novel graph representation model deriving from 130  
131 Transformer-decoder (Vaswani et al., 2017). In 131  
132 our Graph Transformer-decoder, there are implicit 132  
133 edges linking from the nodes in the previous and 133  
134 current hierarchies to the new node. Then, the pre- 134  
135 dicted explicit edges serve as the bias to adjust the 135  
136 attention distribution over the implicit edges so that 136  
137 the model can adaptively select useful structural 137  
138 information. 138

139 Overall, this paper proposes a novel direction 134  
140 – semi-autoregression to deal with parsing prob- 135  
141 lems, distinguished with autoregression and non- 136  
142 autoregression (detailed definitions about them are 137  
143 available in § 2). With the dependencies denoted as 138  
144 the directed edges, the semi-autoregressive pattern 139  
145 unfolds graphs in the ordering of topological hier- 140  
146 archies, which strictly follows the explicit depen- 141  
147 dency relationships defined in dependency graphs. 142  
148 Besides, it alleviates exposure bias in the genera- 143  
149 tion orders as Independent elements are orderless, 144  
150 which promotes models in both quality and effi- 145  
151 ciency. On the other hand, graph transformer has 146  
152 achieved significant progress in the field of classifi- 147  
153 cation (Ying et al., 2021), but rare studies explore 148  
154 its applications in the generation. This paper de- 149  
155 signs a novel graph transformer and adapts it to the 150  
156 semi-autoregressive graph generation to alleviate 151  
157 the edge sparsity problem. On the experimental 152  
158 side, we evaluate SADP on Enhanced Universal 153  
159 Dependencies (EUD) which are non-tree depen- 154  
160 dency graphs. In addition to the official evaluation 155  
161 metric Enhanced Label Attachment Scores (ELAS), 156  
162 we design a graph-level matching score (GMS) to 157  
163 assess the probability of returning an absolutely 158  
164 correct graph. The results show that our model 159  
160 outperforms other baselines significantly. Finally, 160  
161 we introduce multiple model variations to investi- 161  
162 gate the effect of different model components and 162  
163 show that our model is well-designed, especially 163  
164 the parts of discarding imposed orders and adding 164

165 implicit edges.

## 166 2 Related Work

167 **Autoregressive Parser.** Generally, a generator is  
168 in an autoregressive fashion provided its generation  
169 probability at each step is conditional on items it  
170 produces previously. Transition-based parser ob-  
171 viously conforms to the characteristic, as it up-  
172 dates the action probability every step based on  
173 the words, tags and label embeddings previously  
174 put in the buffer and stack (Chen and Manning,  
175 2014). Meanwhile, we note that some mechanisms  
176 commonly used in autoregressive generators are  
177 used to improve transition-based parsers like beam  
178 search and pointer networks (Weiss et al., 2015;  
179 Ma et al., 2018; Fernández-González et al., 2019).  
180 On the other hand, Cheng et al. (2016) proposes a  
181 graph-based autoregressive parser by adding arcs  
182 sequentially with the considerations of previous  
183 parsing decisions. However, it should not be taken  
184 as a rigorous graph generative model as it does not  
185 generate by extending the sub-graph structures. Ac-  
186 tually, instead of dependency graphs, it is more  
187 prevalent that leverage the autoregressive graph  
188 generators to parse Abstract Meaning Representa-  
189 tion (AMR) (Cai and Lam, 2019, 2020; Zhang  
190 et al., 2019b,a). They are all in the (fully) autore-  
191 gressive pattern that an order is imposed to nodes  
192 and edges without topological orderings. In this  
193 paper, we investigate the effects of these imposed  
194 orderings by introducing some variations of the  
195 proposed model. Further studies on AMR will be  
196 available in our future work.

197 **Non-Autoregressive Parser.** In contrast, non-  
198 autoregression implies that all components factor-  
199 ized from the graph are independent, so their prob-  
200 abilities do not affect each other and can be ob-  
201 tained in parallel at any time. A representative  
202 non-autoregressive parser is Deep Biaffine Atten-  
203 tion (BiAtt) (Dozat and Manning, 2016) which  
204 assuming all edges are independent. For the tree-  
205 structure dependency graphs, it is often followed by  
206 a searching algorithm for the Maximum Spanning  
207 Tree (MST). Some heuristic algorithms (Li et al.,  
208 2020; Kiperwasser and Goldberg, 2016) construct  
209 the MST step by step, which yet does not mean they  
210 are in the autoregressive manner because all edge  
211 probabilities are predicted at once and fixed be-  
212 fore the searching. Another confusing models are  
213 higher-order graph-based parsers. Among them, Ji  
214 et al. (2019) incorporates the second-order knowl-

edge into the word representations and still uses the  
BiAtt as the final parser. Wang et al. (2019); Zhang  
et al. (2020) decompose the graph into components  
of different second-order parts. Different from Bi-  
Att that each component is an edge, here some  
components consists of two edges whose joint prob-  
abilities can be calculated as a whole by a trilinear  
function. They still belong to non-autoregressive  
parsers because their components are independent  
of each other and disable to be subdivided.

## 3 Proposed Model

### 3.1 Definitions

**Problem Definition.** Conditional on the source  
sentence  $S = (w_n)_{n=1}^N$ , the task is to generate a  
dependency graph hierarchy by hierarchy. The gen-  
eration process can be denoted as a sequence of  
components:  $(C^{(t)})_{t=0}^T, T \leq N$ . We firstly turn  
dependency graphs to DAGs by deleting the back  
edges in their cycles. It should be mentioned that  
there are only a few graphs with cycles and we  
can add these removed edges back by rules before  
evaluations. Then we can construct Topological  
Hierarchies based on the furthest distance from  
each node to the root node. The initial compo-  
nent  $C^{(0)} = \{v_0\}$  in the 0-th hierarchy only has  
a root node. When  $t > 0$ , the component in the  
 $t$ -th hierarchy is defined as  $C^{(t)} = \{V^{(t)}, E^{(t)}\}$ .  
Let  $\mathcal{V}_t = \bigcup_{j=0}^t V^{(j)}$ , then  $V^{(t)} = \{v_i\}_{i=|\mathcal{V}_{t-1}|}^{|\mathcal{V}_t|-1}$   
is a set of nodes in the  $t$ -th hierarchy. And,  
 $E^{(t)} = \{(v_j, v_i, z_{ji}) | v_j \in \mathcal{V}_{t-1}, v_i \in V^{(t)}\}$  is a  
set of edges pointing from nodes in the previous  
hierarchies to the current nodes, where  $v_j$  is the  
head of  $v_i$  and  $z_{ji}$  is the label on the arc.

**Explicit and Implicit Edge.** We define two  
kinds of edges, namely *explicit* edges and *implicit*  
edges. The former is what we need to really predict.  
Let  $\mathcal{N}_i$  be the explicit first-order neighbours of the  
node  $v_i \in V^{(t)}$  and  $\mathcal{D}_i$  be the implicit neighbours,  
and  $\mathcal{N}_i \cup \mathcal{D}_i = \mathcal{V}_t$ . Notably, nodes in  $\mathcal{N}_i$  can  
not appear in  $V^{(t)}$  according to the definition of  
topological hierarchy. They have uni-directional  
edges pointing to the node  $v_i$  with arc labels, and  
these edge can be found in  $E^{(t)}$ . On the other  
hand, nodes in  $\mathcal{D}_i$  should not have pointed to  $v_i$ ,  
but our model does so because we expect nodes to  
learn structural information adaptively. It should  
be mentioned that nodes in the same component or  
hierarchy also have implicit edges linking to each  
other, i.e.,  $V^{(t)} \subseteq \mathcal{D}_i$ .

**Head and Dependent Representation.** We de-

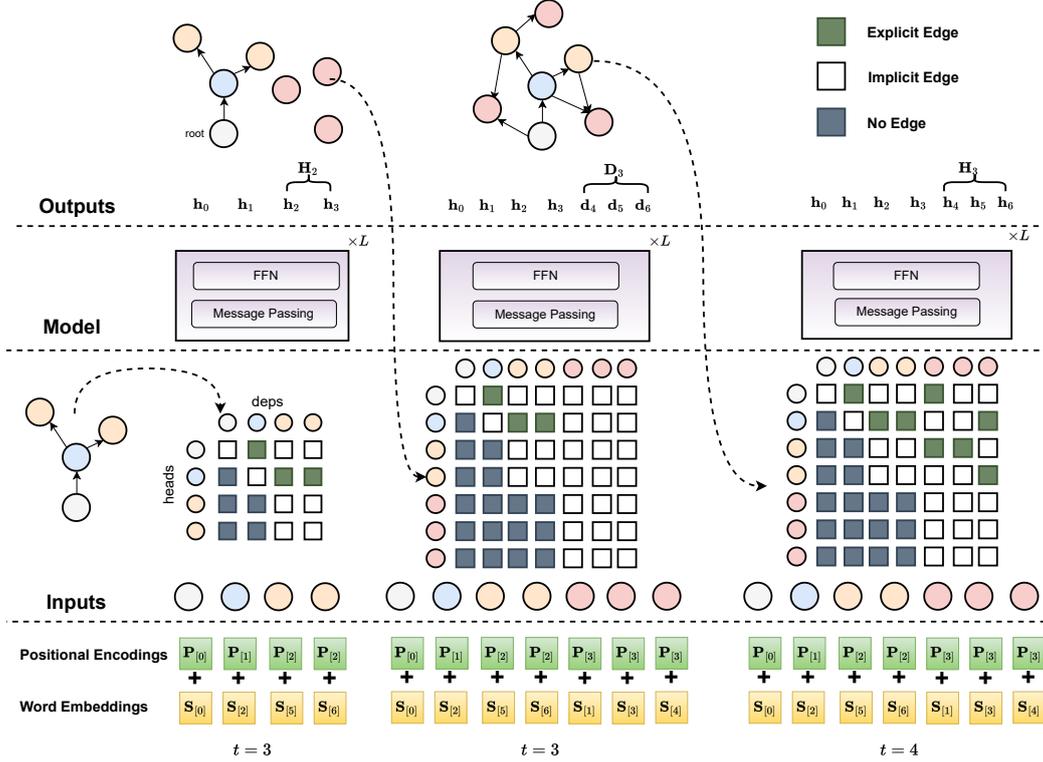


Figure 2: Semi-autoregressive generation process and graph transformer.

fine two representations of the same node with different roles, namely the *head* representation and the *dependent* representation. Each generated node will first be used as a dependent node to calculate its incoming arcs, and then as a head node until the end of generation. We define the head vector of a node  $v_i \in V^{(t)}$  as  $\mathbf{h}_i$  and its dependent vector as  $\mathbf{d}_i$ . For a component, its head matrix  $\mathbf{H}^{(t)} = \mathcal{F}_\theta(\mathcal{V}_t, \mathcal{E}_t, S)$  and dependent matrix  $\mathbf{D}^{(t)} = \mathcal{F}_\theta(\mathcal{V}_t, \mathcal{E}_{t-1}, S)$  are the concatenations of multiple corresponding node representation, where  $\mathcal{E}_t = \bigcup_{j=0}^t E^{(j)}$ . We can see that the difference between them is that the latter inputs lack  $E^{(t)}$ , which means there are no available explicit edges pointing to  $V^{(t)}$  nodes when calculating dependent representations. It should be mentioned that the graph representation model  $\mathcal{F}_\theta(\cdot)$  can represent all components, but we only need to focus on the new component at each generation step because the new component does not affect node representations in the previous components.

**Training Objective.** The objective is to maximize the graph joint probability  $\mathcal{J}$ :

$$\mathcal{J} = \prod_{t=1}^T P(V^{(t)} | \mathcal{V}_{t-1}, \mathcal{E}_{t-1}) P(E^{(t)} | \mathcal{V}_t, \mathcal{E}_{t-1}) \quad (1)$$

$$P(V^{(t)} | \mathcal{V}_{t-1}, \mathcal{E}_{t-1}) = \prod_{v_i \in V^{(t)}} P(v_i | \mathcal{V}_{t-1}, \mathcal{E}_{t-1}) \quad (2)$$

$$P(E^{(t)} | \mathcal{V}_t, \mathcal{E}_{t-1}) = \prod_{e_i \in E^{(t)}} P(e_i | \mathcal{V}_t, \mathcal{E}_{t-1}) \quad (3)$$

$$\mathcal{V}_t = \mathcal{V}_{t-1} \cup V^{(t)}, \quad \mathcal{E}_t = \mathcal{E}_{t-1} \cup E^{(t)} \quad (4)$$

It indicates that we autoregressively generate the new node group  $V^{(t)}$  and the edge group  $E^{(t)}$  based on groups generated previously and the elements in the same group are independent.

### 3.2 Graph Generation Process

Figure 2 presents the generative process from the 3-rd step to the 4-th step. Specifically, At the generation step  $t$ , we firstly update head representations  $\mathbf{H}^{(t-1)}$  for the last-step nodes  $V^{(t-1)}$  using their network structure information  $E^{(t-1)}$ . Notably, although there only generates an intermediate sub-graph of the entire structure, the explicit topological information of  $V^{(t-1)}$  nodes is completed because they would not have incoming arcs from nodes generated later. On the other hand, these

310 sentence words have been represented as a dense  
 311 matrix  $\mathbf{S}$  by a Transformer-encoder. Then, their  
 312 probabilities of being selected are calculated by:

$$313 \quad P(w_{1:N}) = \text{MaxPool} \left[ \sigma \left( \mathbf{H}^{(t-1)} \mathbf{W}_1 \mathbf{W}_2^\top \mathbf{S}^\top \right) \right] \quad (5)$$

$$314 \quad V^{(t)} = \{w_n | P(w_n) > 0.5\} \quad (6)$$

315 where  $\mathbf{W} \in \mathbb{R}^{d \times d}$  is a linear transformation matrix.  
 316 This operation is similar to a multi-label classifi-  
 317 cation. Every source word is assigned with an  
 318 independent probability, and words with probabili-  
 319 ties larger than 0.5 are selected as new nodes  $V^{(t)}$ .  
 320 To represent these new nodes as  $\mathbf{D}^{(t)}$  when their  
 321 network structural information are unknown, we  
 322 suppose that there are implicit edges pointing from  
 323 previous nodes to these nodes. Besides, these new  
 324 nodes are connected to each other by implicit edges.  
 325 Although it is impossible to appear explicit edges  
 326 among them, this operation can further enrich node  
 327 representations. Their connections are illustrated  
 328 by the second adjacency matrix in the middle block  
 329 of Figure 2. Explicit edge connections and types  
 330 are then figured out by Deep Biaffine Attention  
 331 (Dozat and Manning, 2016):  
 332

$$333 \quad E^{(t)} = \text{DeepBiaffine} \left( \left\|_{j=0}^{t-1} \mathbf{H}^{(j)} \right\|, \mathbf{D}^{(t)} \right) \quad (7)$$

334 where  $\left\|_{j=0}^{t-1} \mathbf{H}^{(j)} \right\|$  is achieved by concatenating head  
 335 representations of all nodes in the previous hierar-  
 336 chies. The generation proceeds via repeating the  
 337 aforementioned operations until no words can be  
 338 selected as new nodes.

### 339 3.3 Graph Representation Model

340 Recently, Transformer (Vaswani et al., 2017) has  
 341 made impressive progress in the graph repre-  
 342 sentation field (Ying et al., 2021). In essence,  
 343 Transformer regards inputs as an undirected fully-  
 344 connected graph, thus serving as a special graph  
 345 representation model that can enjoys global percep-  
 346 tion at all layers. Previous works focusing on adapt-  
 347 ing Transformer-encoder to node or graph classifi-  
 348 cation, while this paper modifies Transformer-  
 349 decoder to conduct graph generation.

350 Let  $\mathbf{x}_i^{(l)}$  denote the node  $v_i$  embedding at the  $l$ -th  
 351 layer. If the node  $v_i$  is in the component  $C_t$  and  
 352 copied from the source word  $w_n$ , its initial node  
 353 embedding  $\mathbf{x}_i^{(0)}$  should be the summation of:

$$354 \quad \mathbf{x}_i^{(0)} = \mathbf{S}_{[n]} + \mathbf{P}_{[t]} \quad (8)$$

355 where  $\mathbf{S}, \mathbf{P}$  indicate word embeddings and hierar-  
 356 chical positional encodings respectively, as shown  
 357 in Figure 2. Nodes in the same hierarchy have the  
 358 same hierarchical positional encodings.

359 The message passing layer actually takes the  
 360 position of the masked self-attention layer in the  
 361 decoder. The original decoder self-attention helps  
 362 every word to aggregate left-ward contexts. In con-  
 363 trast, every node in our model not only aggregates  
 364 left-ward contexts (i.e, nodes in previous hierar-  
 365 chies), but also nodes in the same hierarchy. To  
 366 distinguish explicit edges and implicit edges, the  
 367 message vector  $\mathbf{m}_{ji}^{(l)}$  of the node  $v_j$  with an explicit  
 368 edge pointing to the node  $v_i$  should be enriched  
 369 with prior structural knowledge by:

$$370 \quad \mathbf{m}_{ji}^{(l)} = \begin{cases} \mathbf{x}_j^{(l)} + \text{relu} \left( \mathbf{x}_j^{(l)} \mathbf{U}_{z_{ji}} \right), & v_j \in \mathcal{N}_i \\ \mathbf{x}_j^{(l)}, & v_j \in \mathcal{D}_i \end{cases} \quad (9)$$

371 where  $\mathbf{U}_{z_{ji}} \in \mathbb{R}^{d \times d}$  indicates the parametric em-  
 372 bedding matrix of the edge label  $z_{ji}$ . These edge  
 373 embedding metrics are shared across all layers. No-  
 374 tably, we assume that the central node  $v_i$  is self-  
 375 connected implicitly, i.e.  $v_i \in \mathcal{D}_i$ . The reduction  
 376 function is then defined as the multi-head attention:

$$377 \quad \alpha_{ji} = \frac{\exp \left( \mathbf{x}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{m}_{ji}^\top \right)}{\sum_{v_u \in \mathcal{N}_i \cup \mathcal{D}_i} \exp \left( \mathbf{x}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{m}_{ui}^\top \right)} \quad (10)$$

$$378 \quad \mathbf{x}_i^{(l)'} = \left[ \begin{array}{c} H \\ \parallel \\ h=1 \end{array} \left( \sum_{v_j \in \mathcal{N}_i \cup \mathcal{D}_i} \alpha_{ji}^h \mathbf{m}_{ji}^{(l)} \mathbf{W}_V^h \right) \right] \mathbf{W}_O \quad (11)$$

379 We can see that the query is the node embedding  
 380  $\mathbf{x}_i$ , and the keys and values are those message vec-  
 381 tors  $\mathbf{m}_{ji}^{(l)}$ . Its output  $\mathbf{x}_i^{(l)'}$  is then fed into the feed-  
 382 forward layer to enter the next layer:

$$383 \quad \mathbf{x}_i^{(l+1)} = \text{FFN} \left( \mathbf{x}_i^{(l)'} \right) \quad (12)$$

384 The outputs  $\mathbf{x}^{(L)}$  of the final layer are head repre-  
 385 sentations or dependent representations.

386 The edge embedding matrices  $\mathbf{U}$  give the model  
 387 access to prior structural knowledge and enable it  
 388 to select useful prior knowledge adaptively. When  
 389 all structural knowledge is useless (i.e, parameters  
 390 in  $\mathbf{U}$  are trained to be zeros) and each hierarchy  
 391 only contains one node, the graph model degrades  
 392 to a vanilla Transformer decoder.

IWPT 2021		bg	cs	en	et	fi	fr	it	lt	lv	nl	pl	ru	sk	sv	uk	avg
ELAS	BiAtt	92.7	91.0	87.2	87.2	90.6	88.4	92.1	81.9	88.3	90.5	90.2	93.2	91.5	87.3	89.1	89.4
	Tree-Graph	92.8	<b>91.1</b>	87.3	87.1	<b>90.7</b>	88.6	92.3	81.9	88.2	90.5	<b>90.4</b>	93.2	91.6	87.5	89.0	89.5
	Ours	<b>92.9</b>	90.9	<b>87.9</b>	<b>87.3</b>	<b>90.7</b>	<b>89.5</b>	<b>92.8</b>	<b>83.5</b>	<b>88.5</b>	<b>90.9</b>	<b>90.4</b>	<b>93.5</b>	<b>92.1</b>	<b>87.9</b>	<b>89.6</b>	<b>89.9</b>
GMS	BiAtt	47.4	44.8	36.3	37.1	38.7	40.3	43.8	21.0	38.4	46.9	40.8	50.3	51.0	32.2	34.6	40.2
	Tree-Graph	47.8	45.3	36.9	37.0	39.1	41.1	44.7	21.0	37.9	47.0	41.9	50.8	51.5	33.4	34.2	40.6
	Ours	<b>48.8</b>	<b>45.6</b>	<b>40.3</b>	<b>39.2</b>	<b>41.4</b>	<b>45.4</b>	<b>47.1</b>	<b>28.2</b>	<b>42.8</b>	<b>51.3</b>	<b>43.6</b>	<b>54.2</b>	<b>57.8</b>	<b>38.2</b>	<b>39.2</b>	<b>44.2</b>

IWPT 2020		bg	cs	en	et	fi	fr	it	lt	lv	nl	pl	ru	sk	sv	uk	avg
ELAS	Second-order	91.5	90.1	87.1	86.0	89.0	85.3	91.5	78.9	87.6	86.2	84.0	92.3	87.6	84.7	88.0	87.3
	UDify	90.7	87.5	87.2	84.5	89.5	85.9	91.5	77.6	84.9	84.7	84.6	90.7	88.6	85.6	87.2	86.7
	Ours	<b>92.6</b>	<b>90.4</b>	<b>88.2</b>	<b>86.9</b>	<b>90.1</b>	<b>87.4</b>	<b>92.6</b>	<b>82.5</b>	<b>88.5</b>	<b>86.7</b>	<b>86.7</b>	<b>93.2</b>	<b>91.0</b>	<b>87.0</b>	<b>89.0</b>	<b>88.9</b>
GMS	Second-order	43.1	37.7	35.7	31.8	34.4	29.2	44.4	15.1	35.3	31.0	28.6	47.1	38.7	26.5	30.5	33.9
	UDify	41.4	31.4	34.1	31.2	34.5	33.2	41.5	17.8	31.6	23.8	26.1	40.6	43.1	27.1	31.2	32.6
	Ours	<b>48.3</b>	<b>43.5</b>	<b>41.6</b>	<b>36.6</b>	<b>38.4</b>	<b>38.7</b>	<b>47.1</b>	<b>25.6</b>	<b>42.0</b>	<b>34.1</b>	<b>32.9</b>	<b>53.7</b>	<b>55.2</b>	<b>36.4</b>	<b>38.3</b>	<b>40.8</b>

Table 1: Average ELAS and GMS results of 3 calculations on IWPT 2021 and IWPT 2020 datasets. We use  $L = 2$  according to ELAS on the English dev-set.

## 4 Experiment

### 4.1 Datasets

We tune our models primarily on 15 languages that appear in IWPT 2020 dataset and IWPT 2021 dataset (Bouma et al., 2020, 2021). The two shared tasks focus on EUD (Schuster and Manning, 2016) which are non-tree graphs with reentrancies, empty nodes and sparsity cycles. To construct the topological hierarchy, we need to delete the back edges in cycles firstly and add them back by rules at inference time. For the language that has multiple treebanks, we simply concatenate all of its treebanks. Besides, we use gold tokenization and gold sentence segmentation during training and development. At test time, we use the results of tokenization and segmentation provided by the top ranked models.

### 4.2 Baseline Models

Our comparison experiments aim to investigate the performances of models themselves, without considering some learning techniques like ensembling (Grünwald et al., 2021), two-stage training (Shi and Lee, 2021) and automated concatenation of embeddings (Wang et al., 2021). We conclude four strong baselines from top-ranked systems in IWPT 2021 and IWPT 2020, namely Deep Biaffine Attention (Dozat and Manning, 2016), Tree-Graph Parser (Shi and Lee, 2021), Second-order Parser (Wang et al., 2019, 2020) and Language-specific UDify (Kondratyuk and Straka, 2019; Kanerva et al., 2020). Their results are reported after eliminating the effects of learning techniques.

### 4.3 Evaluation Metrics

ELAS results are evaluated by the official script provided by IWPT 2021. Besides, we also define a graph-level matching score (GMS) to investigate whether the model can deal with a few arcs that are difficult to predict properly in a sample. Since we segment UD sentences from raw texts, the numbers of sentences are different for each system. Therefore, GMS is a  $F_1$  score around the number of absolutely matched graphs. The specific GMS calculation is put in App. 1.

### 4.4 Word Embeddings

Similar to the operations in most top-ranked systems, our word embeddings  $S_{[n]}$  are initialized as the weighted summation of the corresponding hidden states in XLM-R layers (Conneau et al., 2020), where the weights are the learned attention distribution over all XLM-R layers. For the word composed of multiple subwords, we extract the hidden states of the last one. We set up the dimension in the graph representation model as  $d = 1024$ , the same as that in the pre-trained models. See App. 2 for more experiment details.

## 5 Results and Analysis

### 5.1 Main Results

The official evaluation metrics ELAS of our models and baselines are shown in Table 1. We note that SADP achieves at least comparable results on all languages. In IWPT 2021, in addition to obtaining the best average ELAS performance (average  $\sim 0.4\%$  points), our model brings significant improvements over multiple languages like Lithuanian ( $\sim 1.6\%$  points), French ( $\sim 0.9\%$  points), English ( $\sim 0.6\%$  points). This enhancement is more

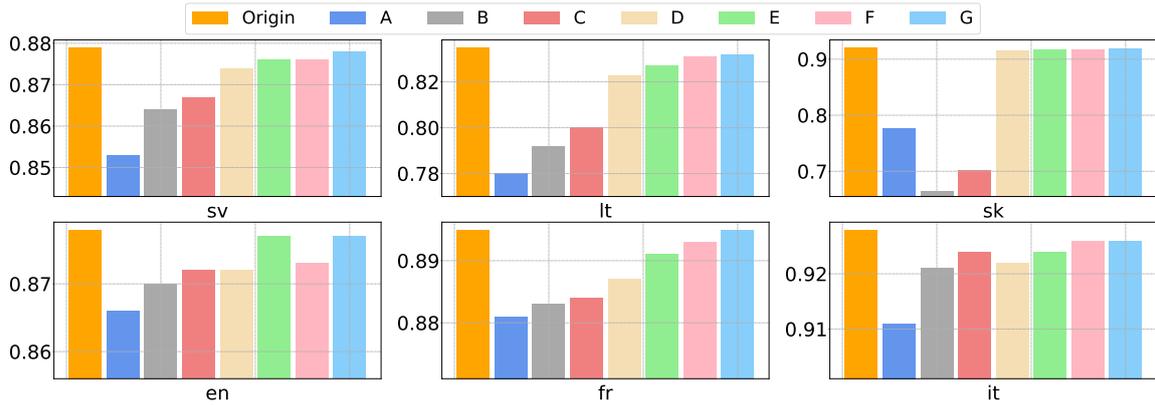


Figure 3: Test-set ELAS results, comparing the origin model with different model variations

significant when comparing our model with the top two models in IWPT 2020 (average  $\sim 1.6\%$  points). Besides, sharper increases appears in GMS of IWPT 2021 (average  $\sim 3.6\%$  points) and IWPT 2020 (average  $\sim 6.9\%$  points), where our model achieves an amazing rising against the baselines in all languages.

It should be mentioned that a higher ELAS does not mean a higher GMS, as shown in the results of Czech (*cs*) language. In other words, some obstinate errors are fixed to make more dependency graphs completely correct, but there appear some samples where more mistakes concentrate. This situation derives from the inherent characteristics of autoregressive generation that the prediction accuracy at one certain step is heavily dependent on that at historical steps. In ideal states, the historical information can calibrate some obstinate errors by the learned dependencies. However, once deviation occurs in an immediate step, it may lead to some mistakes that are too simple to make. This is the essential reason that autoregressive parsers are weaker than non-autoregressive parsers. By comparison, our semi-autoregressive parser mitigates the negative impact of this characteristic by removing some dependency relationships, thus resulting in better performances in both ELAS and GMS.

## 5.2 Model Variant Ablation Studies

To investigate the importance of different model components and input features, we evaluated the following variations of our model.

**A. Autoregressive generation with random orders.** We impose random orders to the sibling nodes, so the model is converted to a fully-autoregressive generator. At each step, the model only generates a new node and its all incoming

edges. The sibling nodes will be re-ordered after a training epoch.

**B. Autoregressive generation with word orders.** The sibling nodes are sorted by the the positions of the node words in the sentence.

**C. Combine random orders and word orders.** The sibling nodes are firstly randomly sorted at the early stage of training and fixed to the word orders at later training.

**D. No implicit edges.** Without the implicit edges, the graph representation model is similar to GAT (Veličković et al., 2017) but the messages are additionally enriched with the arc label information.

**E. No implicit edges in the same hierarchy.** We remove the implicit edges between nodes in the same topological hierarchy. In this case, each node only has the incoming arcs from the nodes in the previous hierarchies.

**F. No explicit edges.** We replace all explicit edges by implicit edges, which is equal to forcing the edge embedding matrix  $U$  to zeros.

**G. No hierarchical positional encodings.** In this case, the model would lose the sequential relationships between hierarchies and fail to locate nodes of different hierarchy.

The ablation results of 6 languages are summarized in Figure 3. We firstly focus on the fully-autoregressive variations, namely the model A, B and C. We can see that there are significant declines in performances when imposing orderings to sibling nodes, indicating that the autoregressive mode heavily suffers from exposure bias in terms of generation orderings. Besides, the extent of declines varies a lot in different languages, ranging from over 30% in the Slovak (*sk*) dataset and within 1% in the English (*en*) and Italian (*it*) datasets. This

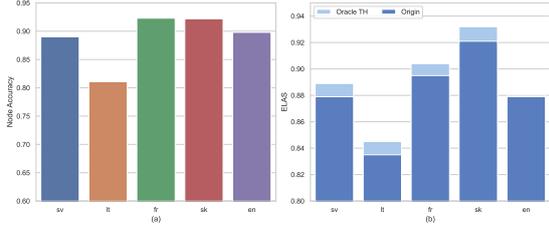


Figure 4: (a) Accuracy of nodes in the correct hierarchy. (b) ELAS results using Oracle Topological Hierarchy.

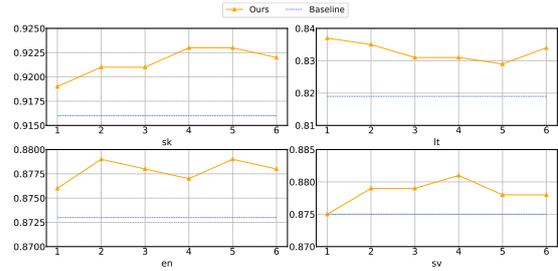


Figure 5: Sensitive analysis of layers on test sets.

proves that the impact of imposed sorting is quite unstable.

Moving to the model D, E and F which are variations with respect to explicit and implicit edges. Although it is not as significant as the negative effects of using autoregressive modes, that of removing implicit or explicit edges cannot be ignored. Generally, implicit edges play a more important role than explicit edges as the performances of model D are often lower than those of model E and F. This validates that the edge sparsity is the major problem of graph generation after the uncertainty of generation orderings. Besides, the implicit edges in the same hierarchy (see model E) and the historical arc label information (see model F) are both compulsory model components because the model always performs worse when dropping them. The final is the model without hierarchical positional encodings. Compared with other variations, its performance is the closest to the origin model, implying that our graph representation model is not very sensitive to the sequential relationships between hierarchies.

### 5.3 Error analysis of Topological Hierarchy

Since a Topological Hierarchy regulates the rough topological structure of a dependency graph, its prediction accuracy is crucial for the whole model. We investigate the node accuracy on 5 languages (see Figure 4.a), and find that about 90% nodes can fall into correct hierarchies. Even the language performing worst under the ELAS evaluation can reach 80% node accuracy. We then provide the model with the Oracle node group at each generation step and plot the comparison results against the origin model in Figure 4.b. There is about a 1% increase of ELAS on most languages when using Oracle Topological Hierarchies. It is surprising that Oracle TH does not bring about improvement to the English dataset, indicating that corrupt topological hierarchies do not always lead to incorrect arcs. Actually, it would cause bad results only when the

dependent node of an arc is generated before its head nodes. It does not matter for corruptions that do not shuffle the orders of heads and dependents.

### 5.4 Sensitive analysis of Layers

We test the sensitivity of ELAS results to different  $L$ . As shown in Figure 5, We select four languages whose ELAS are significantly higher than baseline' when  $L = 2$ . We find that our model still outperforms the baseline whichever  $L$  is used. Besides, it is hard to disclose a trend between model performance and the number of layers from the four plots. This is possibly because graph transformer can capture context information of high-order neighbours even with one layer. Overall, SADP is insensitive to the number of layers.

## 6 Conclusion and Limitation

This paper explores a semi-autoregressive dependency parser that learns the explicit dependencies in dependency graphs. This generation pattern captures the edge dependencies while reducing exposure bias, resulting in a more effective parser. Besides, the paper gives some insights into the two problems of graph generation, namely the ordering uncertainty and edge sparsity.

The limitations of our work fall into the inference speed and the decoding strategy. Semi-autoregressive inference speed is between non-autoregressive and autoregressive, and it is difficult to return graphs immediately like non-autoregressive parsers. Besides, this paper only introduces greedy search as the decoding strategy, which often performs worse than beam search as the latter provides a buffer for exposure bias. It is challenging for the semi-autoregressive beam search because it needs to select variable combinations with the highest probabilities instead of several single variables. We will include it in our future work.

610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666

## References

Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian McAuley. 2021. Rezero is all you need: Fast convergence at large depth. In *Uncertainty in Artificial Intelligence*, pages 1352–1361. PMLR.

Michele Bevilacqua, Rexhina Blloshmi, and Roberto Navigli. 2021. One spring to rule them both: Symmetric amr semantic parsing and generation without a complex pipeline. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12564–12573.

Gosse Bouma, Djamé Seddah, and Daniel Zeman. 2020. [Overview of the IWPT 2020 shared task on parsing into enhanced Universal Dependencies](#). In *Proceedings of the 16th International Conference on Parsing Technologies and the IWPT 2020 Shared Task on Parsing into Enhanced Universal Dependencies*, pages 151–161, Online. Association for Computational Linguistics.

Gosse Bouma, Djamé Seddah, and Daniel Zeman. 2021. [From raw text to enhanced Universal Dependencies: The parsing shared task at IWPT 2021](#). In *Proceedings of the 17th International Conference on Parsing Technologies and the IWPT 2021 Shared Task on Parsing into Enhanced Universal Dependencies (IWPT 2021)*, pages 146–157, Online. Association for Computational Linguistics.

Deng Cai and Wai Lam. 2019. Core semantic first: A top-down approach for amr parsing. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3799–3809.

Deng Cai and Wai Lam. 2020. Amr parsing via graph-sequence iterative inference. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1290–1301.

Danqi Chen and Christopher Manning. 2014. [A fast and accurate dependency parser using neural networks](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 740–750, Doha, Qatar. Association for Computational Linguistics.

Hao Cheng, Hao Fang, Xiaodong He, Jianfeng Gao, and Li Deng. 2016. [Bi-directional attention with agreement for dependency parsing](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2204–2214, Austin, Texas. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the*

*Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics. 667  
668  
669

Timothy Dozat and Christopher D Manning. 2016. [Deep biaffine attention for neural dependency parsing](#). *arXiv preprint arXiv:1611.01734*. 670  
671  
672

Daniel Fernández-González, Carlos Gómez-Rodríguez, and Carlos Gómez-Rodríguez. 2019. [Left-to-right dependency parsing with pointer networks](#). 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)*, pages 710–716, Minneapolis, Minnesota. Association for Computational Linguistics. 673  
674  
675  
676  
677  
678  
679  
680  
681

Stefan Grünewald and Annemarie Friedrich. 2020. [RobertNLP at the IWPT 2020 shared task: Surprisingly simple enhanced UD parsing for English](#). In *Proceedings of the 16th International Conference on Parsing Technologies and the IWPT 2020 Shared Task on Parsing into Enhanced Universal Dependencies*, pages 245–252, Online. Association for Computational Linguistics. 682  
683  
684  
685  
686  
687  
688  
689

Stefan Grünewald, Frederik Tobias Oertel, and Annemarie Friedrich. 2021. [RobertNLP at the IWPT 2021 shared task: Simple enhanced UD parsing for 17 languages](#). In *Proceedings of the 17th International Conference on Parsing Technologies and the IWPT 2021 Shared Task on Parsing into Enhanced Universal Dependencies (IWPT 2021)*, pages 196–203, Online. Association for Computational Linguistics. 690  
691  
692  
693  
694  
695  
696  
697  
698

Tao Ji, Yuanbin Wu, and Man Lan. 2019. [Graph-based dependency parsing with graph neural networks](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2475–2485, Florence, Italy. Association for Computational Linguistics. 699  
700  
701  
702  
703  
704

Jenna Kanerva, Filip Ginter, and Sampo Pyysalo. 2020. [Turku enhanced parser pipeline: From raw text to enhanced graphs in the IWPT 2020 shared task](#). In *Proceedings of the 16th International Conference on Parsing Technologies and the IWPT 2020 Shared Task on Parsing into Enhanced Universal Dependencies*, pages 162–173, Online. Association for Computational Linguistics. 705  
706  
707  
708  
709  
710  
711  
712

Eliyahu Kiperwasser and Yoav Goldberg. 2016. [Easy-first dependency parsing with hierarchical tree lstms](#). *Trans. Assoc. Comput. Linguistics*, 4:445–461. 713  
714  
715

Thomas N Kipf and Max Welling. 2016. [Semi-supervised classification with graph convolutional networks](#). *arXiv preprint arXiv:1609.02907*. 716  
717  
718

Dan Kondratyuk and Milan Straka. 2019. [75 languages, 1 model: Parsing Universal Dependencies universally](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing* 719  
720  
721  
722

723	and the 9th International Joint Conference on	Technologies and the IWPT 2021 Shared Task	779
724	Natural Language Processing (EMNLP-IJCNLP),	on Parsing into Enhanced Universal Dependencies	780
725	pages 2779–2795, Hong Kong, China. Association	(IWPT 2021), pages 189–195, Online. Association	781
726	for Computational Linguistics.	for Computational Linguistics.	782
727	Zuchao Li, Hai Zhao, and Kevin Parnow. 2020. Global	Xinyu Wang, Yong Jiang, and Kewei Tu. 2020. En-	783
728	greedy dependency parsing. In <u>Proceedings of</u>	hanced Universal Dependency parsing with second-	784
729	the AAAI conference on artificial intelligence, vol-	order inference and mixture of training data. In	785
730	ume 34, pages 8319–8326.	Proceedings of the 16th International Conference	786
731	Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng,	on Parsing Technologies and the IWPT 2020	787
732	Graham Neubig, and Eduard Hovy. 2018. <u>Stack-</u>	Shared Task on Parsing into Enhanced Universal	788
733	<u>pointer networks for dependency parsing.</u> In	Dependencies, pages 215–220, Online. Association	789
734	<u>Proceedings of the 56th Annual Meeting of the</u>	for Computational Linguistics.	790
735	<u>Association for Computational Linguistics (Volume</u>	David Weiss, Chris Alberti, Michael Collins, and Slav	791
736	<u>1: Long Papers),</u> pages 1403–1414, Melbourne, Aus-	Petrov. 2015. <u>Structured training for neural net-</u>	792
737	tralia. Association for Computational Linguistics.	<u>work transition-based parsing.</u> In <u>Proceedings of</u>	793
738	Marc’Aurelio Ranzato, Sumit Chopra, Michael	the 53rd Annual Meeting of the Association for	794
739	Auli, and Wojciech Zaremba. 2016. Se-	Computational Linguistics and the 7th International	795
740	quence level training with recurrent neural net-	Joint Conference on Natural Language Processing	796
741	works. In <u>4th International Conference on Learning</u>	(Volume 1: Long Papers), pages 323–333, Beijing,	797
742	<u>Representations, ICLR 2016.</u>	China. Association for Computational Linguistics.	798
743	Sebastian Schuster and Christopher D. Manning. 2016.	Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin	799
744	<u>Enhanced English Universal Dependencies: An</u>	Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-	800
745	<u>improved representation for natural language un-</u>	Yan Liu. 2021. Do transformers really perform	801
746	<u>derstanding tasks.</u> In <u>Proceedings of the Tenth</u>	badly for graph representation? <u>Advances in Neural</u>	802
747	<u>International Conference on Language Resources</u>	<u>Information Processing Systems, 34.</u>	803
748	<u>and Evaluation (LREC’16),</u> pages 2371–2378, Por-	Sheng Zhang, Xutai Ma, Kevin Duh, and Ben-	804
749	toroř, Slovenia. European Language Resources As-	jamin Van Durme. 2019a. Amr parsing as	805
750	sociation (ELRA).	sequence-to-graph transduction. In <u>Proceedings</u>	806
751	Tianze Shi and Lillian Lee. 2021. <u>TGIF: Tree-graph</u>	of the 57th Annual Meeting of the Association for	807
752	<u>integrated-format parser for enhanced UD with two-</u>	Computational Linguistics, pages 80–94.	808
753	<u>stage generic- to individual-language finetuning.</u> In	Sheng Zhang, Xutai Ma, Kevin Duh, and Benjamin	809
754	<u>Proceedings of the 17th International Conference</u>	Van Durme. 2019b. Broad-coverage semantic	810
755	<u>on Parsing Technologies and the IWPT 2021</u>	parsing as transduction. In <u>Proceedings of the</u>	811
756	<u>Shared Task on Parsing into Enhanced Universal</u>	2019 Conference on Empirical Methods in Natural	812
757	<u>Dependencies (IWPT 2021),</u> pages 213–224, Online.	Language Processing and the 9th International	813
758	Association for Computational Linguistics.	Joint Conference on Natural Language Processing	814
759	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	(EMNLP-IJCNLP), pages 3786–3798.	815
760	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	Yu Zhang, Zhenghua Li, and Min Zhang. 2020. <u>Effi-</u>	816
761	Kaiser, and Illia Polosukhin. 2017. Attention is	<u>cient second-order TreeCRF for neural dependency</u>	817
762	all you need. <u>Advances in neural information</u>	<u>parsing.</u> In <u>Proceedings of the 58th Annual Meeting</u>	818
763	<u>processing systems, 30.</u>	of the Association for Computational Linguistics,	819
764	Petar Veličković, Guillem Cucurull, Arantxa Casanova,	pages 3295–3305, Online. Association for Computa-	820
765	Adriana Romero, Pietro Lio, and Yoshua Bengio.	tional Linguistics.	821
766	2017. Graph attention networks. <u>arXiv preprint</u>	<b>A Appendix</b>	822
767	<u>arXiv:1710.10903.</u>	<b>A.1 GMS</b>	823
768	Xinyu Wang, Jingxian Huang, and Kewei Tu. 2019.	$Recall = \frac{\#matched\ graph}{\#gold\ sentences} \quad (13)$	824
769	<u>Second-order semantic dependency parsing with</u>		825
770	<u>end-to-end neural networks.</u> In <u>Proceedings of</u>	$Precision = \frac{\#matched\ graph}{\#system\ sentences} \quad (14)$	826
771	<u>the 57th Annual Meeting of the Association for</u>		827
772	<u>Computational Linguistics,</u> pages 4609–4618, Flo-	$GMS = \frac{2 \times Recall \times Precision}{Recall + Precision} \quad (15)$	828
773	rence, Italy. Association for Computational Linguis-	where # represents <i>The number of</i> .	
774	tics.		
775	Xinyu Wang, Zixia Jia, Yong Jiang, and Kewei Tu. 2021.		
776	<u>Enhanced Universal Dependency parsing with auto-</u>		
777	<u>mated concatenation of embeddings.</u> In <u>Proceedings</u>		
778	<u>of the 17th International Conference on Parsing</u>		

## 829 A.2 Experiment Details

830 We train models directly on each language with  
831 Teacher Forcing Training to output all head (depend-  
832 ent) representations at once. Besides, we truncate  
833 the input sentences to 100 words at training time.  
834 We totally run 100 epochs with 16 batch size and  
835 select the model parameters base on the ELAS on  
836 the development sets. We train our models on a  
837 single v100 with a speed of about 10000 samples  
838 in 10 minutes.

839 We use ReZero (Bachlechner et al., 2021) in  
840 our graph transformer, instead of LayerNorm op-  
841 erations commonly used in Transformer. In this  
842 case, we do not need to use the warm-up learning  
843 schedule, and we use Adam optimizer with the 0.97  
844 decay ratio of the learning rate. We set up the initial  
845 learning rate of pre-trained embeddings as  $2e - 5$ ,  
846 and that of others as  $1e - 3$ . Besides, dropout rates  
847 in the part of pre-training and graph representation  
848 are set to 0.1, while the output layers of nodes and  
849 edges are set to 0.3.

850 We build up the vocabulary on the arc labels for  
851 each language respectively. To shrink the size of  
852 edge vocabularies, we follow the de-lexicalization  
853 operations of arc labels (Grünwald and Friedrich,  
854 2020) and re-lexicalize them before evaluations.