Headed-Span-Based Projective Dependency Parsing

Anonymous ACL submission

Abstract

We propose a new paradigm for projective dependency parsing based on headed spans. In a projective dependency tree, the subtree rooted at each word covers a contiguous sequence (i.e., a span) in the surface order. We call a span marked with a root word headed span. A projective dependency tree can be represented as a collection of headed spans. We decompose the score of a dependency tree into the scores of the headed spans and design a novel $O(n^3)$ dynamic programming algorithm to enable global training and exact inference. The advantages of our headedspan-based dependency parsing include that it captures subtree information more adequately than first-order graph-based methods and performs global optimization in decoding (in contrast to transition-based methods). We evaluate our method on PTB, CTB, and UD and it achieves competitive results in comparison with previous methods.

1 Introduction

Dependency parsing is an important task in natural language processing, which has numerous applications in downstream tasks, such as opinion mining (Zhang et al., 2020a), relation extraction (Jin et al., 2020), named entity recognition (Jie and Lu, 2019), machine translation (Bugliarello and Okazaki, 2020), among others.

There are two main paradigms in dependency parsing: graph-based and transition-based methods. Graph-based methods decompose the score of a tree into the scores of parts. In the simplest first-order graph-based methods (McDonald et al., 2005, *inter alia*), the parts are single dependency arcs. In higher-order graph-based methods (Mc-Donald and Pereira, 2006; Carreras, 2007; Koo and Collins, 2010; Ma and Zhao, 2012), the parts are combinations of multiple arcs. Transition-based methods (Nivre and Scholz, 2004; Chen and Manning, 2014, *inter alia*) read the sentence sequentially and conduct a series of local decisions to build the final parse. Recently, transition-based methods with Pointer Networks (Vinyals et al., 2015) have obtained competitive performance to graph-based methods (Ma et al., 2018; Liu et al., 2019; Fernández-González and Gómez-Rodríguez, 2019; Fernández-González and Gómez-Rodríguez, 2021). 042

043

044

045

047

050

051

053

054

059

060

061

062

063

064

065

066

067

068

069

071

072

074

075

076

077

079

081

A main limitation of first-order graph-based methods is that they independently score each arc based solely on the two words connected by the arc. Ideally, the appropriateness of an arc should depend on the whole parse tree, particularly the subtrees rooted at the two words connected by the arc. Although subtree information could be implicitly encoded (Falenska and Kuhn, 2019) in powerful neural encoders such as LSTMs (Hochreiter and Schmidhuber, 1997) and Transformers (Vaswani et al., 2017), there is evidence that their encoding of such information is inadequate. For example, higher-order graph-based methods, which capture more subtree information by simultaneously considering multiple arcs, have been found to outperform first-order methods despite using powerful encoders (Fonseca and Martins, 2020; Zhang et al., 2020b). In contrast to the line of work on higherorder parsing, we propose a different way to incorporate more subtree information as we will discuss later.

Transition-based methods, on the other hand, can easily utilize information from partially built subtrees, but they have their own shortcomings. For instance, they cannot perform global optimization during decoding and rely on greedy or beam search to find a locally optimal parse; and their sequential decoding may cause error propagation as past decision mistakes will negatively affect the decisions in the future.

To overcome the aforementioned limitations of first-order graph-based and transition-based methods, we propose a new paradigm of projective de-



Figure 1: An example projective dependency parse tree. Each rectangle represents a headed span. A projective parse tree can be treated as a collection of headed spans.

pendency parsing based on so-called headed spans. A projective dependency tree has a nice structural property that the subtree rooted at each word covers a contiguous sequence (i.e., a span) in the sur-086 face order. We call such a span marked with its root word a *headed span*. A projective dependency tree can be treated as a collection of headed spans such that each word corresponds to exactly one 090 headed span, as illustrated in Figure 1. For example, (0, 5, inventory) is a headed span, in which span (0,5) has a head word *inventory*. In this view, projective dependency parsing is similar to constituency parsing as a constituency tree can be treated as a collection of constituent spans. The main difference is that in a binary constituency tree, a constituent span (i, k) is made up by two adjacent spans (i, j) and (j, k), while in a projective dependency tree, a headed span (i, k, x_h) is made up by 100 one or more smaller headed spans and a single word 101 span (h - 1, h). For instance, (0, 5, inventory) is 102 made up by (0, 1, An), (1, 2) and (2, 5, of). There 103 are a few constraints between headed spans to force 104 projectivity ($\S2.3$). These structural constraints are the key to designing an efficient dynamic program-106 ming algorithm for exact inference. 107

Because of the similarity between constituency 108 parsing and our head-span-based view of projec-109 tive dependency parsing, we can draw inspirations 110 from the constituency parsing literature to design 111 our dependency parsing method. Specifically, span-112 based constituency parsers (Stern et al., 2017; Ki-113 taev and Klein, 2018; Zhang et al., 2020c; Xin 114 et al., 2021) decompose the score of a constituency 115 tree into the scores of its constituent spans and 116

use the CYK algorithm (Cocke, 1969; Younger, 1967; Kasami, 1965) for global training and inference. Built upon powerful neural encoders, they have obtained state-of-the-art performance in constituency parsing. Inspired by them, we propose to decompose the score of a projective dependency tree into the scores of headed spans and design a novel $O(n^3)$ dynamic programming algorithm for global training and exact inference, which is on par with the Eisner algorithm (Eisner, 1996) in time complexity for projective dependency parsing. We make a departure from existing graph-based methods since we do not model dependency arcs directly. Instead, the dependency arcs are *induced* from the collection of headed spans ($\S2.3$). Compared with first-order graph-based methods, our method can utilize more subtree information since a headed span contains all children (if any) of the corresponding headword (and all words within the subtree). Compared with transition-based methods, our method allows global training and exact inference and does not suffer from error propagation or exposure bias.

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

Our contributions can be summarized as follows:

- We treat a projective dependency tree as a collection of headed spans, providing a new perspective of projective dependency parsing.
- We design a novel $O(n^3)$ dynamic programming algorithm to enable global training and exact inference for our proposed model.
- We have obtained the state-of-the-art or competitive results on PTB, CTB, and UD v2.2, showing the effectiveness of our proposed method.

2 Model

151

152

153

154

155

157

158

159

161

163

165

166

167

168

170

171

173

174

175

176

177

178

179

180

181

182

184

187

We adopt the two-stage parsing strategy, i.e., we first predict an unlabeled tree and then predict the dependency labels. Given a sentence $x_1, ..., x_n$, its unlabeled projective dependency parse tree ycan be regarded as a collection of headed spans (l_i, r_i, x_i) where $1 \le i \le n$. For each word x_i , we can find exactly one headed span (l_i, r_i, i) (where l_i and r_i are the left and right span boundaries) given parse tree y, so there are totally n headed spans in y as we can see in Figure 1. We can use a simple post-order traversal algorithm to obtain all headed spans in O(n) time.

We then define the score of y as:

$$s(y) = \sum_{i=1,\dots,n} s_{l_i,r_i,i}^{\text{span}}$$

In §2.1, we show how to calculate $s_{l_i,r_i,i}^{\text{span}}$ using neural networks. In §2.2, we present the training objective function and in §2.3, we present the novel $O(n^3)$ parsing algorithm.

2.1 Neural encoding and scoring

We add <bos> (beginning of sentence) at x_0 and <eos> (end of sentence) at x_{n+1} . In the embedding layer, we apply mean-pooling to the last layer of BERT (Devlin et al., 2019) (i.e., taking the mean value of all subword embeddings) to generate dense word-level representation e_i for each token x_i ⁻¹. Then we feed $e_0, ..., e_{n+1}$ into a 3-layer bidirectional LSTM (BiLSTM) to get $c_0, ..., c_{n+1}$, where $c_i = [f_i; b_i]$ and f_i and b_i are the forward and backward hidden states of the last BiLSTM layer at position *i* respectively. We then use $e_{i,j}$ to represent span (i, j):

$$h_k = [f_k, b_{k+1}]$$

 $e_{i,j} = h_j - h_j$

After obtaining the word and span representations, we use deep biaffine function (Dozat and Manning, 2017) to score headed spans:

188
189
190

$$c'_{k} = \text{MLP}_{\text{word}}(c_{k})$$

 $e'_{i,j} = \text{MLP}_{\text{span}}(e_{i,j})$
 $s^{\text{span}}_{i,j,k} = [c'_{k}; 1]^{\top} W^{\text{span}}[e'_{i,j}; 1]$

where MLP_{word} and MLP_{span} are multi-layer perceptrons (MLPs) that project word and span representations into *d*-dimensional spaces respectively; $W^{\text{span}} \in \mathcal{R}^{(d+1)\times(d+1)}$.

Similarly, we use deep biaffine functions to score the labels of dependency arcs for a given gold or predicted tree. In our preliminary experiments, we find that directly calculating the scores based on parent-child word representations leads to a slightly better result:

$$c_i' = \text{MLP}_{\text{parent}}(c_i)$$
 20

191

192

193

194

195

196

197

198

199

200

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

228

229

230

231

$$c'_j = \mathrm{MLP}_{\mathrm{child}}(c_j)$$
 202

$$s_{i,j,r}^{\text{label}} = \begin{bmatrix} c'_i; 1 \end{bmatrix}^\top W_r^{\text{label}} \begin{bmatrix} c'_j; 1 \end{bmatrix}$$
 203

where MLP_{parent} and MLP_{child} are MLPs that map span representations into d'-dimensional spaces; $W_r^{\text{label}} \in \mathcal{R}^{(d'+1)\times(d'+1)}$ for each relation type $r \in R$ in which R is the set of all relation types.

2.2 Training loss

Following previous work, we decompose the training loss into the unlabeled parse loss and arc label loss:

$$L = L_{\text{parse}} + L_{\text{label}}$$

For L_{parse} , we can either design a local spanselection loss (i.e., maximize the probability of the gold span for each word over all feasible spans) which is akin to the head-selection loss (Dozat and Manning, 2017), or use global structural loss. Experimentally, we find that the max-margin loss (Taskar et al., 2004) performs best. The maxmargin loss aims to maximize the margin between the score of the gold tree y and the incorrect tree y'of the highest score:

$$L_{\text{parse}} = \max(0, \max_{y' \neq y} (s(y') + \Delta(y', y) - s(y))$$
(1)

where Δ measures the difference between the incorrect tree and gold tree. Here we let Δ to be the Hamming distance (i.e., the total number of mismatch of headed spans). To calculate Eq. 1, we need to perform loss-augmented inference (Taskar et al., 2005). To achieve that, we update the scores as:

$$s'_{i,j,k} = s_{i,j,k} + \mathbf{1} \left((i,j,k) \notin y \right)$$

where $\mathbf{1}((i, j, k) \notin y)$ means that the headed span (i, j, x_k) does not exists in y. Then we can use

¹For some datasets (e.g., Chinese Treebank), we concatenate the POS tag embedding with the BERT embedding as e_i .



Figure 2: An example subtree.

the parsing algorithm (§2.3) based on the updated
scores to obtain the highest-scoring tree. If it is the
gold tree, then the loss is 0; otherwise, we can put
it back to Eq. 1 to calculate the loss.

Finally, we use cross entropy for L_{label} :

$$L_{\text{label}} = \sum_{(x_i \to x_j, r) \in y} -\log \frac{\exp(s_{i,j,r}^{\text{label}})}{\sum_{r' \in R} \exp(s_{i,j,r'}^{\text{label}})}$$

where $(x_i \to x_j, r) \in y$ denotes every dependency arc from x_i to x_j with label r in y.

2.3 Parsing

241

242

243

244

245

246

247

249

250

251

255

256

263

265

In this section, we detail our proposed $O(n^3)$ inference algorithm. The algorithm is based on the following key observations:

- For a given parent word x_k, if it has any children to the left (right), then all headed spans of its children in this direction should be consecutive and form a larger span, which we refer to as the left (right) child span. The left (right) boundary of the headed span of x_k is the left (right) boundary of the leftmost (rightmost) child span, or k 1 (k) if x_k has no child to the left (right).
- If a parent word xk has children in both directions, then its left span and right span are separated by the single word span (k 1, k).

Figure 2 shows an example subtree. The left child span is (i, j - 1) and the right child span is (j, k). They are separated by the single word span (j - 1, j). The headed span (i, k, j) can be generated by merging the left child span, right child span, and the single word span. Note that the left (right) child span can contain one or more headed spans.

Based on these observations, we design the following dynamic programming chart items:

α_{i,j}: the accumulated score of span (*i*, *j*) serving as a left or right child span.

• $\beta_{i,j,k}$: the accumulated score of the headed span (i, j, k).

Then we can define the following recursive formulas to perform dynamic programming:

$$\beta_{i,i+1,i+1} = s_{i,i+1,i+1}^{\text{span}}$$
 (2) 2

$$\alpha_{i,i} = 0 \tag{3}$$

$$\beta_{i,j,k} = \alpha_{i,k-1} + \alpha_{k,j} + s_{i,j,k}^{\text{span}} \tag{4}$$

$$\alpha_{i,j} = \max(\max_{i < k < j} (\alpha_{i,k} + \alpha_{k,j}),$$

$$\max_{i < h \le j} (\beta_{i,j,h})) \tag{5}$$

We set $\alpha_{i,i} = 0$ for the convenience of calculating $\beta_{i,j,k}$ when x_k does not have children on either side. In Eq. 5, we can see that the child span comes from either multiple smaller consecutive child spans (i.e., $\max_{i < k < j} (\alpha(i, k) + \alpha(k, j)))$ or a single headed span (i.e., $\max_{i < k \leq j} (\beta(i, j, h)))$).

We also maintain the following backtrack points in order to recover the predicted projective tree:

$$B_{i,j} = \begin{cases} 1, & \alpha_{i,j} = \max_{i < h \le j} (\beta_{i,j,h}) \\ 0, & \alpha_{i,j} = \max_{i < k < j} (\alpha_{i,k} + \alpha_{k,j}) \end{cases}$$
288

$$C_{i,j} = \underset{i < k < j}{\arg \max} (\alpha_{i,k} + \alpha_{k,j})$$
290

$$H_{i,j} = \operatorname*{arg\,max}_{i < h \le j}(eta_{i,j,h})$$
 291

The parsing algorithm first computes all the chart items defined above and then recovers the parse tree from top down. The root is $H_{0,n}$. For a given headed span, it finds the best segmentation of left child spans and right child spans, and then adds dependency arcs between the headword of the given headed span and the headword of each child span. Finding the best segmentation is similar to the inference procedure of the semi-Markov CRF model (Sarawagi and Cohen, 2004). Then we apply the same procedure to each child headed span (within the best segmentation) recursively. The whole parsing algorithm is formalized in Algorithm 1.

Time complexity: From Eq. 2 to 5, we can see that at most three variables (i.e., i, j, k) are required to iterate over, so the total time complexity is $O(n^3)$.

297

299

301

302

303

304

305

306

308

271

272

273

274

279

281

282

283

Algorithm 1 Inference algorithm for headed spanbased projective dependency parsing

Require: Input sentence of length n
Calculate all α, β, B, C, H .
$\operatorname{arcs} \leftarrow \{(\operatorname{ROOT} \to H_{0,n})\}$
function FINDARC (i, j)
if $i + 1 = j$ then
return {j}
else if $B_{i,j} = 1$ then
$h \leftarrow H_{i,j}$
if i+1 < h < j then
$L \leftarrow \text{FINDARC}(i, h - 1)$
$R \leftarrow \text{FINDARC}(h, j)$
Children $\leftarrow L \cup R$
else if $h = j$ then
Children \leftarrow FINDARC $(i, j - 1)$
else
Children \leftarrow FINDARC $(i + 1, j)$
end if
for c in Children do
$rcs \leftarrow rcs \cup (h ightarrow c)$
end for
return {h}
else
$c \leftarrow C_{i,j}$
$L \leftarrow \text{FINDARC}(i, c)$
$R \leftarrow FINDARC(c, j)$
return $L \cup R$
end if
end function
FINDARC(0, n)
return arcs

3 Experiments

3.1 Data

310

311

312

313

314 315

316

317

319

320

321

322

324

326

We evaluate our proposed method on Penn Treebank (PTB) 3.0 (Marcus et al., 1993), Chinese Treebank (CTB) 5.1 (Xue et al., 2005) and 12 languages on Universal Dependencies (UD) 2.2. For PTB, we use the Stanford Dependencies conversion software of version 3.3 to obtain dependency trees. For CTB, we use head-rules from Zhang and Clark (2008) and Penn2Malt² to obtain dependency trees. Following Wang and Tu (2020), we use gold POS tags for CTB and UD. We do not use POS tags in PTB. For PTB/CTB, we drop all nonprojective trees during training. For UD, we use MaltParser v1.9.2 3 to adopt the pseudo-projective transformation (Nivre and Nilsson, 2005) to convert nonprojective trees into projective trees when training, and convert back when evaluating.

3.2 Evaluation methods

We report the unlabeled attachment score (UAS) and labeled attachment score (LAS) averaged from three runs with different random seeds. In each run, we select the model based on the performance on the development set. Following Wang and Tu (2020), we ignore all punctuation marks during evaluation. 327

328

330

331

332

333

334

335

336

337

338

339

340

341

343

344

345

346

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

3.3 Implementation details

We use "bert-large-cased" for PTB, "bert-basechinese" for CTB, and "bert-multilingual-cased" for UD, so the dimension of the input BERT embedding is 1024, 768, and 768 respectively. The dimension of POS tag embedding is set to 100 for CTB and UD. The hidden size of BiLSTM is set to 1000. The hidden size of biaffine functions is set to 600 for scoring spans and arcs (used in our reimplemented Biaffine Parser), 300 for scoring labels. We add a dropout layer after the embedding layer, LSTM layers, and MLP layers. The dropout rate is set to 0.33. We use Adam (Kingma and Ba, 2015) as the optimizer with $\beta_1 = 0.9, \beta_2 = 0.9$ to train our model for 10 epochs. The maximal learning rate is lr = 5e - 5 for BERT and lr = 25e - 5 for other components. We linearly warmup the learning rate to the maximal value for the first epoch and gradually decay it to zero for the rest of the epochs. The value of gradient clipping is set to 5. We batch sentences of similar lengths to better utilize GPUs. The token number is 4000 for each batch, i.e., the sum of lengths of sentences is 4000.

3.4 Baselines

- **Biaffine**: Dozat and Manning (2017) are the first use deep biaffine functions to score dependency arcs/labels and use the local head-selection training loss function.
- **TreeCRF2O**: Zhang et al. (2020b) use a deep triaffine function to score sibling factors and use a second-order TreeCRF loss for training.
- MFVI2O: Wang and Tu (2020) use decomposed triaffine functions to score second-order factors (i.e., grandparents and siblings) and unfold mean-field variational inference procedure for end-to-end training.
- **HPSG**: Zhou and Zhao (2019) propose a span-based method to perform joint dependency and constituency parsing by simplifying the head-driven phrase structure grammars (HPSG) (Pollard and Sag, 1994).

²https://cl.lingfil.uu.se/~nivre/ research/Penn2Malt.html

³http://www.maltparser.org/download. html

	P	ГВ	СТВ				
	UAS	LAS	UAS	LAS			
MFVI2O	95.98	94.34	90.81	89.57			
TreeCRF2O	96.14	94.49	-	-			
HierPtr	96.18	94.59	90.76	89.67			
	+BEI	RT _{base}	+BERT _{base}				
RNGTr	96.66	95.01	92.98	91.18			
	+BEF	RT _{large}	+BERT _{base}				
MFVI2O	96.91	95.34	92.55	91.69			
HierPtr	97.01	95.48	92.65	91.47			
$Biaffine + MM^{\dagger}$	97.22	95.71	93.18	92.10			
Ours	97.24	95.73	93.33	92.30			
For reference							
	+XLNet _{large}		+BEF	RT _{base}			
$HPSG^{\flat}$	97.20	95.72	-	-			
$HPSG+LAL^{\flat}$	97.42	96.26	94.56	89.28			

Table 1: Results for different model on PTB and CTB.^b indicate that they use additional annotated constituency trees in training. [†] means our reimplementation.

- **HPSG+LAL**: Mrini et al. (2020) add label attention layer (LAL) upon HPSG.
- **RNGTr**: Mohammadshahi and Henderson (2021) propose an iterative refinement network built upon Transformer.
- **HierPtr**: Fernández-González and Gómez-Rodríguez (2021) improve the mechanisms of transition-based methods with Pointer Nets.

3.5 Main result

Table 1 shows the results on PTB and CTB. Note that Biaffine+MM is our reimplementation of the Biaffine Parser that uses the same setting as in our method, including the use of the max-margin loss instead of the local head-selection loss. Interestingly, we find that Biaffine+MM has already surpassed many strong baselines, and this may due to the proper choices of hyperparameters and the use of the max-margin loss (we observe that using the max-margin loss leads to a better performance compared with the original head-selection loss), so Biaffine+MM is a very strong baseline. It also has the same number of parameters as our methods. Our method surpasses Biaffine+MM on both datasets, showing the competitiveness of our headed-spanbased method in a fair comparison with first-order graph-based parsing. Our method also obtains the state-of-the-art result among methods that only use dependency training data (HPSG+LAL uses additional constituency trees as training data, so it is not directly comparable with the other systems.).

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

Table 2 shows the results on UD. We can see that our reimplemented Biaffine+MM has already surpassed MFVI2O, which utilizes higher-order information. Our method outperforms Biaffine+MM by 0.14 LAS on average, validating the effectiveness of our proposed method in the multilingual scenarios.

4 Analysis

4.1 Influence of training loss function

Table 3 shows the influence of the training loss function. We find that the max-margin loss performs better on both datasets: 0.17 UAS improvement on PTB and 0.05 UAS improvement on CTB comparing to the local span-selection loss, which shows the effectiveness of using global loss.

4.2 Error analysis

As previously argued, first-order graph-based methods are insufficient to model complex subtrees, so they may have difficulties in parsing long sentences and handling long-range dependencies. To verify this, we follow (McDonald and Nivre, 2011) to plot UAS as a function of the sentence length and plot F1 scores as functions of the distance to root and dependency length on the CTB test set. We additionally plot the F1 score of the predicted headed spans against the gold headed spans with different span lengths.

From Figure 3a, we can see that Biaffine+MM has a better UAS score on short sentences (of length <=20), while for long sentences (of length >=30), our headed span-based method has a higher performance, which validates our conjecture.

Figure 3b shows the F1 score for arcs of varying distances to root. Our model is better at predicting arcs of almost all distances to root in the dependency tree, which reveals our model's superior ability to predict complex subtrees.

Figure 3c shows the F1 score for arcs of the varying lengths. Both Biaffine+MM and our model have a very similar performance in predicting arcs of distance < 7, while our model is better at predicting arcs of distance >= 7, which validates the ability of our model at capturing long-range dependencies.

Figure 3d shows the F1 score for headed spans of the varying lengths. We can see that when the span length is small (<=10), Biaffine+MM and our model have a very similar performance. However,

391

397

400

401

402

403

376

377

	bg	ca	cs	de	en	es	fr	it	nl	no	ro	ru	Avg
TreeCRF2O MFVI2O	90.77 90.53	91.29 92.83	91.54 92.12	80.46 81.73	87.32 89.72	90.86 92.07	87.96 88.53	91.91 92.78	88.62 90.19	91.02 91.88	86.90 85.88	93.33 92.67	89.33 90.07
+BERT _{multilingual}													
MFVI2O	91.30	93.60	92.09	82.00	90.75	92.62	89.32	93.66	91.21	91.74	86.40	92.61	90.61
$Biaffine + MM^{\dagger}$	90.30	94.49	92.65	85.98	91.13	93.78	91.77	94.72	91.04	94.21	87.24	94.53	91.82
Ours	91.10	94.46	92.57	85.87	91.32	93.84	91.69	94.78	91.65	94.28	87.48	94.45	91.96

Table 2: Labeled Attachment Score (LAS) on twelves languages in UD 2.2. We use ISO 639-1 codes to represent languages. † means our implementation.



Figure 3: Error analysis on the CTB test set.

	P	ГВ	СТВ		
	UAS	LAS	UAS	LAS	
max-margin loss span-selection loss	97.24 97.07	95.73 95.50	93.33 93.28	92.30 92.20	

Table 3: The influence of training loss function on PTB and CTB.

our model is much better in predicting longer spans (especially for spans of length >30).

4.3 Parsing speed

453

454

455

456

457

458

One of the advantages of transition-based methods is that they have a low complexity in parsing a sentence. In contrast, first-order graph-based methods require $O(n^2)$ time to produce an unrestricted 459 tree by using the maximum spanning tree (MST) 460 algorithm and require $O(n^3)$ time to produce a pro-461 jective tree by using the Eisner algorithm. As we 462 discussed in §2.3, our proposed parsing algorithm 463 also has a $O(n^3)$ time complexity, which seems 464 slow. Luckily, inspired by Zhang et al. (2020b) 465 and Rush (2020) who independently propose to 466 batchify the Eisner algorithm using Pytorch, we 467 batchify our proposed method so that $O(n^2)$ out of 468 $O(n^3)$ can be computed in parallel, which greatly 469 accelerates parsing. We achieve a similar parsing 470 speed of our method to the fast implementation of 471 the Eisner algorithm by Zhang et al. (2020b): it 472 takes 20 seconds to parse the entire PTB test set 473

474 475

476

using BERT as the encoder under a single TITAN RTX GPU.

5 Related work

Dependency parsing with more complex sub-477 tree information: There has always been an in-478 terest to incorporate more complex subtree infor-479 mation into graph-based and transition-based meth-480 ods since their invention. Before the deep learning 481 era, it is difficult to incorporate sufficient contex-482 tual information in first-order graph-based parsers. 483 To mitigate this, researchers develop higher-order 484 dependency parsers to capture more contextual 485 486 information (McDonald and Pereira, 2006; Carreras, 2007; Koo and Collins, 2010; Ma and Zhao, 487 2012). However, incorporating more complex fac-488 tors worsens inference time complexity. For ex-489 ample, exact inference for third-order projective 490 dependency parsing has a $O(n^4)$ time complexity 491 and exact inference for higher-order non-projective 492 dependency parsing is NP-hard (McDonald and 493 Pereira, 2006). To decrease inference complex-494 ity, researchers use approximate parsing methods. 495 Smith and Eisner (2008) use belief propagation 496 (BP) framework for approximate inference to trade 497 accuracy for efficiency. They show that third-order 498 parsing can be done in $O(n^3)$ time using BP. Gorm-499 ley et al. (2015) unfold the BP process and use gra-500 dient descent to train their parser in an end-to-end manner. Wang and Tu (2020) extend their work by 502 using neural scoring functions to score factors. For 503 higher-order non-projective parsing, researchers re-504 sort to dual decomposition algorithm (e.g., AD^3) 505 for decoding (Martins et al., 2011, 2013). They observe that the approximate decoding algorithm often obtains exact solutions. Fonseca and Mar-508 tins (2020) combine neural scoring functions and their decoding algorithms for non-projective higher-510 order parsing. Zheng (2017) propose a incremental 511 graph-based method to utilize higher-order infor-512 mation without hurting the advantage of global 513 inference. Ji et al. (2019) use a graph attention 514 network to incorporate higher-order information 515 into the Biaffine Parser. Zhang et al. (2020b) en-516 hance the Biaffine Parser by using a deep triaffine 517 function to score sibling factors. Mohammadshahi 518 and Henderson (2021) propose an iterative refine-519 ment network that injects the predicted soft trees 520 from the previous iteration to the self-attention lay-521 ers to predict the soft trees of the next iteration, so that information of the whole tree is considered in parsing. As for transition-based methods, de Lhoneux et al. (2019) explore the impact of the way of subtree composition, Ma et al. (2018); Liu et al. (2019); Fernández-González and Gómez-Rodríguez (2021) incorporate sibling and grandparent information into transition-based parsing with Pointer Networks.

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

566

567

568

570

571

572

Span-based constituency parsing: Span-based parsing is originally proposed in continuous constituency parsing (Stern et al., 2017; Kitaev and Klein, 2018; Zhang et al., 2020c; Xin et al., 2021). Span-based constituency parsers decompose the score of a constituency tree into the scores of its constituents. Recovering the highest-scoring tree can be done via the exact CYK algorithm or greedy top-down approximate inference algorithm (Stern et al., 2017). Kitaev and Klein (2018) propose a self-attentive network to improve the parsing accuracy. They separate content and positional attentions and show the improvement. Zhang et al. (2020c) use a two-stage bracketing-then-labeling framework and replace the max-margin loss with the TreeCRF loss (Finkel et al., 2008). Xin et al. (2021) recently propose a recursive semi-Markov model, incorporating sibling factor scores into the score of a tree to explicitly model n-ary branching structures. Corro (2020) adapts span-based parsing to discontinuous constituency parsing and obtains the state-of-the-art result.

6 Conclusion

In this work, we have presented a headed-spanbased method for projective dependency parsing. Our proposed method can utilize more subtree information and meanwhile enjoy global training and exact inference. Experiments show that our proposed method has a high parsing accuracy in PTB, CTB, and twelve languages in UD v2.2. In addition to its empirical competitiveness, we believe our work provides a novel perspective of projective dependency parsing and could lay the foundation for further theoretical and algorithmic advancements in the future.

References

Emanuele Bugliarello and Naoaki Okazaki. 2020. Enhancing machine translation with dependency-aware self-attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1618–1627, Online. Association for Computational Linguistics.

 Xavier Carreras. 2007. Experiments with a higherorder projective dependency parser. In Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL), pages 957–961, Prague, Czech Republic. Association for Computational Linguistics.

573

574

580

582

588

590

594

596

599

602

603

612

613

615

617

618

619

- 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.
- J. Cocke. 1969. Programming languages and their compilers: Preliminary notes.
- Caio Corro. 2020. Span-based discontinuous constituency parsing: a family of exact chart-based algorithms with time complexities from O(n⁶) down to O(n³). In *Proceedings of the 2020 Conference* on Empirical Methods in Natural Language Processing (EMNLP), pages 2753–2764, Online. Association for Computational Linguistics.
 - Miryam de Lhoneux, Miguel Ballesteros, and Joakim Nivre. 2019. Recursive subtree composition in LSTM-based dependency parsing. 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 1566–1576, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Timothy Dozat and Christopher D. Manning. 2017. Deep biaffine attention for neural dependency parsing. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. Open-Review.net.
- Jason M. Eisner. 1996. Three new probabilistic models for dependency parsing: An exploration. In COL-ING 1996 Volume 1: The 16th International Conference on Computational Linguistics.
- Agnieszka Falenska and Jonas Kuhn. 2019. The (non-)utility of structural features in BiLSTM-based dependency parsers. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 117–128, Florence, Italy. Association for Computational Linguistics.

Daniel Fernández-González 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. 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

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

- Daniel Fernández-González and Carlos Gómez-Rodríguez. 2021. Dependency parsing with bottom-up hierarchical pointer networks. *CoRR*, abs/2105.09611.
- Jenny Rose Finkel, Alex Kleeman, and Christopher D. Manning. 2008. Efficient, feature-based, conditional random field parsing. In *Proceedings of ACL-*08: HLT, pages 959–967, Columbus, Ohio. Association for Computational Linguistics.
- Erick Fonseca and André F. T. Martins. 2020. Revisiting higher-order dependency parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8795– 8800, Online. Association for Computational Linguistics.
- Matthew R. Gormley, Mark Dredze, and Jason Eisner. 2015. Approximation-aware dependency parsing by belief propagation. *Transactions of the Association for Computational Linguistics*, 3:489–501.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Tao Ji, Yuanbin Wu, and Man Lan. 2019. Graphbased 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.
- Zhanming Jie and Wei Lu. 2019. Dependency-guided LSTM-CRF for named entity recognition. 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 3862–3872, Hong Kong, China. Association for Computational Linguistics.
- Lifeng Jin, Linfeng Song, Yue Zhang, Kun Xu, Wei-Yun Ma, and Dong Yu. 2020. Relation extraction exploiting full dependency forests. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI* 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 8034– 8041. AAAI Press.
- T. Kasami. 1965. An efficient recognition and syntaxanalysis algorithm for context-free languages.

797

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.

688

689

698

710

711

712

713

714

715

716

718

719

720

721

722

723

724

726

727

728

729 730

731

732

733

734

735

736

737

738

739

740

- Nikita Kitaev and Dan Klein. 2018. Constituency parsing with a self-attentive encoder. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.
- Terry Koo and Michael Collins. 2010. Efficient thirdorder dependency parsers. In Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, pages 1–11, Uppsala, Sweden. Association for Computational Linguistics.
 - Linlin Liu, Xiang Lin, Shafiq Joty, Simeng Han, and Lidong Bing. 2019. Hierarchical pointer net 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 1007– 1017, Hong Kong, China. Association for Computational Linguistics.
- Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng, Graham Neubig, and Eduard Hovy. 2018. Stackpointer networks for dependency parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1403–1414, Melbourne, Australia. Association for Computational Linguistics.
- Xuezhe Ma and Hai Zhao. 2012. Fourth-order dependency parsing. In Proceedings of COLING 2012: Posters, pages 785–796, Mumbai, India. The COL-ING 2012 Organizing Committee.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- André Martins, Miguel Almeida, and Noah A. Smith. 2013. Turning on the turbo: Fast third-order nonprojective turbo parsers. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 617–622, Sofia, Bulgaria. Association for Computational Linguistics.
- André Martins, Noah Smith, Mário Figueiredo, and Pedro Aguiar. 2011. Dual decomposition with many overlapping components. In Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, pages 238–249, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Ryan McDonald, Koby Crammer, and Fernando Pereira. 2005. Online large-margin training of dependency parsers. In *Proceedings of the 43rd*

Annual Meeting of the Association for Computational Linguistics (ACL'05), pages 91–98, Ann Arbor, Michigan. Association for Computational Linguistics.

- Ryan McDonald and Joakim Nivre. 2011. Analyzing and integrating dependency parsers. *Computational Linguistics*, 37(1):197–230.
- Ryan McDonald and Fernando Pereira. 2006. Online learning of approximate dependency parsing algorithms. In 11th Conference of the European Chapter of the Association for Computational Linguistics, Trento, Italy. Association for Computational Linguistics.
- Alireza Mohammadshahi and James Henderson. 2021. Recursive non-autoregressive graph-to-graph transformer for dependency parsing with iterative refinement. *Trans. Assoc. Comput. Linguistics*, 9:120– 138.
- Khalil Mrini, Franck Dernoncourt, Quan Hung Tran, Trung Bui, Walter Chang, and Ndapa Nakashole. 2020. Rethinking self-attention: Towards interpretability in neural parsing. In *Findings of the Association for Computational Linguistics: EMNLP* 2020, pages 731–742, Online. Association for Computational Linguistics.
- Joakim Nivre and Jens Nilsson. 2005. Pseudoprojective dependency parsing. In Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05), pages 99–106, Ann Arbor, Michigan. Association for Computational Linguistics.
- Joakim Nivre and Mario Scholz. 2004. Deterministic dependency parsing of English text. In *COLING* 2004: Proceedings of the 20th International Conference on Computational Linguistics, pages 64–70, Geneva, Switzerland. COLING.
- Carl Pollard and Ivan A. Sag. 1994. *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago.
- Alexander Rush. 2020. Torch-struct: Deep structured prediction library. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 335– 342, Online. Association for Computational Linguistics.
- Sunita Sarawagi and William W. Cohen. 2004. Semimarkov conditional random fields for information extraction. In Advances in Neural Information Processing Systems 17 [Neural Information Processing Systems, NIPS 2004, December 13-18, 2004, Vancouver, British Columbia, Canada], pages 1185– 1192.
- David Smith and Jason Eisner. 2008. Dependency parsing by belief propagation. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 145–156, Honolulu, Hawaii. Association for Computational Linguistics.

Mitchell Stern, Jacob Andreas, and Dan Klein. 2017. A Bo Zhang, Yue Zhang, Rui Wang, Zhenghua Li, and Min minimal span-based neural constituency parser. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 818-827, Vancouver, Canada. Association for Computational Linguistics.

799

801

810

811

812

813

815

816

817

818

819

823

824

827

832

833

839

841

846

847

849

- Ben Taskar, Dan Klein, Mike Collins, Daphne Koller, and Christopher Manning. 2004. Max-margin parsing. In Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing, pages 1-8, Barcelona, Spain. Association for Computational Linguistics.
- Benjamin Taskar, Vassil Chatalbashev, Daphne Koller, and Carlos Guestrin. 2005. Learning structured prediction models: a large margin approach. In Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005, volume 119 of ACM International Conference Proceeding Series, pages 896-903. ACM.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998-6008.
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. In Advances in Neural Junru Zhou and Hai Zhao. 2019. Head-Driven Phrase Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, pages 2692-2700.
 - Xinyu Wang and Kewei Tu. 2020. Second-order neural dependency parsing with message passing and endto-end training. In Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, pages 93-99, Suzhou, China. Association for Computational Linguistics.
- Xin Xin, Jinlong Li, and Zeqi Tan. 2021. N-ary constituent tree parsing with recursive semi-Markov model. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 2631-2642, Online. Association for Computational Linguistics.
- Naiwen Xue, Fei Xia, Fu-Dong Chiou, and Martha Palmer. 2005. The penn chinese treebank: Phrase structure annotation of a large corpus. Nat. Lang. Eng., 11(2):207–238.
 - D. Younger. 1967. Recognition and parsing of contextfree languages in time n³. Inf. Control., 10 : 189 – -208.

Zhang. 2020a. Syntax-aware opinion role labeling with dependency graph convolutional networks. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 3249–3258, Online. Association for Computational Linguistics.

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

- Yu Zhang, Zhenghua Li, and Min Zhang. 2020b. Efficient second-order TreeCRF for neural dependency parsing. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 3295-3305, Online. Association for Computational Linguistics.
- Yu Zhang, Houquan Zhou, and Zhenghua Li. 2020c. Fast and accurate neural CRF constituency parsing. In Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020, pages 4046-4053. ijcai.org.
- Yue Zhang and Stephen Clark. 2008. A tale of two parsers: Investigating and combining graph-based and transition-based dependency parsing. In Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing, pages 562-571, Honolulu, Hawaii. Association for Computational Linguistics.
- Xiaoqing Zheng. 2017. Incremental graph-based neural dependency parsing. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 1655–1665, Copenhagen, Denmark. Association for Computational Linguistics.
- Structure Grammar parsing on Penn Treebank. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2396–2408, Florence, Italy. Association for Computational Linguistics.