### A Transition-Based Directed Acyclic Graph Parser for Universal Conceptual Cognitive Annotation

Daniel Hershcovich, Omri Abend and Ari Rappoport



Technion June 29, 2017

#### Linguistic Annotation Schemes

• Dependency grammar (Tesniére, 1959; Nivre, 2005):

Purely syntactic analysis.

Semantic representation (Abend and Rappoport, 2017):

• Semantic dependencies (Oepen et al., 2015):

Coupled with syntactic representation.

• AMR (Banarescu et al., 2013), UCCA (Abend and Rappoport, 2013):

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

 $\dots \text{ showered } = \boxed{\dots \text{ took a shower}}$  $\dots \text{ 's war against crime } = \boxed{\dots \text{ fights crime}}$ 

#### **Dependency Parsing**

- Graph representation of syntactic structure.
- Bilexical tree: edges are only between tokens.
- Fast and accurate parsers (e.g. *transition-based*).



#### **Dependency Parsing**

- Graph representation of syntactic structure.
- Bilexical tree: edges are only between tokens.
- Fast and accurate parsers (e.g. *transition-based*).



Non-projectivity (discontinuity) is a challenge (Nivre, 2009).



#### Semantic Dependency Parsing

- Representation of predicate-argument relationships.
- Bilexical graph: allows reentrancy (and discontinuity).
- Various formalisms.



DELPH-IN MRS-derived bi-lexical dependencies (DM).

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

#### Semantic Dependency Parsing

- Representation of predicate-argument relationships.
- Bilexical graph: allows reentrancy (and discontinuity).
- Various formalisms.



DELPH-IN MRS-derived bi-lexical dependencies (DM).



Prague Dependency Treebank tectogrammatical layer (PSD).

# Universal Conceptual Cognitive Annotation (UCCA) — primary edge



#### The UCCA Semantic Representation Scheme

Cross-linguistically applicable (Abend and Rappoport, 2013). Stable in translation (Sulem et al., 2015).



#### UCCAApp

## Rapid and intuitive annotation interface (Abend et al., 2017). Usable by non-experts. http://ucca-demo.cs.huji.ac.il

Linker (L) i Ground (G) i Participant (A) i	William Bradley Pitt was born in Shawnee , Oklahoma , to William Alvin Pitt , who ran a trucking company , and Jane Etta ( née Hillhouse ), school counsellor . The family scon moved to Springfield , Missouri , where he lived together with his younger siblings , Douglas ( born 1966 Julie Neal ( born 1969 ). Born into a conservative household , he was raised as Southern Baptist , but has since stated Wat he does not " hi great relationship with religion " and that he " oscillates between agnosticism and atheism ." Pitt has described Springfield as " Mark Twain o , Jesse James country", having arown up with " a tot of hills . a lot of fakes ".	a 5) and ave a country							
State (S) i									
Process (P) i	1 H William Bradley Pitt was born in Shawnee , Oklahoma + C	9 F 🗙							
Adverbial (D) i	1-1 A William Bradley Pitt + 🛡	9 F X							
Time (T) i	1-2 F was + 🗉	J F X							
Center (C) i									
Elaborator (E) i	1-3 P born + 🛡	) F 🗙							
Connector (N) i	E 1-4 A in Shawnee , Oklahoma + C	9 F 🗙							
Relator (R) i	1-4-1 <b>R</b>   in + 🛡	9 f <b>x</b>							
Uncertain (UNC) i	1-4-2 C UNA Shawnee , Oklahoma + C	) F <b>x</b>							
Unanalyzable (UI i									
Function (F) i									

#### HUME

## UCCA facilitates semantics-based human evaluation of machine translation (Birch et al., 2016). http://ucca.cs.huji.ac.il/mteval

Für	leic	ht fettleibige Diabetiker kann die Gewichtsschaden. OP <mark>hilfreich</mark> sein	
Forn	nildly	y obese diabetics , weight loss surgery may be helpfu	A B 🌒 🗭 🗭 🗙
Θ	l Fo	r mildly obese diabetics . veight lass augery may be helpful Für leicht fettleibige Diabetiker kann die Gewichtsschaden - 02 hilfreich sein	A B 🔵 🔴 🗙
0	-	For mildly obese diabetics Fir leicht fettleibige Diabetiker kann	A B 🔵 🔴 🗙
	ſ	For <i>Riv</i>	🧧 🔴 🗮 🗙
		nuldiy leicht	<b>•</b> • ×
		Cobese fettleibige	<b>⊡ ● ●</b> ×
	L	diabetics Diabetiker kann	<b>⊡</b> ●●×
6		weight loss surgery Gewichtsschaden -	A B 🔵 🔴 🗙
	e	ewight loss Gewichtsschaden	A B 😑 😑 🗮 🗙
		weight Gewichtsschaden	● ● ×
		loss	● ● ● ×
	L	surgery -	● ● ×
		may	● ● ● ×
6	•	be helpful hilfeich sein	A B 💽 🔴 🗙
		be sein	● ● ● ×
	L	helpful hilfreich sein	● ● ● ×

#### Graph Structure

UCCA forms a directed acyclic graph (DAG). Tokens are terminals. Structural properties:

1. Non-terminal nodes



You want to take a long bath

- 日本 本語 本 本 田 本 王 本 田 本

#### Graph Structure

UCCA forms a directed acyclic graph (DAG). Tokens are terminals. Structural properties:

- 1. Non-terminal nodes
- 2. Reentrancy



You want to take a long bath

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ○ ○ ○

### Graph Structure

UCCA forms a directed acyclic graph (DAG). Tokens are terminals. Structural properties:

- 1. Non-terminal nodes
- 2. Reentrancy
- 3. Discontinuity



You want to take a long bath

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ○ ○ ○

• Parse text  $w_1 \dots w_n$  to graph  $G = (V, E, \ell)$  incrementally.

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Classifier determines transition to apply at each step.
- Trained by an oracle based on gold-standard graph.

- Parse text  $w_1 \dots w_n$  to graph  $G = (V, E, \ell)$  incrementally.
- Classifier determines transition to apply at each step.
- Trained by an oracle based on gold-standard graph.

Initial state:

stack

buffer

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

You want to take a long bath	You
------------------------------	-----

- Parse text  $w_1 \dots w_n$  to graph  $G = (V, E, \ell)$  incrementally.
- Classifier determines transition to apply at each step.
- Trained by an oracle based on gold-standard graph.

Initial state:

stack

buffer

	You	want	to	take	a	long	bath
--	-----	------	----	------	---	------	------

TUPA transitions:

{Shift, Reduce, Node<sub>X</sub>, Left-Edge<sub>X</sub>, Right-Edge<sub>X</sub>, Left-Remote<sub>X</sub>, Right-Remote<sub>X</sub>, Swap, Finish}

Support non-terminal nodes, reentrancy and discontinuity.

stack						buffer
• You	want	to	take	a	long	bath
graph						

#### $\Rightarrow \mathrm{Right}\text{-}\mathrm{Edge}_A$





 $\Rightarrow$  Swap



#### $\Rightarrow \mathrm{Right}\text{-}\mathrm{Edge}_\mathrm{P}$



 $\Rightarrow$  Reduce







 $\Rightarrow \mathrm{Node}_\mathrm{F}$ 



 $\Rightarrow$  Reduce







 $\Rightarrow \mathrm{Node}_{\mathrm{C}}$ 



 $\Rightarrow$  Reduce



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○



 $\Rightarrow \mathrm{Right}\text{-}\mathrm{Edge}_\mathrm{P}$ 





 $\Rightarrow \mathrm{Right}\text{-}\mathrm{Edge}_\mathrm{F}$ 



 $\Rightarrow$  Reduce




$\Rightarrow$  Swap



 $\Rightarrow \mathrm{Right}\text{-}\mathrm{Edge}_\mathrm{D}$ 



◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ 少々ぐ

 $\Rightarrow$  Reduce



 $\Rightarrow$  Swap



#### $\Rightarrow \operatorname{Right-Edge}_A$



 $\Rightarrow$  Reduce



 $\Rightarrow$  Reduce



 $\Rightarrow$  Shift



 $\Rightarrow$  Shift



#### $\Rightarrow \text{Left-Remote}_A$



 $\Rightarrow$  Shift



◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ 少々ぐ

#### $\Rightarrow$ Right-Edge<sub>C</sub>



▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ ▲圖 - 釣��

 $\Rightarrow$  Finish



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○

#### Word Embeddings

Represent discrete features by dense vectors (Goldberg, 2016).



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Word Embeddings

## Feed-Forward Neural Network (MLP)

Learns representation and classification by optimizing weights.



#### Recurrent Neural Network (RNN)

Applied to sequences, updates state given input and previous state.



▲ロ▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへで

#### Long Short-Term Memory (LSTM)

Memory cell to avoid vanishing gradients in RNNs.



◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

#### **TUPA Model**

Greedy parsing, experimenting with three classifiers:

Sparse	Perceptron with sparse features.
MLP	Embeddings $+$ feedforward NN classifier.
BiLSTM	Embeddings + deep bidirectional LSTM + MLP.

Features: words, POS, syntactic dependencies, existing edge labels from the stack and buffer + parents, children, grandchildren; ordinal features (height, number of parents and children)

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・





◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ 三臣 - ∽ � � �

#### **Experimental Setup**

- UCCA Wikipedia corpus  $\begin{pmatrix} train \\ 4268 \\ + 454 \\ + 503 \\ train \\ 503 \\ sentences \end{pmatrix}$ .
- Out-of-domain: English part of English-French parallel corpus, *Twenty Thousand Leagues Under the Sea* (506 sentences).



#### **Evaluation**

Comparing graphs over the same sequence of tokens,

- Match edges by their terminal yield and label.
- Calculate labeled precision, recall and F1 scores.
- Separate primary and remote edges.



#### Results

		Primary	ý	Remote		
	LP LR LF			LP	LR	LF
Sparse	64.5	63.7	64.1	19.8	13.4	16
MLP	65.2	64.6	64.9	23.7	13.2	16.9
BiLSTM	74.4	72.7	73.5	47.4	51.6	49.4

Results on the Wiki test set.

#### Results

		Primary	y	Remote		
	LP LR LF			LP	LR	LF
Sparse	64.5	63.7	64.1	19.8	13.4	16
MLP	65.2	64.6	64.9	23.7	13.2	16.9
BiLSTM	74.4	72.7	73.5	47.4	51.6	49.4

Results on the Wiki test set.

		Primary	y	Remote			
	LP LR LF		LP	LR	LF		
Sparse	59.6	59.9	59.8	22.2	7.7	11.5	
MLP	62.3	62.6	62.5	20.9	6.3	9.7	
BiLSTM	68.7	68.5	68.6	38.6	18.8	25.3	

Results on the 20K Leagues out-of-domain set.

#### **Bilexical Graph Approximation**

No existing UCCA parsers  $\Rightarrow$  compare to bilexical parsers:

- 1. Convert UCCA to bilexical dependencies.
- 2. Train bilexical parsers and apply to test sentences.
- 3. Reconstruct UCCA graphs and compare with gold standard.



Bilexical DAG approximation.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

#### Baselines

Bilexical DAG parsers:

- DAGParser (Ribeyre et al., 2014): transition-based.
- TurboParser (Almeida and Martins, 2015): graph-based.

Tree parsers (all transition-based):

- MaltParser (Nivre et al., 2007): bilexical tree parser.
- LSTM Parser (Dyer et al., 2015): bilexical tree parser.
- UPARSE (Maier, 2015): allows non-terminals, discontinuity.



Bilexical tree approximation.

★ Ξ → Ξ

Conversion to trees is just removing remote edges.

#### Results

 $\mathsf{TUPA}_{\mathsf{BiLSTM}}$  obtains the highest F-scores in all metrics:

	Primary			Remote		
	LP	LR	LF	LP	LR	LF
<b>TUPA</b> Sparse	64.5	63.7	64.1	19.8	13.4	16
TUPA <sub>MLP</sub>	65.2	64.6	64.9	23.7	13.2	16.9
TUPA <sub>Bilstm</sub>	74.4	72.7	73.5	47.4	51.6	49.4
Bilexical DAG			(91)			(58.3)
DAGParser	61.8	55.8	58.6	9.5	0.5	1
TurboParser	57.7	46	51.2	77.8	1.8	3.7
Bilexical tree			(91)			_
MaltParser	62.8	57.7	60.2	-	-	-
LSTM Parser	73.2	66.9	69.9	-	_	-
Tree			(100)			-
UPARSE	60.9	61.2	61.1	-	-	-

Results on the Wiki test set.

#### Results

Similar on out-of-domain test set:

	Primary			Remote		
	LP	LR	LF	LP	LR	LF
TUPA <sub>Sparse</sub>	59.6	59.9	59.8	22.2	7.7	11.5
TUPA <sub>MLP</sub>	62.3	62.6	62.5	20.9	6.3	9.7
TUPA <sub>Bilstm</sub>	68.7	68.5	68.6	38.6	18.8	25.3
Bilexical DAG			(91.3)			(43.4)
DAGParser	56.4	50.6	53.4	-	0	0
TurboParser	50.3	37.7	43.1	100	0.4	0.8
Bilexical tree			(91.3)			-
MaltParser	57.8	53	55.3	-	—	-
LSTM Parser	66.1	61.1	63.5	-	-	-
Tree			(100)			-
UPARSE	52.7	52.8	52.8	-	—	-

Results on the 20K Leagues out-of-domain set.

#### Conclusion

- UCCA's semantic distinctions require a graph structure including challenging structural properties.
- TUPA is a transition-based parser suitable for UCCA, achieving high accuracy with BiLSTM model.
- Outperforms conversion-based parsing with a variety of strong bilexical DAG and tree baselines.

Corpora: http://www.cs.huji.ac.il/~oabend/ucca.html Code: https://github.com/danielhers/tupa Demo: https://rebrand.ly/tupa

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

#### Future Work

- Beam search, training with exploration.
- More languages (German corpus construction is underway).

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Parsing other schemes, such as AMR.
- Different target representations for conversion.
- Application to text simplification and other tasks.

# Thank you

#### References I

Abend, O. and Rappoport, A. (2013).

Universal Conceptual Cognitive Annotation (UCCA). In *Proc. of ACL*, pages 228–238.

Abend, O. and Rappoport, A. (2017). The state of the art in semantic representation. In *Proc. of ACL*. to appear.

Abend, O., Yerushalmi, S., and Rappoport, A. (2017). UCCAApp: Web-application for syntactic and semantic phrase-based annotation. In Proc. of ACL: System Demonstration Papers. to appear.

Almeida, M. S. C. and Martins, A. F. T. (2015). Lisbon: Evaluating TurboSemanticParser on multiple languages and out-of-domain data. In Proc. of SemEval, pages 970-973.

 Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., Knight, K., Palmer, M., and Schneider, N. (2013).
Abstract Meaning Representation for sembanking. In Proc. of the Linguistic Annotation Workshop.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Birch, A., Abend, O., Bojar, O., and Haddow, B. (2016). HUME: Human UCCA-based evaluation of machine translation. In Proc. of EMNLP, pages 1264–1274.

Dyer, C., Ballesteros, M., Ling, W., Matthews, A., and Smith, N. A. (2015). Transition-based dependeny parsing with stack long short-term memory. In Proc. of ACL, pages 334–343.

#### References II

Goldberg, Y. (2016).

A primer on neural network models for natural language processing.

#### Maier, W. (2015).

Discontinuous incremental shift-reduce parsing. In *Proc. of ACL*, pages 1202–1212.

Nivre, J. (2005).

Dependency grammar and dependency parsing.

#### Nivre, J. (2009).

Non-projective dependency parsing in expected linear time. In *Proc. of ACL*, pages 351–359.

- Nivre, J., Hall, J., Nilsson, J., Chanev, A., Eryigit, G., Kübler, S., Marinov, S., and Marsi, E. (2007). MaltParser: A language-independent system for data-driven dependency parsing. *Natural Language Engineering*, 13(02):95–135.
- Oepen, S., Kuhlmann, M., Miyao, Y., Zeman, D., Cinková, S., Flickinger, D., Hajič, J., and Urešová, Z. (2015). SemEval 2015 task 18: Broad-coverage semantic dependency parsing. In Proc. of SemEval, pages 915–926.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- Ribeyre, C., Villemonte de la Clergerie, E., and Seddah, D. (2014). Alpage: Transition-based semantic graph parsing with syntactic features. In Proc. of SemEval, pages 97–103.
- Sulem, E., Abend, O., and Rappoport, A. (2015). Conceptual annotations preserve structure across translations: A French-English case study. In *Proc. of S2MT*, pages 11–22.
- Tesniére, L. (1959). *Elements de syntaxe structuralle.* Klincksieck, Paris, 2 edition.
## UCCA Corpora

		Wiki		20K
	Train	Dev	Test	Leagues
# passages	300	34	33	154
# sentences	4268	454	503	506
# nodes	298,993	33,704	35,718	29,315
% terminal	42.96	43.54	42.87	42.09
% non-term.	58.33	57.60	58.35	60.01
% discont	0 5/	0 53	0 44	0.81
/0 discont.	0.34	0.55	0.44	0.01
% reentrant	2.38	1.88	2.15	2.03
% reentrant # edges	<b>2.38</b> 287,914	<b>1.88</b> 32,460	<b>2.15</b> 34,336	<b>2.03</b> 27,749
% reentrant # edges % primary	<b>2.38</b> 287,914 98.25	<b>1.88</b> 32,460 98.75	<b>2.15</b> 34,336 98.74	<b>2.03</b> 27,749 97.73
% reentrant # edges % primary % remote	<b>2.38</b> 287,914 98.25 1.75	1.88 32,460 98.75 1.25	<b>2.15</b> 34,336 98.74 1.26	<b>2.03</b> 27,749 97.73 2.27
% reentrant # edges % primary % remote Average per non-	<b>2.38</b> 287,914 98.25 1.75 terminal noo	1.88 32,460 98.75 1.25 de	<b>2.15</b> 34,336 98.74 1.26	<b>2.03</b> 27,749 97.73 2.27

Corpus statistics.

## Evaluation

Mutual edges between predicted graph  $G_p = (V_p, E_p, \ell_p)$  and gold graph  $G_g = (V_g, E_g, \ell_g)$ , both over terminals  $W = \{w_1, \dots, w_n\}$ :

$$M(G_{p}, G_{g}) = \left\{ (e_{1}, e_{2}) \in E_{p} \times E_{g} \mid y(e_{1}) = y(e_{2}) \wedge \ell_{p}(e_{1}) = \ell_{g}(e_{2}) \right\}$$

The yield  $y(e) \subseteq W$  of an edge e = (u, v) in either graph is the set of terminals in W that are descendants of v.  $\ell$  is the edge label.

Labeled precision, recall and F-score are then defined as:

$$LP = \frac{|M(G_p, G_g)|}{|E_p|}, \quad LR = \frac{|M(G_p, G_g)|}{|E_g|},$$
$$LF = \frac{2 \cdot LP \cdot LR}{LP + LR}.$$

Two variants: one for primary edges, and another for remote edges.