

# The CLaC Discourse Parser at CoNLL-2015

Majid Laali\* Elnaz Davoodi† Leila Kosseim‡  
 Department of Computer Science and Software Engineering,  
 Concordia University, Montreal, Canada

## I. Summary

- Focus on the treatment of explicit discourse relations.
- Overall F<sub>1</sub> measure of 17.38%, ranking in 6<sup>th</sup> place out of the 17 parsers submitted to CoNLL 2015.
- Architecture similar to the End-to-End Discourse parser.
- The CLaC Discourse parser is based on the UIMA framework.
- Uses ClearTK to add machine learning functionality.
- Written in Java and its source code is available at “https://github.com/mjlaali/CLaCDiscourseParser.git”.

## IV. Argument Labeler

### Algorithm:

- Calculates the *Connective-Root path nodes*, the nodes that appear in the path from the discourse connective to the root of the sentence.
- Labels all constituents that are directly connected to one of the *Connective-Root path nodes* with ‘part of ARG1’, ‘part of ARG2’ or ‘NON’.
- Uses a classifier with nine features (i.e. F<sub>1</sub>-F<sub>9</sub>).
- Merges all constituents which were tagged as part of ARG1 or as part of ARG2 to obtain the actual boundaries of ARG1 and ARG2.
- If no constituent was labeled as a part of ARG1, the whole text of the previous sentence is considered as ARG1.

### Results:

- The F<sub>1</sub> score of the *Argument Labeler*:

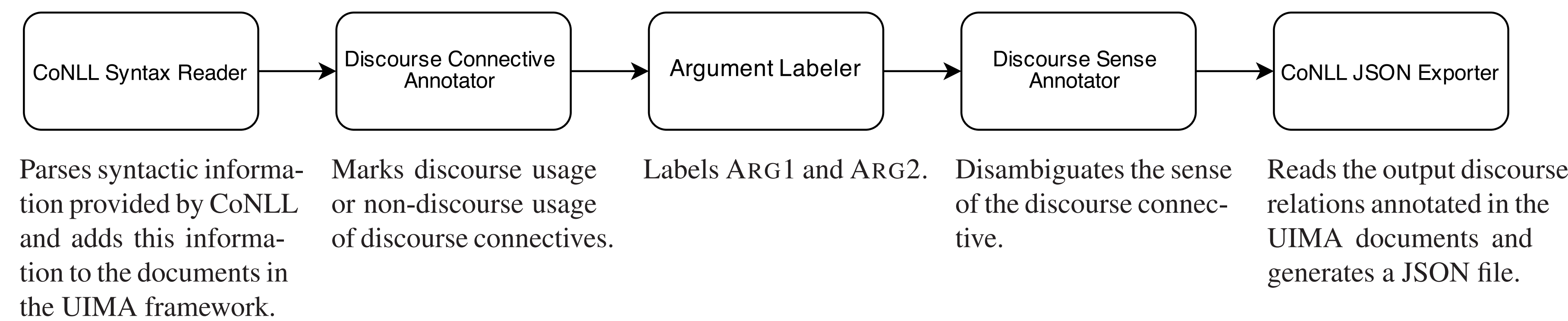
	Arg1	Arg2
Best Result	49.68%	74.29%
<b>CLaC Discourse Parser</b>	<b>45.18%</b>	<b>69.18%</b>
Average	30.77%	50.91%
Std. deviation	15.31%	20.58%

- Results show that the identification of ARG1 is more difficult than ARG2.

### Error Analysis:

- *Attribute spans*:
  - But the RTC also requires “working” capital to maintain the bad assets of thrifts that are sold until the assets can be sold separately.
- Subordinate and coordinate clauses:
  - We would have to wait until we have collected on those assets before we can move forward.

## II. Architecture

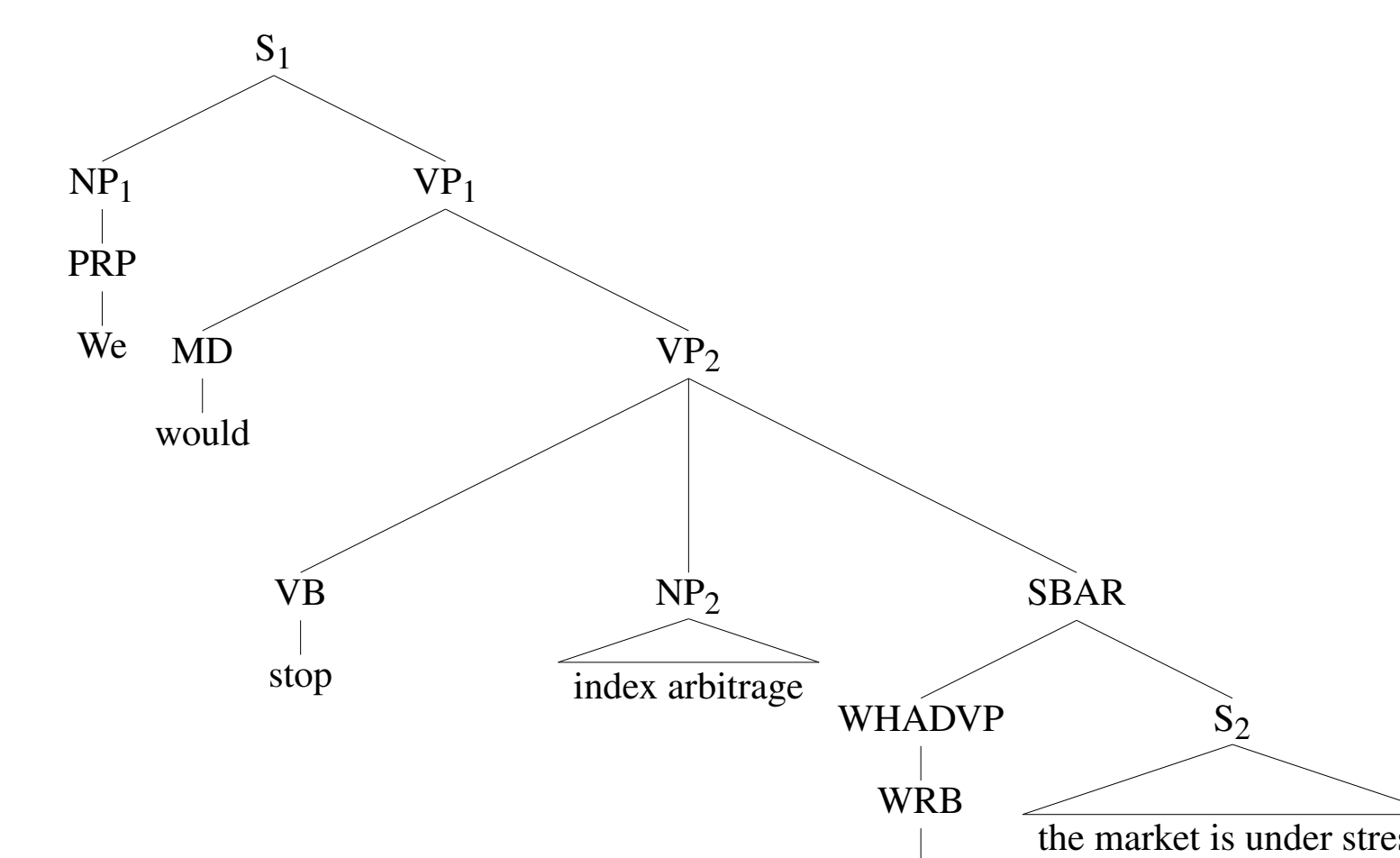


## V. Features

Category	Description
Connective Features	F <sub>1</sub> : Discourse connective text in lower-case.
	F <sub>2</sub> : Categorization of the case of the connective: <i>all lowercase</i> , <i>all uppercase</i> and <i>initial uppercase</i>
	F <sub>3</sub> : Highest node in the parse tree that covers the connective words but nothing more
	F <sub>4</sub> : Parent of <i>SelfCat</i>
	F <sub>5</sub> : Left sibling of <i>SelfCat</i>
	F <sub>6</sub> : Right sibling of <i>SelfCat</i>
Syntactic Node Features	F <sub>7</sub> Path from the node to the <i>SelfCat</i> node in the parse tree
	F <sub>8</sub> : Context of the node in the parse tree. The context of a node is defined by its label the label of its parent, the label of left and right sibling in the parse tree.
	F <sub>9</sub> : Position of the node relative to the <i>SelfCat</i> node: <i>left</i> or <i>right</i>

## VI. Example

We would stop index arbitrage when the market is under stress.



F <sub>1</sub> = <i>when</i>	F <sub>2</sub> = <i>all lowercase</i>
F <sub>3</sub> = <i>WRB</i>	F <sub>4</sub> = <i>WHADVP</i>
F <sub>5</sub> = <i>null</i>	F <sub>6</sub> = <i>S</i>
F <sub>7</sub> = <i>S ↑ SBAR ↓ WHADVP</i>	F <sub>8</sub> = <i>S-SBAR-WHADVP-null</i>
F <sub>9</sub> = <i>left</i>	

## VIII. Overall Results

- The F<sub>1</sub> scores of the CLaC discourse parser and the individual performance of its components:

	Discourse Connective Classifier	Argument Labeler	Discourse Parsing (explicit only)	Discourse Parsing (explicit and implicit)
Best Result	91.86%	41.35%	30.58%	24.00%
<b>CLaC Discourse Parser</b>	<b>90.19%</b>	<b>36.60%</b>	<b>27.32%</b>	<b>17.38%</b>
Average	74.20%	23.89%	18.28%	13.25%
Standard deviation	23.24%	13.01%	9.93%	6.41%

## III. Discourse Connective Annotator

### Algorithm:

- Searches the input texts for terms that match a predefined list of discourse connectives (was built solely from the CoNLL training dataset).
- Checks each match of discourse connective to see if it occurs in discourse usage or not.
  - Uses a binary classifier with six features (i.e. F<sub>1</sub>-F<sub>6</sub>)

### Results:

- F<sub>1</sub> = 90.19%.

## VII. Discourse Sense Annotator

- Uses the naïve approach that labels each discourse connective with its most frequent. The most frequent relation for discourse connectives is mined from the CoNLL training dataset.

## IX. Conclusion

- CLaC Discourse Parser was developed from scratch for CoNLL 2015.
- 3 person-month effort focused on *Discourse Connective Classification* and *Argument Labeler*.
- Naïve approach for sense labelling and consider only explicit relations.
- Yet, good results.

## References

- Steven Bethard, Philip Ogren, and Lee Becker. ClearTK 2.0: Design patterns for machine learning in UIMA. LREC, 2014.
- David Ferrucci and Adam Lally. UIMA: An architectural approach to unstructured information processing in the corporate research environment. *Natural Language Engineering*, 10(3-4):327–348, 2004.
- Fang Kong, Hwee Tou Ng, and Guodong Zhou. A Constituent-Based Approach to Argument Labeling with Joint Inference in Discourse Parsing. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 68–77, Doha, Qatar, October 2014.
- Ziheng Lin, Hwee Tou Ng, and Min-Yen Kan. A PDTB-styled end-to-end discourse parser. *Natural Language Engineering*, 20(02):151–184, 2014.