

	w2v-wiki			ELMo			No embeddings		
Model	Pearson's $r$	$R^2$	$R^2_{adj}$	Pearson's $r$	$R^2$	$R^2_{adj}$	Pearson's $r$	$R^2$	$R^2_{adj}$
Embeddings	0.430	0.064	0.046	0.404	0.063	0.044	-	-	-
+MI	0.464	0.158	0.138	0.458	0.153	0.133	0.376	0.083	0.079
+D.O.	0.575	0.245	0.227	0.530	0.202	0.183	0.301	0.029	0.025
+diag.	0.449	0.125	0.104	0.440	0.138	0.117	0.268	0.027	0.023
+MI. +D.O.	0.586	0.278	0.258	0.607	0.297	0.277	0.466	0.165	0.158
+MI. +diag.	0.515	0.205	0.182	0.512	0.193	0.170	0.436	0.141	0.134
+diag. +D.O.	0.572	0.265	0.245	0.535	0.224	0.203	0.392	0.114	0.107
+all	0.624	0.330	0.309	0.609	0.304	0.281	0.516	0.215	0.206

Table 4: Ablation results from the two best models and a non-embedding features only-model (10-fold cross validation).

Test tasks	metric	Sentence encoder pretraining tasks			
		Random	Arg	Arg fullsent	Arg fullsent 3-way
SRL-CoNLL2005 (WSJ)	$F_1$	81.7	<b>83.9***</b>	<b>84.7***</b>	<b>84.5***</b>
SRL-CoNLL2012 (OntoNotes)	$F_1$	77.3	<b>80.2***</b>	<b>80.4***</b>	<b>80.7***</b>
PP attachment (?)	$acc.$	87.5	87.6	88.2	87.0

Table 5: Gains over random initialization from pretraining sentence encoders on PP argumenthood tasks. (\*\*\*) :  $p < .001$

may have impacted the results, and ran additional experiments using higher-dimensional embeddings that were publicly available. Higher-dimensional embeddings did indeed lead to performance improvements, even though the actual inputs given to the models were all PCA-reduced to  $d = 5$ . From this observation, we could further improve upon the initial ELMo results. Results from the best model (w2v-wiki) are given in addition to the set of results using the same embeddings as the models in Experiment 1. This model uses 1000- $d$  word2vec features with additional interaction features (multiplicative, subtractive) that improved dev set performance.

**Ablation** We conducted ablation experiments with the two best-performing models to examine the contribution of non-embedding features discussed in Section 4.2. Table 4 indicates that any linguistic feature contributes positively towards performance, with the direct object feature helping both word2vec and ELMo models the most. This supports our initial hypothesis that adding the direct object feature would help reduce noise in the data. When only the linguistic features are used without embeddings as base features, mutual information is the most informative. This suggests that there is some (but not complete) redundancy in information captured by word embeddings and mutual information. The diagnostics feature is informative but is a comparatively weak predictor, which aligns with the current state of diagnostic acceptability tests—they are sometimes useful but not always, especially with respect to syntactically oblique items such as PPs. This behavior of the diagnostics predictor adds credibility to our data collection protocol.

## 5 Why Is This a Useful Standalone Task?

In motivating our tasks, we suggested that PP argumenthood information could improve existing NLP task performance such as SRL and parsing. We investigate whether this is a grounded claim by testing two separate hypotheses: (1) whether the task is indeed useful, and if so, (2) whether it is useful as a standalone task. We leave the issue of gradient argumenthood to future work for now, since the dataset is currently small and the notion of gradient argumenthood is not yet compatible with formulations of many NLP tasks.

### 5.1 Improving Representations with Pretraining

We first test the utility of the binary argumenthood task in improving performances on existing NLP tasks. We selected three tasks that may benefit from PP argumenthood information: SRL on Wall Street Journal (WSJ) data (CoNLL 2005; ? ?), SRL on OntoNotes Corpus (CoNLL 2012 data; ? ?)<sup>5</sup>, and PP attachment disambiguation on WSJ (?).

We follow ? (?)’s setup to pretrain and evaluate sentence encoders<sup>6</sup>. If learning to make correct PP argumenthood distinction teaches models knowledge that is generalizable to the new tasks, the classifier trained on top of the fixed-weights encoder will perform better on those tasks compared to a classifier trained on top of an encoder with randomly initialized weights. Improvements over the randomly initialized setup from pretraining on our main PP argumenthood task (Arg) and its full-sentence variants (Arg fullsent and Arg fullsent 3-way; see Section 3.1 for details) are shown in Table 5. Only statistically significant ( $p < .05$ ) improvements over the random encoder model are bolded, with significance levels calculated via Approximate Randomization (?) ( $R = 1000$ ). The models trained on PP ar-

<sup>5</sup>Tasks are labeling only, as described in ? (?).

<sup>6</sup>[github.com/jsalt18-sentence-repl/jiant](https://github.com/jsalt18-sentence-repl/jiant)

gumenthood tasks perform significantly better than the random initialization model in both SRL tasks, which supports our initial claim that argumenthood tasks can be useful for SRL. Although not all errors made by the models were interpretable, we found interesting improvements such as the model trained on the PP argumenthood task being slightly more accurate than the random initialization model on AM-DIR, AM-LOC, and AM-MNR labels. However, we did not observe significant improvements for the PP attachment disambiguation task. We speculate that since the task as formulated in ? (?) requires the model to understand PP dependents of NPs as well as VPs, our tasks that focus on verbal dependents may not provide the full set of linguistic knowledge necessary to solve this task. Nevertheless, our models are not significantly worse than the baseline, and the accuracy of the Arg fullsent model (88.2%) was comparable to a model that uses an encoder directly trained on PP attachment (88.7%).

Secondly, we discuss whether it is indeed useful to formulate PP argumenthood prediction as a separate task. The questions that need to be answered are (1) whether it would be the same or better to use a different pretraining task that would provide similar information (e.g., PP attachment disambiguation), and (2) whether the performance gain can be attributed to simply seeing more datapoints at train time rather than to the regularities we hope the models would learn through our task. Table 6 addresses both questions; we compare models pretrained on argumenthood tasks to a model pretrained directly on the PP attachment task listed in Table 5. All models trained on PP argumenthood prediction outperform the model trained on PP attachment, despite the fact that the latter has advantage for SRL2005 since the tasks share the same source text (WSJ). Furthermore, the variance in the sizes of the datasets indicates that the reported performance gains cannot solely be due to the increased number of datapoints seen during training.

	PP att.	Arg	Arg full	Arg full 3-way
Size	32k	19k	58k	87k
SRL2005	80.2	<b>83.9***</b>	<b>84.7***</b>	<b>84.5***</b>
SRL2012	79.8	<b>80.2***</b>	<b>80.3***</b>	<b>80.7***</b>

Table 6: Comparison against using PP attachment directly as a pretraining task (\*\*\*:  $p < .001$ ).

## 6 Conclusion

We have proposed two different tasks—binary and gradient—for predicting PP argumenthood, and reported results on each using four different types of word embeddings as base predictors. We obtain 95.5 accuracy and 95.4  $F_1$  in the binary classification task with BiLSTM and ELMo, and  $r = 0.624$  for the gradient human judgment prediction task. Our overall contribution is threefold: first, we have demonstrated that a principled prediction of both binary and gradient argumenthood judgments is possible with informed selection of lexical features; second, we justified the utility of

our binary PP argumenthood classification as a standalone task by reporting performance gains on multiple end-tasks through encoder pretraining. Finally, we have conducted a proof-of-concept study with a novel gradient argumenthood prediction task, paired with a new public dataset<sup>7</sup>.

### 6.1 Future Work

The pretraining approach holds promise in understanding and improving neural network models of language. Especially for end-to-end models, this method has an advantage over architecture engineering or hyperparameter tuning in terms of interpretability. That is, we can attribute the source of the performance gain on end tasks to the knowledge necessary to do well on the pretraining task. For instance, in Section 5 we can infer that that knowing how to make correct PP argumenthood distinction helps models encode representations that are more useful for SRL. Furthermore, we believe it is important to contribute to the recent efforts for designing better probing tasks to understand what machines really know about natural language (as opposed to directly taking downstream performances as metrics of better models). We hope to scale up our preliminary experiments and will continue to work on developing a set of linguistically informed probing and pretraining tasks for higher-quality, better-generalizable sentence representations.

### Acknowledgments

This research was supported by NSF INSPIRE (BCS-1344269) and DARPA LORELEI. We thank C. Jane Lutken, Rachel Rudinger, Lilia Rissman, Géraldine Legendre and the anonymous reviewers for their constructive feedback. Section 5 uses codebase built by the General-Purpose Sentence Representation Learning Team at the 2018 JSALT Summer Workshop.

Dolor inventore exercitationem molestias, officia rerum id repudiandae voluptas ratione natus corporis?Minima ratione omnis pariat, consequuntur ducimus nihil beatae sed?Molestiae maxime modi perferendis laudantium accusamus tempore, quas laboriosam officiis fugit eius obcaecati, voluptates magnam alias nobis dolores odio, aliquam vitae pariat reiciendis?Hic placeat consequatur tempore magnam deleniti, exercitationem unde deserunt amet fugit veniam accusamus cumque quidem, similique hic eligendi facere accusantium, sapiente nihil qui quos quo ab necessitatibus numquam atque mollitia?Quidem vero ipsa excepturi delectus nisi deserunt, aspernatur minus atque iusto, odio pariat modi facilis fugiat voluptatum reiciendis nihil nam cum neque expedita, amet saepe distinctio vero officiis asperiores culpa, sed reiciendis nam hic expedita soluta dicta?Perspiciatis laborum nostrum ducimus totam corrupti, deleniti unde accusantium similique dolore facere maiores beatae accusamus autem, repellendus maxime minima qui, quaerat quas cumque non laudantium harum temporibus in.Placet tempore nihil perspiciatis omnis praesentium dolores dolorum deleniti, unde similique error voluptas aliquid qui neque distinctio repellendus officiis?Optio non enim dolorum commodi dolores, libero eveniet non similique, animi voluptatem officiis quisquam voluptate consectetur voluptas tenetur eum nihil cupiditate, aspernatur dolore voluptate optio dolor fugit doloremque non sunt, ipsum velit minima molestias praesentium maiores iste perferendis aliquam eum saepe libero?Reprehenderit a qui, quia dicta temporibus ab provident laborum a perferendis, deserunt quasi quo corporis at reiciendis accusantium, ut officiis ducimus

<sup>7</sup>To be released at: [decomp.io](https://decomp.io)