



Zero-resource Dependency Parsing: Boosting Delexicalized Cross-lingual Transfer with Linguistic Knowledge

Lauriane Aufrant, Guillaume Wisniewski, François Yvon
December 13, 2016

COLING'16

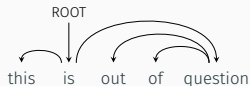


Zero-resource scenario:

{ No annotated data
No parallel data \implies delexicalized cross-lingual transfer
No raw data

Zero-resource scenario:

{ No annotated data
No parallel data \implies delexicalized cross-lingual transfer
No raw data



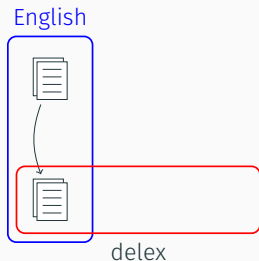
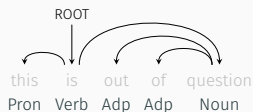
English



Zero-resource scenario:

{ No annotated data
No parallel data
No raw data

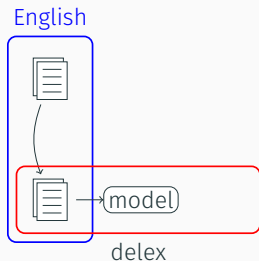
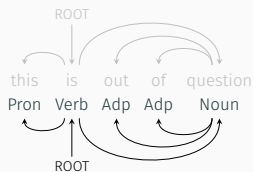
⇒ delexicalized cross-lingual transfer



Zero-resource scenario:

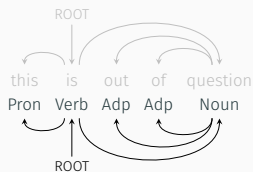
{ No annotated data
No parallel data
No raw data

⇒ delexicalized cross-lingual transfer

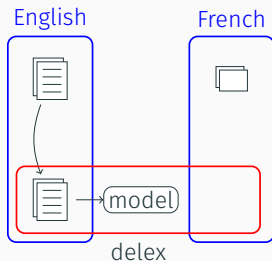


Zero-resource scenario:

{ No annotated data
No parallel data
No raw data } \Rightarrow delexicalized cross-lingual transfer

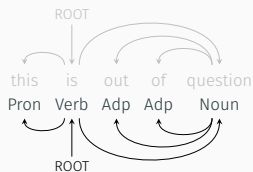


l' autre rive est hors de portée .

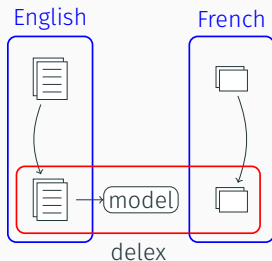


Zero-resource scenario:

{ No annotated data
No parallel data
No raw data } \Rightarrow delexicalized cross-lingual transfer

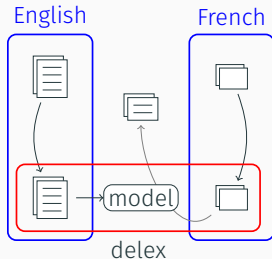
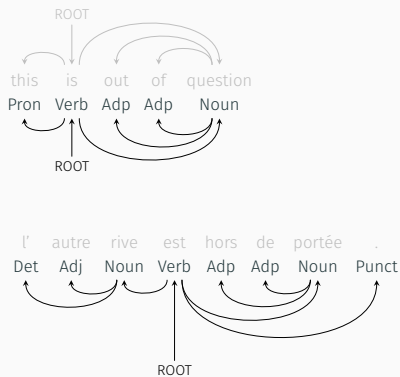


l' autre rive est hors de portée .
Det Adj Noun Verb Adp Adp Noun Punct



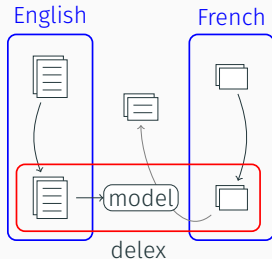
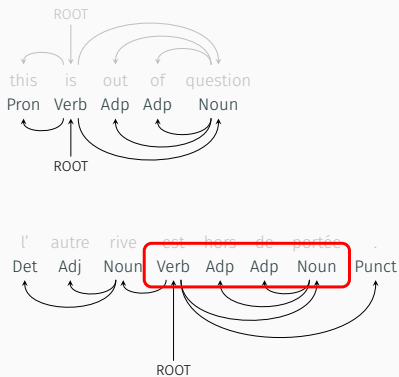
Zero-resource scenario:

{ No annotated data
No parallel data \Rightarrow delexicalized cross-lingual transfer
No raw data



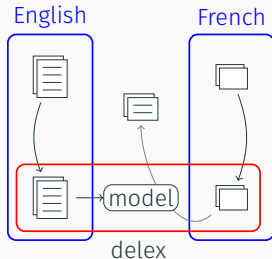
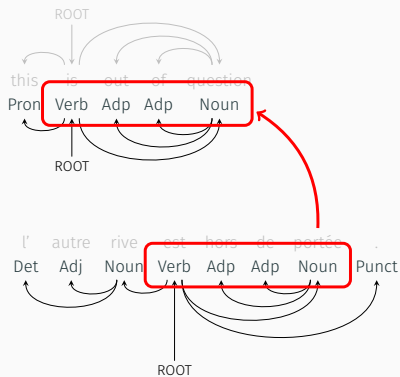
Zero-resource scenario:

{ No annotated data
No parallel data \Rightarrow delexicalized cross-lingual transfer
No raw data



Zero-resource scenario:

{ No annotated data
No parallel data \implies delexicalized cross-lingual transfer
No raw data



An adjective close to a noun depends on this noun.

An adjective close to a noun depends on this noun.

An adjective close to a noun depends on this noun.

True in...

✓ English

✓ French

✓ Hebrew

✓ Bulgarian

An adjective close to a noun depends on this noun.

True in...

✓ English

✓ French

✓ Hebrew

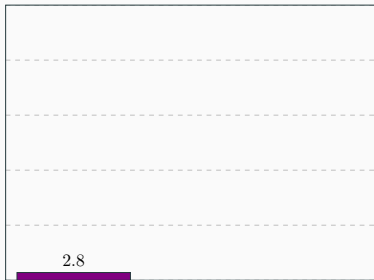
✓ Bulgarian

Hebrew → Hebrew



NOUN
↓
ADJ

Hebrew → Bulgarian



NOUN
↓
ADJ

An adjective close to a noun depends on this noun.

True in...

✓ English

✓ French

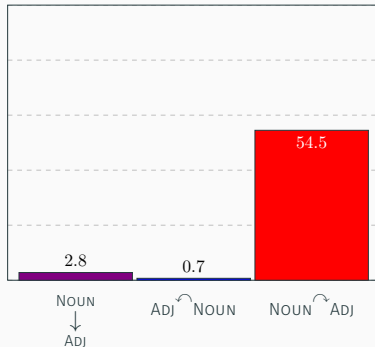
✓ Hebrew

✓ Bulgarian

Hebrew → Hebrew



Hebrew → Bulgarian



A step towards scaling to the 7,000 languages in the world

↔ zero-resource dependency transfer

Many transfer errors are easy to avoid

↔ regular divergences between source and target

↔ word order issues

Our approach: leveraging previous works in linguistics (WALS)

↔ +3 UAS on average

↔ very efficient on some error types: up to +90 UAS

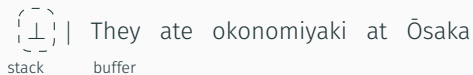
A fine-grained analysis across various language pairs

↔ 6,000+ experiments

1. Transition-based dependency parsing
2. Impact of word order on transfer
3. Leveraging WALS
4. Leveraging raw data
5. Wrap-up

Transition-based dependency parsing

Parsing with the ArcEager system



They ate okonomiyaki at Ōsaka

CLASSIFIER

Parsing with the ArcEager system

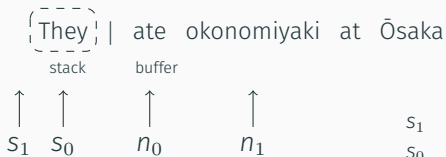
(They) | ate okonomiyaki at Ōsaka
stack buffer

They ate okonomiyaki at Ōsaka

CLASSIFIER

SHIFT

Parsing with the ArcEager system



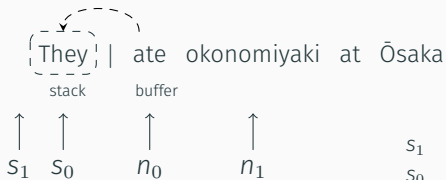
$S_1 = \emptyset$
 $S_0 = \text{THEY}$
 $n_0 = \text{ATE}$
 $n_1 = \text{OKONOMIYAKI}$

↓

CLASSIFIER

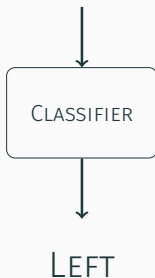
They ate okonomiyaki at Ōsaka

Parsing with the ArcEager system

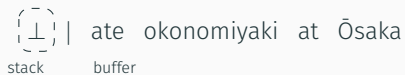


$s_1 = \emptyset$
 $s_0 = \text{THEY}$
 $n_0 = \text{ATE}$
 $n_1 = \text{OKONOMIYAKI}$

They ate okonomiyaki at Ōsaka



Parsing with the ArcEager system

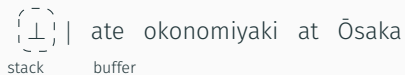
 | ate okonomiyaki at Ōsaka
stack buffer


They ate okonomiyaki at Ōsaka

CLASSIFIER

LEFT

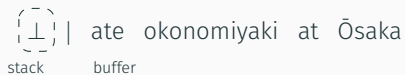
Parsing with the ArcEager system

 | ate okonomiyaki at Ōsaka
stack buffer


They ate okonomiyaki at Ōsaka

CLASSIFIER

Parsing with the ArcEager system

 | ate okonomiyaki at Ōsaka
stack buffer


They ate okonomiyaki at Ōsaka

CLASSIFIER

SHIFT

Parsing with the ArcEager system

{ate} | okonomiyaki at Ōsaka
stack buffer

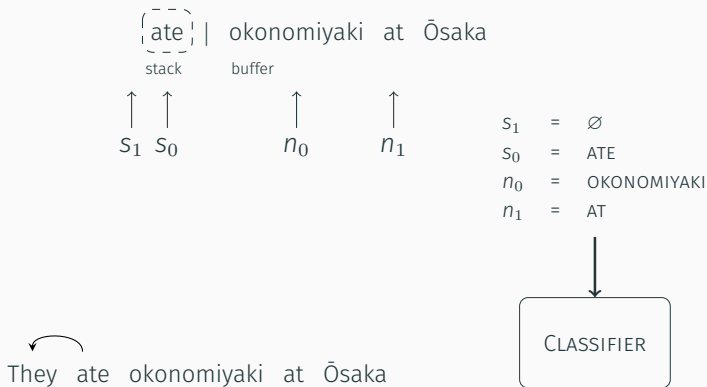
They ate okonomiyaki at Ōsaka



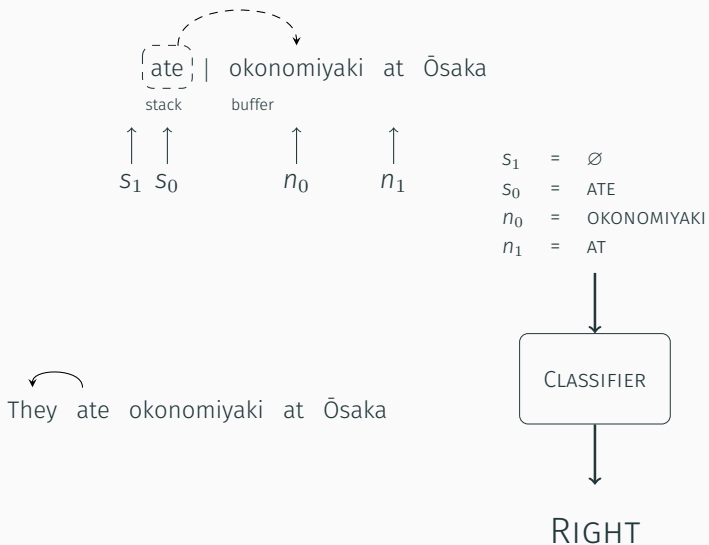
CLASSIFIER

SHIFT

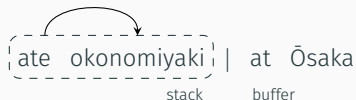
Parsing with the ArcEager system



Parsing with the ArcEager system



Parsing with the ArcEager system

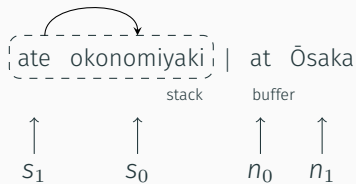


They ate konomiyaki at Ōsaka

CLASSIFIER

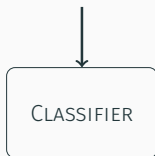
RIGHT

Parsing with the ArcEager system

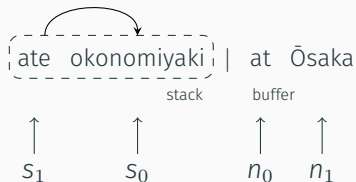


S_1 = ATE
 S_0 = OKONOMIYAKI
 n_0 = AT
 n_1 = ŌSAKA

They ate okonomiyaki at Ōsaka

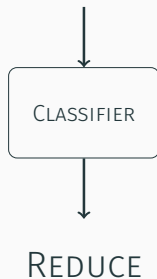


Parsing with the ArcEager system



S_1 = ATE
 S_0 = OKONOMIYAKI
 n_0 = AT
 n_1 = ŌSAKA

They ate okonomiyaki at Ōsaka



Parsing with the ArcEager system

{ate} | at Ōsaka
stack buffer

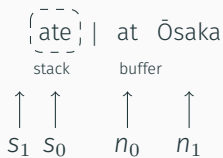
They ate okonomiyaki at Ōsaka



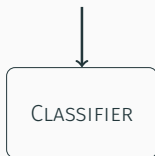
CLASSIFIER

REDUCE

Parsing with the ArcEager system

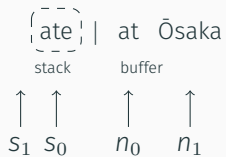


$s_1 = \emptyset$
 $s_0 = \text{ATE}$
 $n_0 = \text{AT}$
 $n_1 = \text{ŌSAKA}$



They ate okonomiyaki at Ōsaka

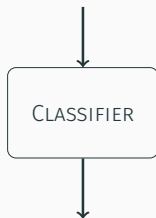
Parsing with the ArcEager system



They ate okonomiyaki at Ōsaka

Diagram showing the sentence "They ate okonomiyaki at Ōsaka" with arcs indicating dependencies between "ate" and "okonomiyaki", and "at" and "Ōsaka".

$s_1 = \emptyset$
 $s_0 = \text{ATE}$
 $n_0 = \text{AT}$
 $n_1 = \text{ŌSAKA}$



SHIFT

Parsing with the ArcEager system

{ate at} | Ōsaka
stack buffer

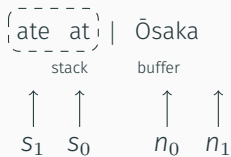
They ate okonomiyaki at Ōsaka



CLASSIFIER

SHIFT

Parsing with the ArcEager system

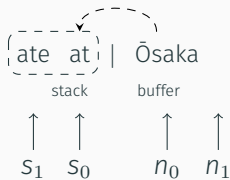


$s_1 = \text{ATE}$
 $s_0 = \text{AT}$
 $n_0 = \text{ŌSAKA}$
 $n_1 = \emptyset$

↓
CLASSIFIER

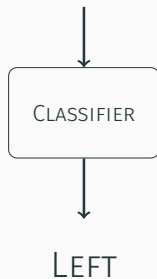
They ate okonomiyaki at Ōsaka

Parsing with the ArcEager system



They ate okonomiyaki at Ōsaka

s_1 = ATE
 s_0 = AT
 n_0 = ŌSAKA
 n_1 = \emptyset



Parsing with the ArcEager system

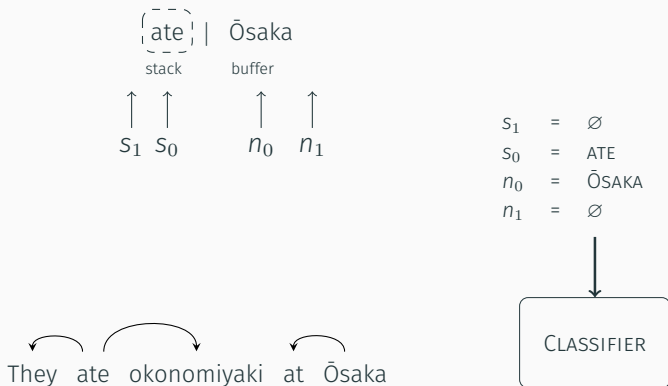
{ate} | Ōsaka
stack buffer

They ate okonomiyaki at Ōsaka

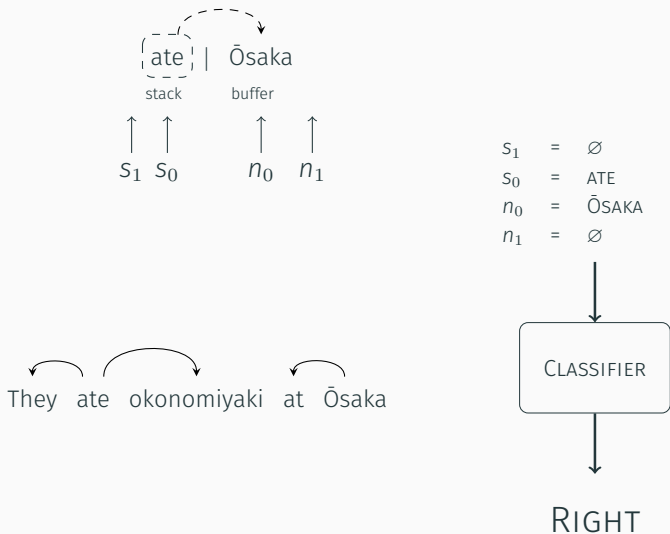
CLASSIFIER

LEFT

Parsing with the ArcEager system



Parsing with the ArcEager system



Parsing with the ArcEager system



CLASSIFIER

RIGHT

Parsing with the ArcEager system



CLASSIFIER

Parsing with the ArcEager system



CLASSIFIER

REDUCE

Parsing with the ArcEager system

(ate) | ⊥
stack buffer

They ate okonomiyaki at Ōsaka



CLASSIFIER

REDUCE

Parsing with the ArcEager system

(ate) | ⊥
stack buffer

They ate okonomiyaki at Ōsaka



CLASSIFIER

Parsing with the ArcEager system

(ate) | ⊥
stack buffer



CLASSIFIER

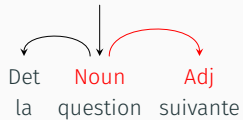
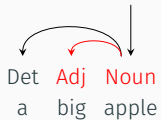
Impact of word order on transfer

Towards language-independent data representations

- **Words:** delexicalized, common tagset
 - ✓ UPOS [Petrov et al., 2012]
- **Dependencies:** common guidelines
 - ✓ Universal Dependencies [Nivre et al., 2016]
- ... other regular divergences?

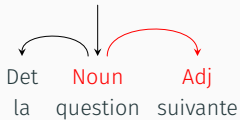
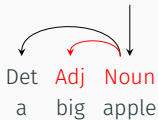
Impact of word order

At data level:



Impact of word order

At data level:



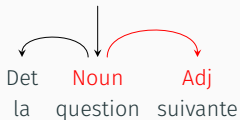
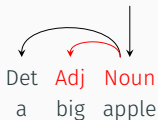
At model level:

$(s_0 = \text{ADJ} \wedge n_0 = \text{NOUN}) \Rightarrow \text{LEFT}$

$(s_0 = \text{NOUN} \wedge n_0 = \text{ADJ}) \Rightarrow \text{RIGHT}$

Impact of word order

At data level:



At model level:

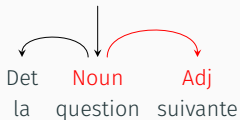
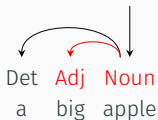
different features + different decisions

$(s_0 = \text{ADJ} \wedge n_0 = \text{NOUN}) \Rightarrow \text{LEFT}$

$(s_0 = \text{NOUN} \wedge n_0 = \text{ADJ}) \Rightarrow \text{RIGHT}$

Impact of word order

At data level:



At model level:

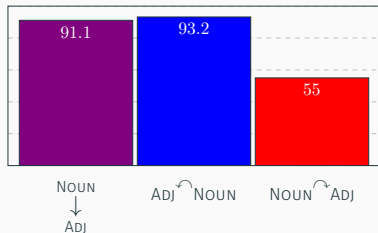
different features + different decisions

$(s_0 = \text{ADJ} \wedge n_0 = \text{NOUN}) \Rightarrow \text{LEFT}$

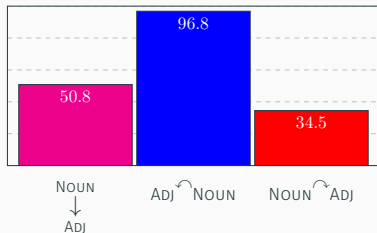
$(s_0 = \text{NOUN} \wedge n_0 = \text{ADJ}) \Rightarrow \text{RIGHT}$

On accuracy:

English \rightarrow English

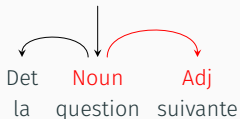
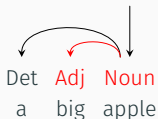


English \rightarrow French



Impact of word order

At data level:



At model level:

different features + different decisions

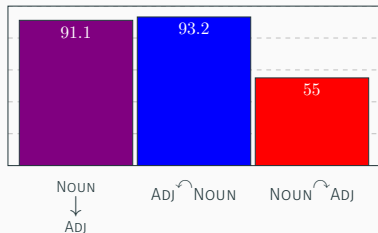
$(s_0 = \text{ADJ} \wedge n_0 = \text{NOUN}) \Rightarrow \text{LEFT}$

$(s_0 = \text{NOUN} \wedge n_0 = \text{ADJ}) \Rightarrow \text{RIGHT}$

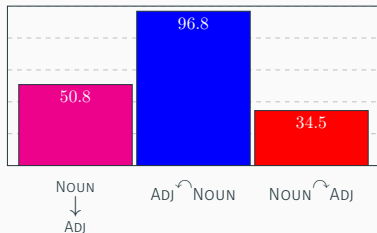
On accuracy:

no knowledge sharing + transfer errors

English \rightarrow English



English \rightarrow French

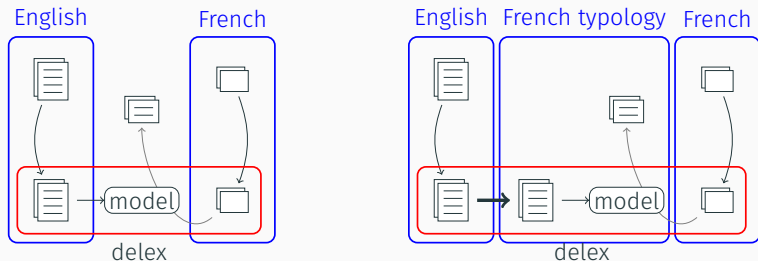


Leveraging WALs

Reshaping training instances

Our proposal:

preprocess source data, to match target regularities



↔ same dependency annotations

↔ same training process

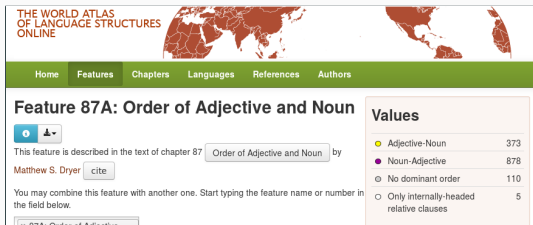
Reshaping training instances... a few examples

English – training data		French – desired output
Baseline	Our proposal	
<p>Det Adj Noun a big apple</p>	⇒ <p>Det Noun Adj a apple big</p>	<p>Det Noun Adj la question suivante</p>
<p>Det Adj Noun the whole world</p>	⇒ <p>Det Noun Adj the world whole</p>	<p>Det Noun Adj la flotte romaine</p>
<p>Det Noun Noun an investment firm</p>	⇒ <p>Det Noun Noun an firm investment</p>	<p>Det Noun Adp Noun la plate-forme de distribution</p>

The World Atlas of Language Structures

WALS: a database of typological features for 2,679 languages
[<http://wals.info>]

↔ Over 1,000 languages with word order features



THE WORLD ATLAS OF LANGUAGE STRUCTURES ONLINE

Home Features Chapters Languages References Authors

Feature 87A: Order of Adjective and Noun

This feature is described in the text of chapter 87 [Order of Adjective and Noun](#) by [Matthew S. Dryer](#) [cite](#)

You may combine this feature with another one. Start typing the feature name or number in the field below.

Values	
<input checked="" type="radio"/> Adjective-Noun	373
<input type="radio"/> Noun-Adjective	878
<input type="radio"/> No dominant order	110
<input type="radio"/> Only internally-headed relative clauses	5

English	<input checked="" type="radio"/> Adjective-Noun		<input type="radio"/>	<input type="checkbox"/>
French	<input type="radio"/> Noun-Adjective	Harris 1988: 227	<input type="radio"/>	<input type="checkbox"/>

Heuristic rule extraction for **switching** and **deleting** tokens

87A $\left\{ \begin{array}{l} [\text{English}] \text{ Adjective-Noun} \\ [\text{French}] \text{ Noun-Adjective} \end{array} \right.$

\implies [English \rightarrow French] switch (most?) ADJ-NOUN into NOUN-ADJ

Heuristic rule extraction for **switching** and **deleting** tokens

$$87A \quad \left\{ \begin{array}{l} [\text{English}] \text{ Adjective-Noun} \\ [\text{French}] \text{ Noun-Adjective} \end{array} \right.$$

⇒ [English→French] switch (most?) ADJ-NOUN into NOUN-ADJ

- ✓ easily extensible rule templates
 - ↔ switch tokens, but also delete, insert, replace...
- ✓ accepts free order
- ✓ mostly conservative for related languages

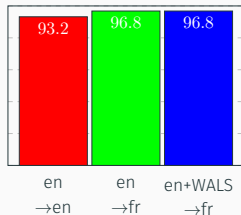
- ✓ readily available for 1,000 languages
- ✓ most work already done by linguists

Experimental results on UD 1.3

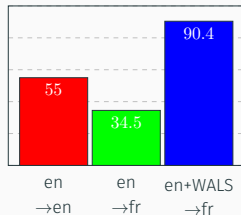
English → French

Hebrew → Bulgarian

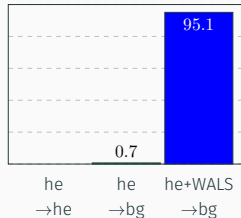
ADJ → NOUN



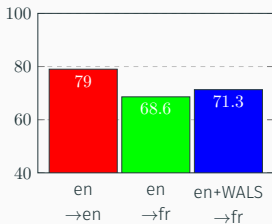
NOUN → ADJ



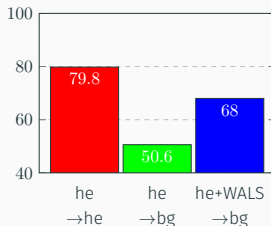
ADJ → NOUN



Overall score: +2.7 UAS



Overall score: +17.4 UAS



Experimental results on UD 1.3... over 40 languages

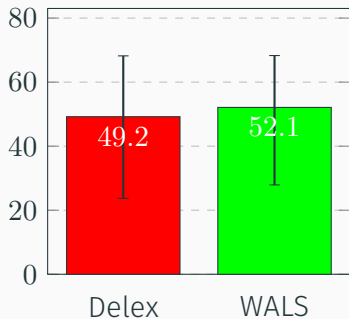
↪ 6,000+ experiments

65.9	+1.2
78.4	+1.8
79.9	+4.4
62.0	+5.7
48.4	+5.2
53.0	+3.7
66.0	+1.8
70.0	+2.3
77.7	+2.6
73.9	+0.9

Target	Differential			Mean scores									Mean score					
	PSLM selection			PSLM evolution			WALS score rules			WALS			Index	WALS				
	min	med	max	min	med	max	min	med	max	min	med	max						
ar	5.1	43.2	86.9	36.1	4.8	43.0	87.2	32.7	18.9	45.1	87.2	+5.6	12.2	47.6	90.9	+5.7	57.3	89.0
bg	20.4	97.5	78.9	89.6	26.5	97.5	78.9	-0.1	33.5	60.8	75.5	+0.5	27.2	67.6	80.9	+1.0	76.0	79.1
ca	20.4	62.3	76.5	57.8	27.9	62.6	76.7	-0.1	33.5	60.8	75.5	+0.2	30.4	66.2	76.6	+1.9	79.0	79.1
cs	29.7	96.1	70.0	56.2	26.6	74.0	-0.0	37.1	58.4	68.8	+4.4	94.8	79.3	73.8	+1.4	64.0	73.8	
de	22.6	58.5	76.7	53.9	22.6	58.7	75.4	+0.9	34.2	60.2	70.0	+3.8	24.3	60.2	70.7	+1.6	57.7	76.7
el	20.0	64.2	76.3	36.2	26.5	63.8	75.2	-0.2	40.1	61.0	70.4	-0.0	28.6	64.6	74.5	+1.4	67.1	75.2
en	36.2	61.2	76.5	57.7	55.9	61.0	70.0	-0.2	45.4	61.1	69.1	+1.3	43.0	61.8	70.9	+1.7	68.7	68.7
es	21.0	51.0	67.8	46.8	28.9	60.2	-0.0	33.5	51.6	64.9	-0.0	33.5	51.6	64.9	+1.7	67.0	66.2	
eu	33.1	56.5	65.8	52.8	32.5	56.2	65.5	-0.2	38.4	57.0	63.8	+6.2	32.2	58.4	65.9	+1.7	65.7	66.0
fi	30.0	63.9	76.5	58.9	28.5	64.4	79.5	-0.0	37.6	62.8	76.1	+0.3	31.5	66.6	80.6	+1.7	79.2	79.1
fr	20.4	53.1	69.4	31.5	26.0	52.9	69.6	-0.2	36.1	57.0	67.8	+0.7	37.1	57.9	68.7	+0.4	68.4	69.1
gl	20.7	45.2	57.8	46.2	26.9	49.4	57.7	-0.1	24.5	34.6	46.1	+4.6	24.6	32.0	45.3	+5.7	35.7	49.0
he	17.9	45.3	56.1	48.3	37.6	45.9	56.0	-0.1	38.7	46.3	56.8	+1.1	35.5	48.4	58.8	+5.2	61.4	63.3
hr	27.4	48.1	62.1	46.6	27.4	48.2	61.8	-0.0	32.4	39.2	50.8	+2.3	32.4	39.2	50.2	+5.7	66.0	65.7
hu	30.9	64.0	76.1	59.0	28.9	63.9	79.3	-0.2	35.0	64.4	76.8	+0.5	34.2	64.0	80.2	+1.9	79.1	79.5
id	18.4	56.0	65.8	56.1	26.5	62.4	66.4	-0.1	26.6	36.2	44.6	+1.6	28.8	39.0	48.7	+4.9	47.4	62.2
it	33.0	49.1	67.5	48.6	32.8	49.5	67.1	-0.1	32.1	43.0	60.2	+1.5	35.0	43.7	51.0	+2.6	48.7	49.6
ja	20.4	78.0	75.7	54.5	26.6	77.5	77.4	-0.0	34.0	60.0	66.2	+3.2	28.2	59.9	73.9	+0.9	72.7	73.9
ko	32.5	53.9	57.5	50.7	29.8	53.9	57.8	-0.1	39.0	53.9	58.1	+6.8	32.4	53.9	57.8	+0.0	62.0	60.1
ku	20.1	53.8	68.0	49.9	19.8	54.2	67.7	-0.0	36.2	54.1	63.6	+0.2	31.6	55.4	65.8	+1.6	71.2	68.7
la	11.0	27.1	66.5	32.5	11.1	28.9	65.8	-0.2	22.0	37.5	61.6	+6.7	19.8	33.8	60.9	+5.4	37.1	44.2
lt	11.4	49.2	79.1	48.6	18.2	49.0	79.1	-0.0	37.7	56.2	66.1	+0.5	23.1	43.8	69.6	+5.7	78.0	73.9
lv	27.9	52.7	67.8	50.8	27.1	53.1	68.2	-0.0	49.4	56.5	65.4	+4.4	40.1	55.4	68.3	+1.4	63.0	64.4
nl	17.6	18.8	52.6	36.7	5.2	18.8	52.6	-0.1	15.7	19.9	70.6	+0.2	14.0	19.8	82.3	+2.2	83.2	82.9
no	31.0	67.1	82.6	61.7	39.4	69.5	82.1	-0.1	38.1	67.0	80.7	+1.2	34.0	70.8	82.3	+1.4	81.1	80.0
pl	14.4	49.9	66.1	47.1	14.5	50.0	65.9	-0.0	28.5	41.4	61.8	+1.9	21.3	63.5	62.1	+6.1	66.1	60.0
pt	18.7	33.0	56.3	32.4	19.0	34.0	54.3	+0.2	17.9	32.2	52.8	+1.4	20.0	38.5	59.6	+4.9	53.9	54.6
ro	22.8	49.1	66.1	47.1	20.9	48.2	67.9	-0.1	32.4	48.3	56.8	+2.4	28.4	45.4	60.4	+0.5	62.1	60.7
ru	22.8	49.1	66.1	47.1	20.9	48.2	67.9	-0.1	32.4	48.3	56.8	+2.4	28.4	45.4	60.4	+0.5	62.1	60.7
sk	25.9	64.0	76.6	57.1	25.6	63.9	76.7	-0.1	37.2	58.5	69.2	+2.2	28.5	63.8	76.2	+1.5	78.6	76.5
sl	25.7	62.1	77.9	59.2	25.6	63.7	77.9	-0.1	34.0	63.6	76.1	+0.4	29.5	57.4	77.4	+1.0	78.0	77.9
sv	30.8	62.8	76.5	56.7	30.2	63.3	75.9	+0.0	35.7	60.0	73.5	-0.9	32.0	63.5	76.4	+1.4	75.5	75.7
th	18.8	57.5	66.2	51.6	16.6	55.7	66.2	-0.1	31.5	45.7	67.5	+1.1	24.5	53.5	64.8	+1.0	72.0	73.4
tr	20.8	53.9	69.0	51.3	26.1	54.0	68.9	-0.0	34.8	55.1	67.8	+2.3	30.9	59.7	68.7	+4.2	71.0	70.4
uk	30.8	65.2	80.4	58.4	28.4	65.1	80.1	-0.0	42.8	64.8	77.4	+3.7	38.9	64.9	80.4	+1.4	80.4	80.4
vi	20.4	62.5	76.5	58.9	29.4	62.3	75.9	-0.2	39.0	60.1	70.6	+0.5	38.4	63.2	74.9	+2.3	71.1	73.5
wa	9.1	36.3	56.5	36.5	9.1	36.6	56.5	-0.0	19.9	43.7	64.5	+6.2	19.1	43.8	65.6	+1.3	65.3	65.9
xh	14.1	35.3	67.0	38.8	14.7	35.4	67.0	-0.1	19.5	39.5	64.5	+2.0	21.5	40.8	67.8	+3.6	58.5	58.5
zh	17.6	32.8	63.1	37.8	17.5	32.8	63.0	-0.1	18.8	35.5	61.9	+1.3	20.1	38.6	64.8	+3.6	63.2	62.2
Avg	23.7	52.0	68.2	42.2	23.5	52.0	68.1	-0.1	31.8	53.5	65.8	+2.3	27.9	55.2	66.3	+2.0	66.9	67.4

Experimental results on UD 1.3... over 40 languages

↔ 6,000+ experiments



Average improvement:

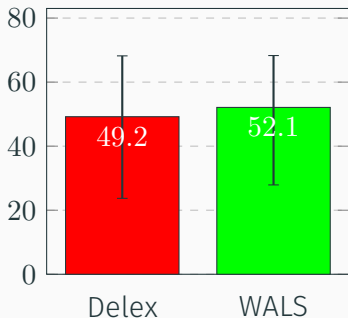
+2.9 UAS

+0.1 on best sources

+4.2 on worst sources

Experimental results on UD 1.3... over 40 languages

↔ 6,000+ experiments



Average improvement:

+2.9 UAS

+0.1 on best sources

+4.2 on worst sources

- ✓ Very efficient on nearly deterministic reorderings
- ✓ Specific error correction: many pairs (21%) have a very large (50+) improvement on at least one frequent tag pair
- ✓ Rarely detrimental

Leveraging raw data

Extract target regularities from raw data

Contrastive experiment:

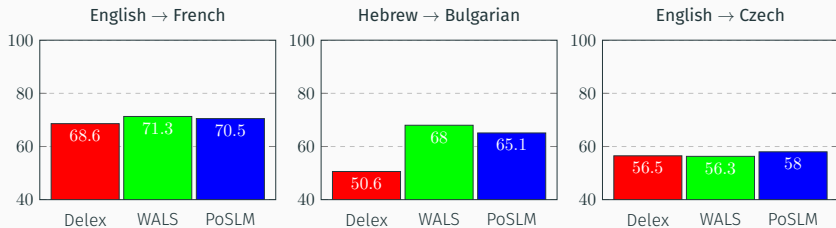
knowledge-driven method \iff data-driven method

\hookrightarrow use a target language model to reorder the PoS sequence

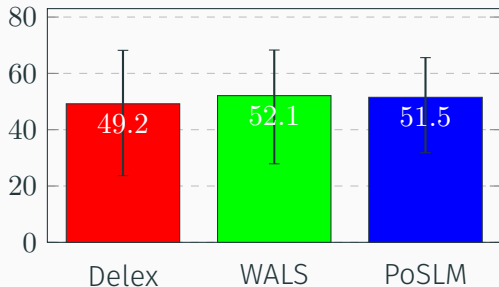
Procedure:

- train a PoSLM on delexicalized raw target data
- generate **local reordering lattices** for source data
- score them with the PoSLM
- keep the best projective reordering

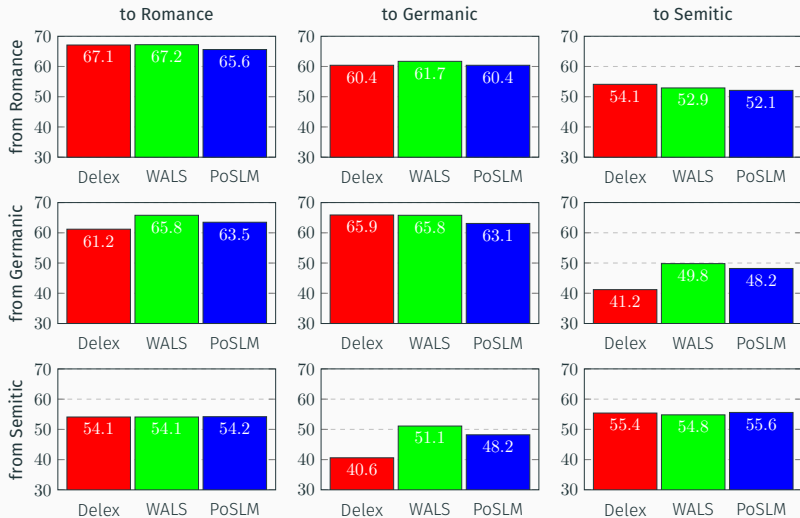
Comparative results



Average improvement: +2.3 UAS



Language family analysis



PoSLM approach: pros and cons

- ✗ can only reorder tokens
- ✗ does not acknowledge free order
- ✓ captures more diverse reordering patterns

- ✗ not *really* zero-resource
- ✓ *in theory* applicable to any language
- ✗ much more involved method

- ✗ too hard constraints for closely related sources
- ✗ sometimes out of control
- ✓ combine with known typology?

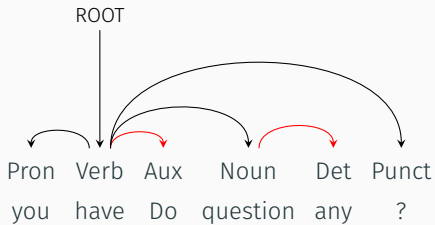
Wrap-up

- **Zero-resource** scenario: target 1,000 languages, delexicalized cross-lingual transfer
- Despite common guidelines, **regular divergences** still affect representations, at model level
⇒ poor knowledge sharing, transfer errors
- We reshape training examples, using **target word order**
⇒ large improvements on specific errors
- Two setups: the naive **knowledge-rich approach** outperforms the involved **data-driven method**

- **Zero-resource** scenario: target 1,000 languages, delexicalized cross-lingual transfer

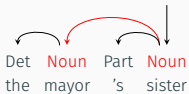
- Despite common guidelines, **regular divergences** still affect representations, at model level
⇒ poor knowledge sharing, transfer errors

- We reshape training examples, using **target word order**
⇒ large improvements on specific errors
- Two setups: the naive **knowledge-rich approach** outperforms the involved **data-driven method**



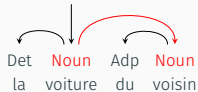
Appendix: more examples

English



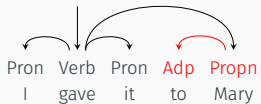
$(s_0 = \text{NOUN} \wedge n_0 = \text{NOUN}) \Rightarrow \text{LEFT}$

→ French



$(s_0 = \text{NOUN} \wedge n_0 = \text{NOUN}) \Rightarrow \text{RIGHT}$

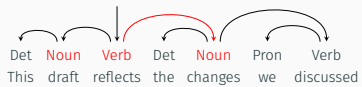
English



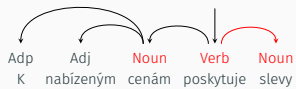
→ Japanese



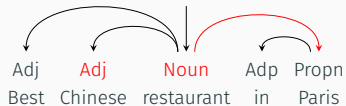
English



→ Czech



Limits of feature engineering



Full abstraction from word order is not desirable