

CSSL: Contrastive Self-Supervised Learning for Dependency Parsing on Relatively Free Word Ordered and Morphologically Rich Low Resource Languages

Anonymous ACL submission

Abstract

Significant advancements have been made in the domain of dependency parsing, with researchers introducing novel architectures to enhance parsing performance. However, the majority of these architectures have been evaluated predominantly in languages with a fixed word order, such as English. Consequently, little attention has been devoted to exploring the robustness of these architectures in the context of relatively free word-ordered languages. In this work, we examine the robustness of graph-based parsing architectures on 7 relatively free word order languages. We focus on investigating essential modifications such as data augmentation and the removal of position encoding required to adapt these architectures accordingly. To this end, we propose a contrastive self-supervised learning method to make the model robust to word order variations. Furthermore, our proposed modification demonstrates a substantial average gain of 3.03/2.95 points in 7 relatively free word order languages, as measured by the Unlabelled/Labelled Attachment Score metric when compared to the best performing baseline.

1 Introduction

Morphologically rich languages (MRLs) tend to have sentences which follow a relatively free word order. Instead of relying on the word ordering, such languages prefer encoding the structural information of a sentence using inflectional morphology. Majority of the pretrained models tend to include (relative or absolute) position encoding in their pre-training stage, which may not be ideal for several of the MRLs. Moreover, simply dropping the position encoding of the encoder for such models, during fine-tuning, often would lead to sub-optimal performances for parsing tasks (Krishna et al., 2019; Ghosh et al., 2024). In this work, we propose a self-supervised contrastive learning based module that

makes a model agnostic to word order variations within a sentence.

We propose a novel Contrastive Self-Supervised Learning (CSSL) module, inspired by He et al. (2020), to accommodate variations in word order within the model architecture. Moreover, the modular nature of our approach allows for seamless integration with any encoder architecture, without necessitating alterations to pretraining decisions. In self supervised contrastive learning, for a given input, one needs to find positive samples, whose embedding level similarity with the input needs to be increased, and negative samples, whose embedding similarity with the original input, needs to be decreased. As shown in Figure 1, the original sentence serves as an anchor point, while its permutations represent positive examples, juxtaposed with randomly generated sentences serving as negative examples.



Figure 1: The Contrastive Loss minimizes the distance between an anchor (blue) and a positive (green), both of which have a similar meaning, and maximizes the distance between the anchor and a negative (red) of a different meaning.

The self-supervised contrastive learning objective aims to minimize the distance between positive examples and the anchor point, while simul-

taneously maximizing the distance from negative examples. In essence, this objective fosters the robustness of the encoder to accommodate word order variations. Our approach, to the best of our knowledge, is the first to use a contrastive learning technique for dependency parsing to overcome challenges caused by a lack of set word order and limited data resources.

MRLs rely less on word order and instead use morphological markers to encode structural information of a sentence. Given the comprehensive morphological marking system inherent in MRLs, the core semantic essence of the sentence remains unaltered, rendering the permuted counterpart as a suitable positive pairing for contrastive learning. Several MRLs have demonstrated that permutations of word order following weak projectivity generally retain semantic equivalence of the original (Sapir, 1921; Kulkarni et al., 2015; Kuboň et al., 2013; Ghosh et al., 2024).

Moreover, preference for certain word order topology in these languages is often not due to the limitations of the morphology, but are attributed to cognitive, psycho-linguistic, and information theoretic aspects of communication (Krishna et al., 2019; Clark et al., 2023; Dyer et al., 2023; Xu and Futrell, 2024). For instance, Sanskrit, a classical language, predominantly consists of sentences written as verses in its pre-classic and classic literature. Here, such sentences prefer to adherence to metrical constraints in prosody over any word ordering constructions, resulting in arbitrary word orderings (Krishna et al., 2020, §2). In our experiments, hence we treat the permutations of a given sentences as their semantic equivalents.

Substantial progress has been made in dependency parsing, including for low-resource languages and MRLs (Ji et al., 2021a,b; Dozat and Manning, 2017; Kulmizev et al., 2019), aimed at augmenting parsing efficacy. Our proposed approach is agnostic of the encoder architecture and does not necessitate the need for changes in pre-training. Moreover, our objective is to leverage recent advancements in parsing literature and further augment them by adding our CSSL module that would make these models more robust to word order variations. In this work, we start by examining the robustness of graph-based parsing architectures (Ji et al., 2019; Mohammadshahi and Henderson, 2020, 2021) on 7 relatively free word order languages. We believe, graph-based parsing architectures could be a natural choice to model flexible

word order. We then focus on investigating essential modifications such as data augmentation (Şahin and Steedman, 2018) and the removal of position encoding required to adapt these architectures accordingly. We finally show the efficacy of our approach on the best baseline (Mohammadshahi and Henderson, 2021, RNGTr) model by integrating CSSL with it and report an average performance gain of 3.03/2.95 points (UAS/LAS) improvement over 7 MRLs.

Our main contributions are as follows:

- We propose a novel contrastive self-supervised learning (CSSL) module to make dependency parsing robust for free word order languages.
- Empirical evaluations of CSSL module affirm its efficacy for 7 free word-ordered languages
- We demonstrate statistically significant improvements with an average gain of 3.03/2.95 points over the best baseline on 7 MRLs.

2 Contrastive Self-Supervised Learning

CSSL enables joint learning of representation, via contrastive learning, with the standard classification loss for dependency parsing. Here, via CSSL, we identify sentences which are word-level permutations of each other as similar sentences, and others as dissimilar sentences. The similar sentences are brought closer while pushing dissimilar examples apart (van den Oord et al., 2019; Tian et al., 2020). For a given input, when selecting a dissimilar sample, we choose a random sentence that clearly differs significantly from any permutation of the given sentence.

Formally, for a sentence X_i (anchor example), its representation should be similar to the permuted instance X_i^+ as permutation does not alter the meaning of a sentence belonging to MRL. However, the representation will differ from a random sentence X_i^- (negative example). Therefore, the distance between the appropriate representations of X_i and X_i^+ is expected to be small. Thus, we can develop a contrastive objective by considering (X_i, X_i^+) a positive pair and $N - 1$ negative pairs (X_i, X_i^-) :

$$\mathcal{L}_{\text{CSSL}} = -\log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_{i^+} / \tau)}{\sum_{a \in N} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)} \quad 158$$

where N represents a batch, \mathbf{z}_i represents the representation vector of the anchor sample, \mathbf{z}_i^+ denotes

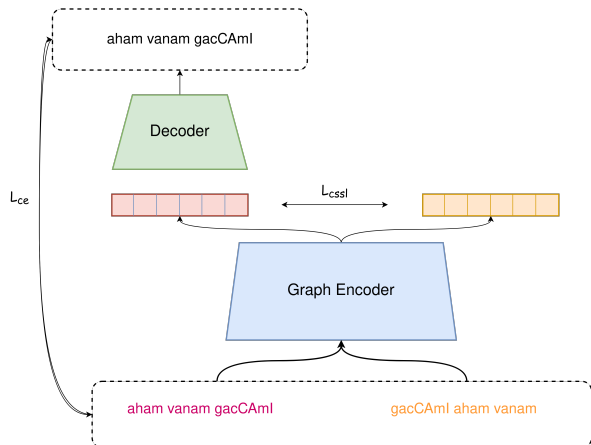


Figure 2: Schematic illustration of the proposed approach. Starting from an input sentence (bottom) (Translation: “I am going to the forest.”), two embeddings are produced: (1) original and (2) permuted sentence. Self-contrastive loss is imposed on the embeddings (**center**). A decoder uses a cross-entropy objective for predicting the dependency tree.

the representation vector for the positive sample (permuted sample), z_a represents the representation vector for a sample in the batch (N different samples), and τ is a temperature parameter that controls the concentration of the distribution. We employ pooled sentence embedding of the original and permuted sentences for CSSL loss. Therefore, our final loss is:

$$\mathcal{L} = \mathcal{L}_{\text{cssl}} + \mathcal{L}_{\text{ce}} \quad (1)$$

The classification loss L_{ce} is applied only to token-level labels of the original training input.

3 Experiment

3.1 Dataset and metric

As our primary benchmark dataset, we utilize the Sanskrit Treebank Corpus (Kulkarni, 2013, STBC). From STBC, we use a train and dev split of 2,800 and 1,000, respectively. Further, we employ a test set comprising 300 sentences, drawn from the classical Sanskrit work, *Śisūpāla-vadha* (Ryali, 2016).

Moreover, from Universal Dependencies (de Marneffe et al., 2021, UD-2.13), we choose 6 additional languages, namely, Turkish, Telugu, Gothic, Hungarian, Ancient Hebrew, and Lithuanian.¹ Please note that all the seven languages are chosen from diverse language families and are typologically diverse. Our experiments are

¹the statistics of each of the treebanks used for our experiments is mentioned in Table 4 in the Appendix.

primarily focused on a low-resource setting. Here, the largest training set size we use is of 3,435 sentences for Turkish. For Turkish, we simulate a low-resource scenario by considering from the Turkish-IMST treebank. We also experiment with English which is a fixed-ordered high-resource language. Here, we use a training set of 12,544 sentences. We use standard UAS/LAS metrics (McDonald and Nivre, 2011) for evaluation.

Model	UAS	LAS
G2GTr (Transition-based)	85.75	82.21
GNN (Graph-based)	88.01	82.8
RNGTr (Graph-based)	89.62	87.43
RNGTr (NoPos)	80.78	78.37
RNGTr (DA)	90.38	88.46
Prop. System (CSSL)	91.86	89.38
CSSL + DA	92.43	90.18

Table 1: Comparison of graph-based parsers on Sanskrit STBC dataset. We modify the best baseline RNGTr by integrating the proposed method (CSSL) to compare against variants, removing position encoding (NoPos) and data augmentation (DA). The best performances are bold-faced. The results (CSSL vs DA) and (CSSL vs DA+CSSL) are statistically significant as per the t-test with a p-value < 0.01 for the LAS metric.

Baselines: We utilize Mohammadshahi and Henderson (2020, G2GTr), a transition-based dependency parser. Furthermore, we explore Ji et al. (2019, GNN) a graph neural network-based model that captures higher-order relations in dependency trees. Finally, we examine Graph-to-Graph Non-Autoregressive Transformer proposed by Mohammadshahi and Henderson (2021, RNGTr) which iteratively refines arbitrary graphs through recursive operations.

Hyper-parameters: We implement our CSSL module in RNGTr architecture which uses a pre-trained mBERT model (110M parameters) from Huggingface transformers (Wolf et al., 2020). For RNGTr model, we use the same architecture as Mohammadshahi and Henderson (2020) with pre-trained mBERT as the encoder and an MLP and biaffine followed by softmax for the decoder. We adopt the RNGTr codebase with hyperparameter settings as follows: the batch size is 16, the learning rate as $2e-5$, the number of transformer blocks as 12 and for the decoder 2 Feed Forward Layers with dropout as 0.33 having bi-affine attention, and the remaining hyperparameters are the same.

Language	RNGTr		RNGTr + DA		RNGTr + CSSL	
	UAS	LAS	UAS	LAS	UAS	LAS
Turkish-IMST	72.86	71.99	74.18	72.96	78.21	74.69
Telugu-MTG	90.02	80.34	91.86	81.51	93.79	85.67
Gothic-POIEL	86.59	81.28	88.61	82.93	89.15	84.19
Hungarian-SZEGED	88.13	84.93	90.02	86.65	91.65	87.28
Ancient Hebrew-PTNK	90.76	86.42	91.43	87.12	92.35	88.68
Lithuanian-ALKSNIS	87.63	83.27	88.41	84.79	89.82	86.45
English-EWT	<i>92.08</i>	<i>90.23</i>	93.76	92.16	<i>93.19</i>	<i>90.71</i>

Table 2: Performance comparison on the RNGTr model, RNGTr + DA (Data Augmentation) and RNGTr + CSSL module. The best performances are bold-faced. Our results (CSSL) are statistically significant compared to both RNGTr and RNGTr + DA for each language as per the t-test with a p-value < 0.01 for the LAS metric

3.2 Results

In Table 1, we benchmark graph-based parsers on the Sanskrit STBC dataset. Our proposed contrastive loss module is standalone and could be integrated with any parser.² Thus, we modify the best baseline RNGTr by integrating the proposed method (CSSL) and comparing it against variants, removing position encoding (NoPos), and augmenting data augmentation (DA). Table 1 illustrates that the proposed framework adds a complementary signal making robust word order representations to RNGTr by improving 2.24/1.95 points in UAS/LAS scores. The performance significantly drops (8.8/9.0 UAS/LAS) when position embeddings are removed (vs. Pos kept) from RNGTr due to train-test mismatch in pretraining and fine-tuning steps. Moreover, our method outperforms data augmentation technique (Şahin and Steedman, 2018) by 1.48/0.92 points (UAS/LAS) when integrated with the RNGTr baseline. We integrate CSSL on top of an RNGTr+DA system and observe statistically significant improvements of 0.57/0.80 points (UAS/LAS), suggesting the proposed method complements the data-augmentation technique.

Results on multilingual experiments: In this section, we investigate the efficacy of CSSL module in multi-lingual settings. Table 2 reports results on 6 other morphologically rich languages in low-resource settings. Our approach averages 3.16/3.12 higher UAS/LAS scores than the usual cross-entropy-based RNGTr baseline. Our system outperforms the rotation-based DA technique with an average increase of 1.74/1.83 in UAS/LAS scores. Here, as expected, our proposed CSSL approach outperforms the standard RNGTr and DA

²Refer to Appendix A.1 for empirical evidence.

approaches for all the languages, except English. English is not an MRL and it relies heavily on configurational information of the words to understand sentence structure. The DA approach performs better by 0.57/1.45 UAS/LAS scores than our framework. However, it is interesting to note that CSSL still outperforms the RNGTr baseline by 1.11/0.48 UAS/LAS, possibly due to robustness of permutation invariant representation learning we employ in CSSL. As illustrated in Table 1, it is evident that combining CSSL with DA surpasses CSSL alone by approximately 0.5 points, exhibiting a 2-point enhancement over DA.

4 Conclusion

In this work, we investigated the robustness of graph-based parsing architectures across 7 languages characterized by relatively flexible word order. We introduced a self-supervised contrastive learning module aimed at making encoders insensitive to variations in word order within sentences. Additionally, the modular nature of our approach enables seamless integration with any encoder architecture without necessitating modifications to pretraining decisions. To the best of our knowledge, our approach represents the first utilization of contrastive learning techniques for dependency parsing to address challenges arising from variable word order in low-resource settings. Finally, we demonstrate the effectiveness of our approach by integrating it with the RNGTr model (Mohammadshahi and Henderson, 2021), reporting an average performance improvement of 3.03/2.95 points (UAS/LAS) across the 7 MRLs.

Limitations We could not evaluate on complete UD due to limited available compute resources (single GPU); hence, we selected 7 representative

291	languages for our experiments.	
292	Ethics Statement We do not foresee any eth-	
293	ical concerns with the work presented in this	
294	manuscript.	
295	References	
296	Thomas Hikaru Clark, Clara Meister, Tiago Pimentel,	
297	Michael Hahn, Ryan Cotterell, Richard Futrell, and	
298	Roger Levy. 2023. A Cross-Linguistic Pressure for	
299	Uniform Information Density in Word Order . <i>Trans-</i>	
300	<i>actions of the Association for Computational Linguis-</i>	
301	<i>tics</i> , 11:1048–1065.	
302	Marie-Catherine de Marneffe, Christopher D. Man-	
303	ning, Joakim Nivre, and Daniel Zeman. 2021. Uni-	
304	versal Dependencies . <i>Computational Linguistics</i> ,	
305	47(2):255–308.	
306	Timothy Dozat and Christopher D. Manning. 2017.	
307	Deep biaffine attention for neural dependency pars-	
308	ing . In <i>5th International Conference on Learning</i>	
309	<i>Representations, ICLR 2017, Toulon, France, April</i>	
310	<i>24-26, 2017, Conference Track Proceedings</i> . Open-	
311	Review.net.	
312	William Dyer, Charles Torres, Gregory Scontras, and	
313	Richard Futrell. 2023. Evaluating a Century of	
314	Progress on the Cognitive Science of Adjective Or-	
315	dering . <i>Transactions of the Association for Computa-</i>	
316	<i>tional Linguistics</i> , 11:1185–1200.	
317	Hongchao Fang, Sicheng Wang, Meng Zhou, Jiayuan	
318	Ding, and Pengtao Xie. 2020. Cert: Contrastive self-	
319	supervised learning for language understanding .	
320	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021.	
321	SimCSE: Simple contrastive learning of sentence em-	
322	beddings . In <i>Proceedings of the 2021 Conference</i>	
323	<i>on Empirical Methods in Natural Language Process-</i>	
324	<i>ing</i> , pages 6894–6910, Online and Punta Cana, Do-	
325	minican Republic. Association for Computational	
326	Linguistics.	
327	Poulami Ghosh, Shikhar Vashishth, Raj Dabre, and	
328	Pushpak Bhattacharyya. 2024. A morphology-based	
329	investigation of positional encodings .	
330	John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader.	
331	2021. DeCLUTR: Deep contrastive learning for un-	
332	supervised textual representations . In <i>Proceedings</i>	
333	<i>of the 59th Annual Meeting of the Association for</i>	
334	<i>Computational Linguistics and the 11th International</i>	
335	<i>Joint Conference on Natural Language Processing</i>	
336	<i>(Volume 1: Long Papers)</i> , pages 879–895, Online.	
337	Association for Computational Linguistics.	
338	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and	
339	Ross Girshick. 2020. Momentum contrast for unsu-	
340	pervised visual representation learning .	
341	Tao Ji, Yong Jiang, Tao Wang, Zhongqiang Huang, Fei	
342	Huang, Yuanbin Wu, and Xiaoling Wang. 2021a. A	
	unified encoding of structures in transition systems .	343
	In <i>Proceedings of the 2021 Conference on Empiri-</i>	344
	<i>cal Methods in Natural Language Processing</i> , pages	345
	4121–4133, Online and Punta Cana, Dominican Re-	346
	public. Association for Computational Linguistics.	347
	Tao Ji, Yong Jiang, Tao Wang, Zhongqiang Huang, Fei	348
	Huang, Yuanbin Wu, and Xiaoling Wang. 2021b.	349
	Word reordering for zero-shot cross-lingual struc-	350
	tured prediction . In <i>Proceedings of the 2021 Confer-</i>	351
	<i>ence on Empirical Methods in Natural Language Pro-</i>	352
	<i>cessing</i> , pages 4109–4120, Online and Punta Cana,	353
	Dominican Republic. Association for Computational	354
	Linguistics.	355
	Tao Ji, Yuanbin Wu, and Man Lan. 2019. Graph-based	356
	dependency parsing with graph neural networks . In	357
	<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	358
	<i>ciation for Computational Linguistics</i> , pages 2475–	359
	2485, Florence, Italy. Association for Computational	360
	Linguistics.	361
	Amrith Krishna, Bishal Santra, Ashim Gupta, Pavanku-	362
	mar Satuluri, and Pawan Goyal. 2020. A graph-based	363
	framework for structured prediction tasks in Sanskrit .	364
	<i>Computational Linguistics</i> , 46(4):785–845.	365
	Amrith Krishna, Vishnu Sharma, Bishal Santra, Aishik	366
	Chakraborty, Pavankumar Satuluri, and Pawan Goyal.	367
	2019. Poetry to prose conversion in Sanskrit as a lin-	368
	earisation task: A case for low-resource languages .	369
	In <i>Proceedings of the 57th Annual Meeting of the As-</i>	370
	<i>sociation for Computational Linguistics</i> , pages 1160–	371
	1166, Florence, Italy. Association for Computational	372
	Linguistics.	373
	Vladislav Kuboň, Markéta Lopatková, and Jiří	374
	Mírovský. 2013. A case study of a free word order .	375
	In <i>Proceedings of the 27th Pacific Asia Conference on</i>	376
	<i>Language, Information, and Computation (PACLIC</i>	377
	<i>27)</i> , pages 222–231, Taipei, Taiwan. Department of	378
	English, National Chengchi University.	379
	Amba Kulkarni. 2013. A deterministic dependency	380
	parser with dynamic programming for Sanskrit . In	381
	<i>Proceedings of the Second International Conference</i>	382
	<i>on Dependency Linguistics (DepLing 2013)</i> , pages	383
	157–166, Prague, Czech Republic. Charles Univer-	384
	sity in Prague, Matfyzpress, Prague, Czech Republic.	385
	Amba Kulkarni, Preethi Shukla, Pavankumar Satuluri,	386
	and Devanand Shukl. 2015. How free is free word	387
	order in sanskrit. <i>The Sanskrit Library, USA</i> , pages	388
	269–304.	389
	Artur Kulmizev, Miryam de Lhoneux, Johannes	390
	Gontrum, Elena Fano, and Joakim Nivre. 2019. Deep	391
	contextualized word embeddings in transition-based	392
	and graph-based dependency parsing - a tale of two	393
	parsers revisited . In <i>Proceedings of the 2019 Confer-</i>	394
	<i>ence on Empirical Methods in Natural Language Pro-</i>	395
	<i>cessing and the 9th International Joint Conference</i>	396
	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	397
	pages 2755–2768, Hong Kong, China. Association	398
	for Computational Linguistics.	399

400 Lajanugen Logeswaran and Honglak Lee. 2018. [An](#)
401 [efficient framework for learning sentence representa-](#)
402 [tions](#).

403 Ryan McDonald and Joakim Nivre. 2011. [Analyzing](#)
404 [and integrating dependency parsers](#). *Computational*
405 *Linguistics*, 37(1):197–230.

406 Alireza Mohammadshahi and James Henderson. 2020. [Graph-to-graph transformer for transition-based de-](#)
407 [pendency parsing](#). In *Proceedings of the 2020 Con-*
408 *ference on Empirical Methods in Natural Language*
409 *Processing: Findings*, pages 3278–3289, Online. As-
410 sociation for Computational Linguistics.

412 Alireza Mohammadshahi and James Henderson. 2021. [Recursive non-autoregressive graph-to-graph trans-](#)
413 [former for dependency parsing with iterative refine-](#)
414 [ment](#). *Transactions of the Association for Computa-*
415 *tional Linguistics*, 9:120–138.

417 Anupama Ryali. 2016. [Challenges in developing](#)
418 [sanskrit e-readers:semi-automatically using online](#)
419 [analyser saMsAdhanI:with special reference to](#)
420 [ŚiŚupĀlavadhā of mĀgha](#). In *Workshop on Bridging*
421 *4797the Gap Between Sanskrit CL Tools Manage-*
422 *ment of Sanskrit DL, ICON2016*.

423 Gözde Gül Şahin and Mark Steedman. 2018. [Data](#)
424 [augmentation via dependency tree morphing for low-](#)
425 [resource languages](#). In *Proceedings of the 2018 Con-*
426 *ference on Empirical Methods in Natural Language*
427 *Processing*, pages 5004–5009, Brussels, Belgium.
428 Association for Computational Linguistics.

429 Edward Sapir. 1921. An introduction to the study of
430 speech. *Language*, 1:15.

431 Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2020.
432 [Contrastive multiview coding](#).

433 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019.
434 [Representation learning with contrastive predictive](#)
435 [coding](#).

436 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
437 Chaumond, Clement Delangue, Anthony Moi, Pier-
438 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,
439 Joe Davison, Sam Shleifer, Patrick von Platen, Clara
440 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le
441 Scao, Sylvain Gugger, Mariama Drame, Quentin
442 Lhoest, and Alexander M. Rush. 2020. [Transformers:](#)
443 [State-of-the-art natural language processing](#). In
444 *Proceedings of the 2020 Conference on Empirical*
445 *Methods in Natural Language Processing: System*
446 *Demonstrations*, pages 38–45, Online. Association
447 for Computational Linguistics.

448 Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa,
449 Fei Sun, and Hao Ma. 2020. [Clear: Contrastive](#)
450 [learning for sentence representation](#).

451 Weijie Xu and Richard Futrell. 2024. [Syntactic depen-](#)
452 [dency length shaped by strategic memory allocation](#).
453 In *Proceedings of the 6th Workshop on Research in*
454 *Computational Linguistic Typology and Multilingual*

NLP, pages 1–9, St. Julian’s, Malta. Association for
455 Computational Linguistics. 456

457 Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang,
458 Wei Wu, and Weiran Xu. 2021. [Consert: A con-](#)
459 [trastive framework for self-supervised sentence rep-](#)
460 [resentation transfer](#).

461 Dejjiao Zhang, Shang-Wen Li, Wei Xiao, Henghui Zhu,
462 Ramesh Nallapati, Andrew O. Arnold, and Bing Xi-
463 ang. 2022. [Pairwise supervised contrastive learning](#)
464 [of sentence representations](#).

A Appendix 465

A.1 Integration of CSSL with another encoder 466

467 The modular nature of CSSL framework allows for
468 seamless integration with any encoder architecture,
469 without necessitating alterations to pretraining de-
470 cisions. We have shown its effectiveness for the
471 best-performing baseline. We are also showing re-
472 sults with one more baseline (for Sanskrit). Our
473 supplementary results indicate that activating con-
474 trastive loss for the G2GTr baseline on the STBC
475 treebank for Sanskrit leads to an approximate 2-
476 point enhancement in performance measured by
477 UAS/LAS.

	CE		CSSL	
	UAS	LAS	UAS	LAS
G2GTr	87.16	85.68	89.05	87.05

Table 3: Contrastive Loss with G2GTr on STBC dataset.

A.2 Treebank Statistics 478

479 Table 4 provides the detailed statistics for the lan-
480 guages used in the experiments. 481

A.3 Related Work 482

483 Contrastive learning has been the pinnacle of re-
484 cent successes in sentence representation learn-
485 ing. In order to optimize the appropriately de-
486 signed contrastive loss functions, (Gao et al., 2021;
487 Zhang et al., 2022) uses the entailment sentences
488 in NLI as positive pairs, significantly improving
489 upon the prior state-of-the-art results. To this end,
490 a number of methods have been put forth recently
491 in which the augmentations are obtained through
492 back-translation (Fang et al., 2020), dropout (Yan
493 et al., 2021; Gao et al., 2021), surrounding con-
494 text sampling (Logeswaran and Lee, 2018; Giorgi
495 et al., 2021), or perturbations carried out at dif-
496 ferent semantic-level (Wu et al., 2020; Yan et al.,
497 2021).

Treebank	Language Family	train	dev	test
Sanskrit-STBC	Indo-Aryan	2,800	1,000	300
UD-Turkish_IMST	Turkic	3,435	1,100	1,100
UD-Gothic_Proeil	Germanic	3,387	985	1,029
UD-Telugu_MTG	Dravidian	1,051	131	146
UD-Hungarian_Szeged	Uralic	910	441	449
UD-Ancient_Hebrew_PTNIK	Semitic	730	439	410
UD-Lithuanian_ALKSNIS	Baltic	2,341	617	684
UD-English_EWT	Roman	12,544	2,001	2,077

Table 4: Treebank Statistics. The number of sentences in train, dev and test for each language.