



# Dependency parsing

Benoît Sagot

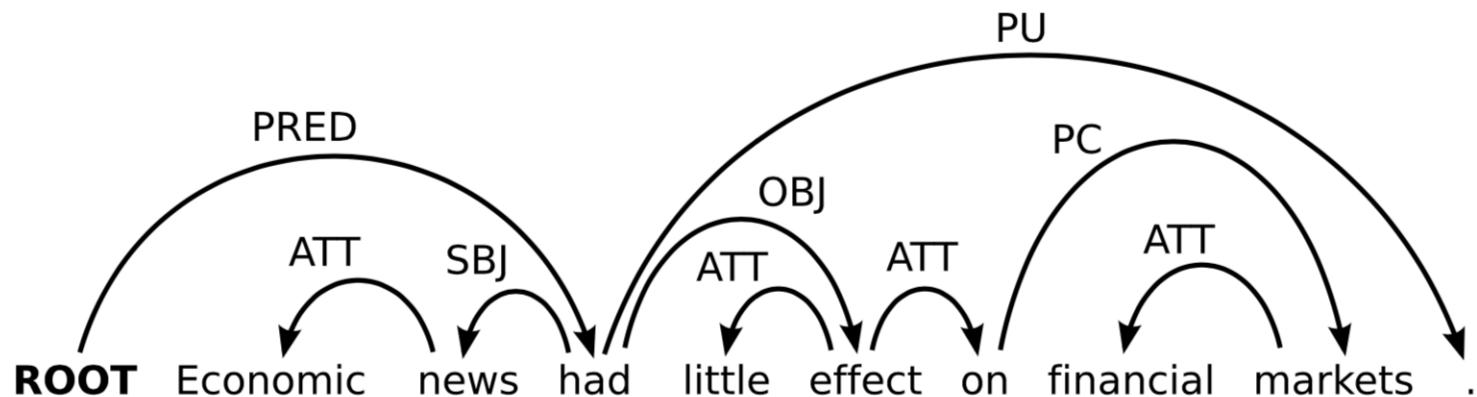
Inria (ALMAnaCH), Paris (France)

MVA — Speech and Language Processing — Class #6 — 26<sup>th</sup> February, 2018

Credit and disclaimer: some of the following slides are taken from, illustrated or inspired by presentations and article figures by Nivre, Ballesteros, Rasooli and Tetreault

# Introduction

- **Syntactic parsing of natural language**
  - Building the structure of natural language sentences
- **Dependency-based syntactic representations**
  - Long tradition in descriptive and theoretical linguistics
  - Have become popular in computational linguistics



# Strategies for dependency parsing

- Graph-based parsing
- Transition-based parsing
- Other strategies

# Graph-based parsing

- MSTParser (McDonald et al. 2005)
  - <http://www.seas.upenn.edu/~strctlrn/MSTParser/MSTParser.html>
- Simplified version of the underlying idea:
  - Create all possible dependencies
  - Weigh them
  - Extract the optimal dependency tree
    - I.e. the tree that covers all words and minimises the overall weight of all retained dependencies

# Arc-standard Transition-Based Parsing

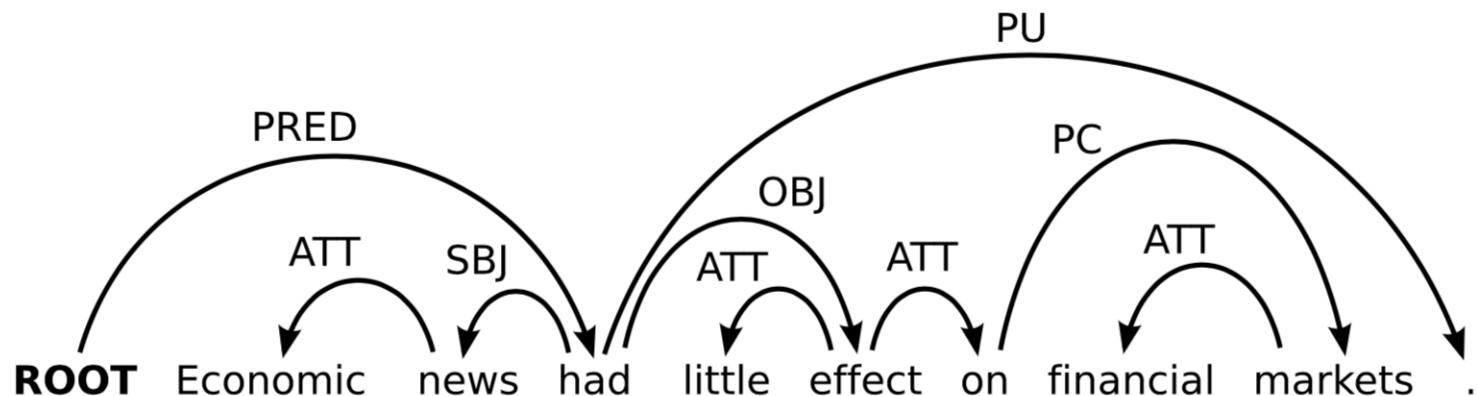


# Starting point

- **The basic idea:**
  - Define a transition system for dependency parsing
  - Learn a model for scoring possible transitions
  - Parse by searching for the optimal transition sequence
- **Advantages:**
  - Highly efficient parsing with low complexity
  - Rich history-based feature models for disambiguation
- Cf. Nivre (et al.)
  - <http://www.maltparser.org>

# Formalising dependency trees

- A **dependency tree** is a **labelled directed tree**  $T$  with
  - a set  $V$  of nodes, labelled with wordforms (including the special “wordform” **ROOT**)
  - a set  $A$  of arcs, labelled with dependency types
  - a linear precedence order  $<$  on  $V$
- **Notation:**
  - Arc  $(w_i, l, w_j)$  connects head  $w_i$  to dependent  $w_j$  with label  $l$
  - Node  $w_0$  (labeled **ROOT**) is the unique root of the tree



# Parser configurations

- A **parser configuration** is a triple  $c=(S, Q, A)$ , where
  - $S =$  a stack  $[ \dots, w_i ]_S$  of partially processed nodes,
  - $Q =$  a queue  $[w_j, \dots]_Q$  of remaining input nodes,
  - $A =$  a set of labelled arcs  $(w_i, l, w_j)$ .
- **Initialisation:**  
 $([w_0]_S, [w_1, \dots, w_n]_Q, \{\})$   
(recall that  $w_0 = \text{ROOT}$ )
- **Termination:**  $([w_0]_S, []_Q, A)$

# Transitions for the “arc-standard algorithm”

- **Left-Arc( $l$ )**

$$\begin{array}{c} ([\dots, w_i, w_j]_S, Q, A) \\ \hline ([\dots, w_j]_S, Q, A \cup \{(w_j, l, w_i)\}) \end{array} \quad [i \neq 0]$$

- **Right-Arc( $l$ )**

$$\begin{array}{c} ([\dots, w_i, w_j]_S, Q, A) \\ \hline ([\dots, w_i]_S, Q, A \cup \{(w_i, l, w_j)\}) \end{array}$$

- **Shift**

$$\begin{array}{c} ([\dots]_S, [w_i, \dots]_Q, A) \\ \hline ([\dots, w_i]_S, [\dots]_Q, A) \end{array}$$

# Example Transition Sequence

[ROOT]<sub>s</sub> [Economic, news, had, little, effect, on, financial, markets, .]<sub>Q</sub>

**ROOT** Economic news had little effect on financial markets .

# Example Transition Sequence

[ROOT, Economic]<sub>s</sub> [news, had, little, effect, on, financial, markets, .]<sub>Q</sub>

action: Shift

**ROOT** Economic news had little effect on financial markets .

# Example Transition Sequence

[ROOT, Economic, news]<sub>s</sub> [had, little, effect, on, financial, markets, .] <sub>Q</sub>

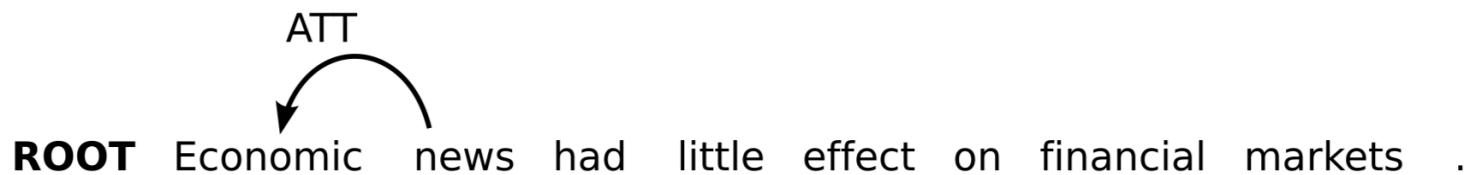
action: Shift

**ROOT** Economic news had little effect on financial markets .

# Example Transition Sequence

[ROOT, Economic, news]<sub>s</sub> [had, little, effect, on, financial, markets, .]\_<sub>Q</sub>

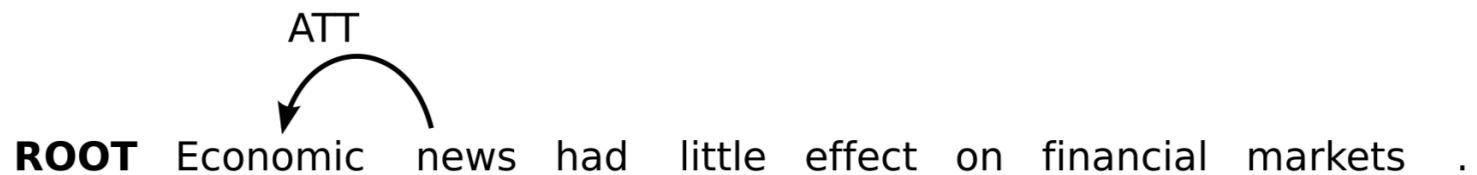
action: Left-Arc(ATT)



# Example Transition Sequence

[ROOT, news, had]<sub>S</sub> [little, effect, on, financial, markets, .]<sub>Q</sub>

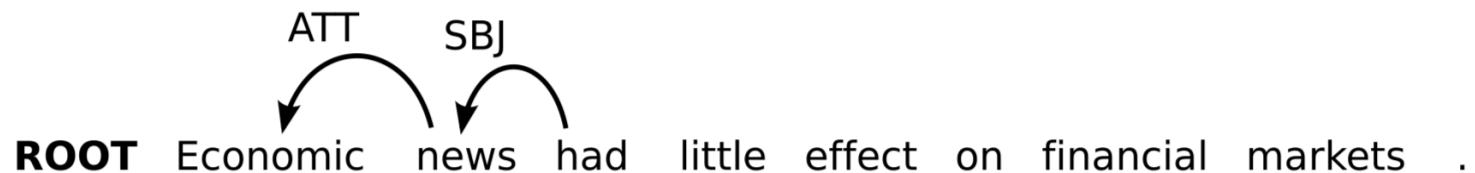
**action: Shift**



# Example Transition Sequence

[ROOT, news, had]<sub>S</sub> [little, effect, on, financial, markets, .]\_{Q}

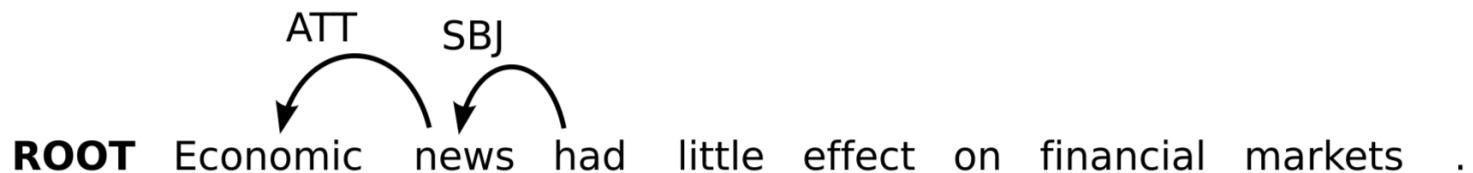
action: Left-Arc(SBJ)



# Example Transition Sequence

[ROOT, had, little]<sub>S</sub> [effect, on, financial, markets, .]<sub>Q</sub>

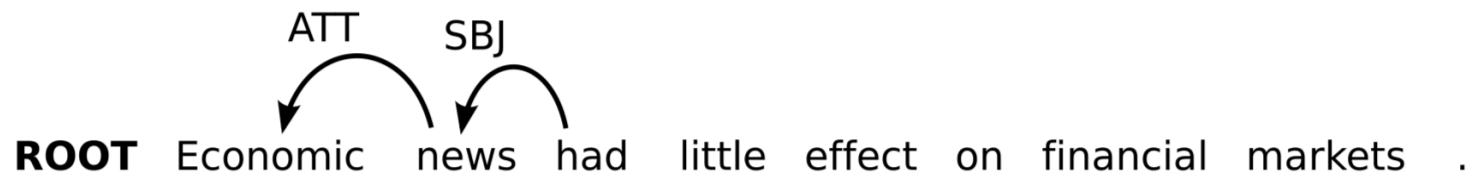
action: Shift



# Example Transition Sequence

[ROOT, had, little, effect]<sub>S</sub> [on, financial, markets, .]Q

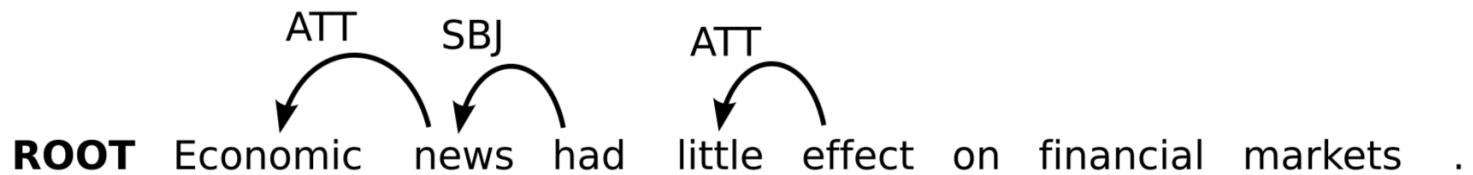
action: Shift



# Example Transition Sequence

[ROOT, had, little, effect]<sub>S</sub> [on, financial, markets, .]\_<sub>Q</sub>

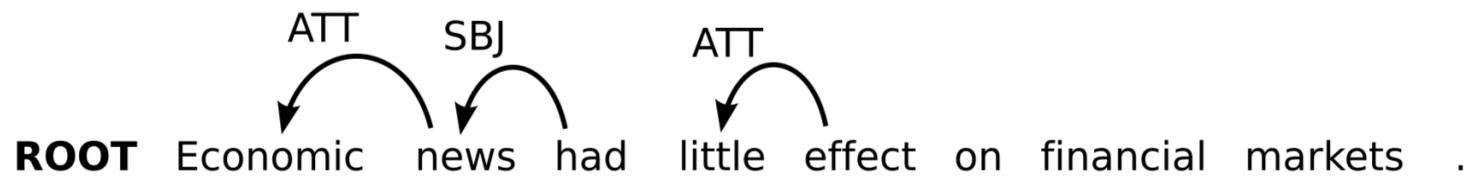
action: Left-Arc(ATT)



# Example Transition Sequence

[ROOT, had, effect, on]<sub>s</sub> [financial, markets, .]<sub>Q</sub>

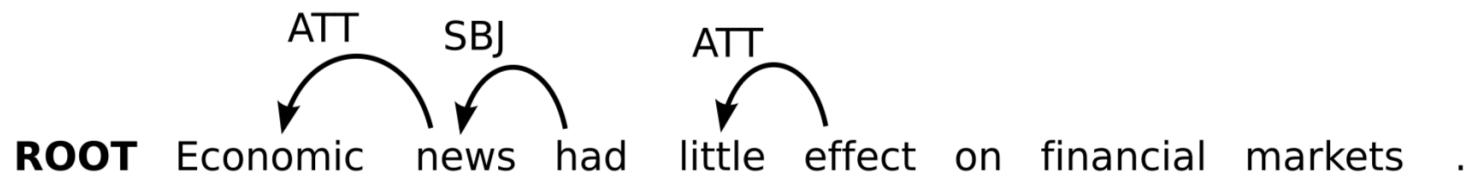
action: Shift



# Example Transition Sequence

[ROOT, had, effect, on, financial]<sub>S</sub> [markets, .]<sub>Q</sub>

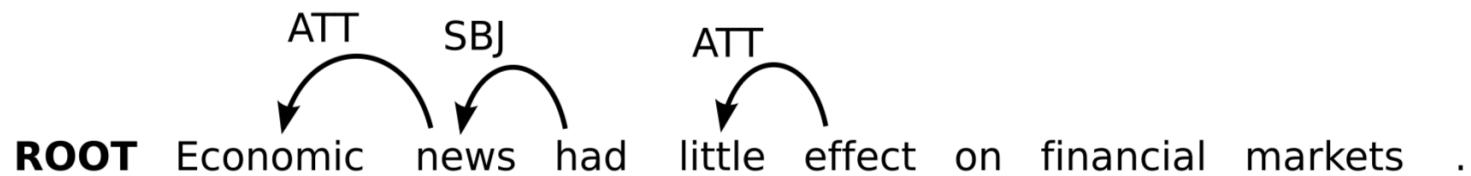
action: Shift



# Example Transition Sequence

[ROOT, had, effect, on, financial, markets]<sub>s</sub> [.]<sub>Q</sub>

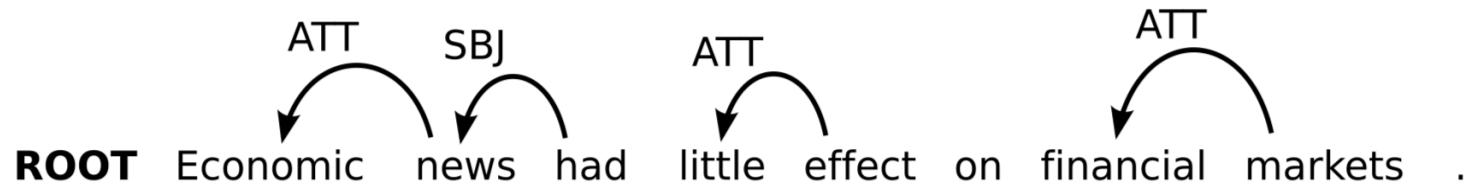
action: Shift



# Example Transition Sequence

[ROOT, had, effect, on, financial, markets]<sub>s</sub> [.]<sub>Q</sub>

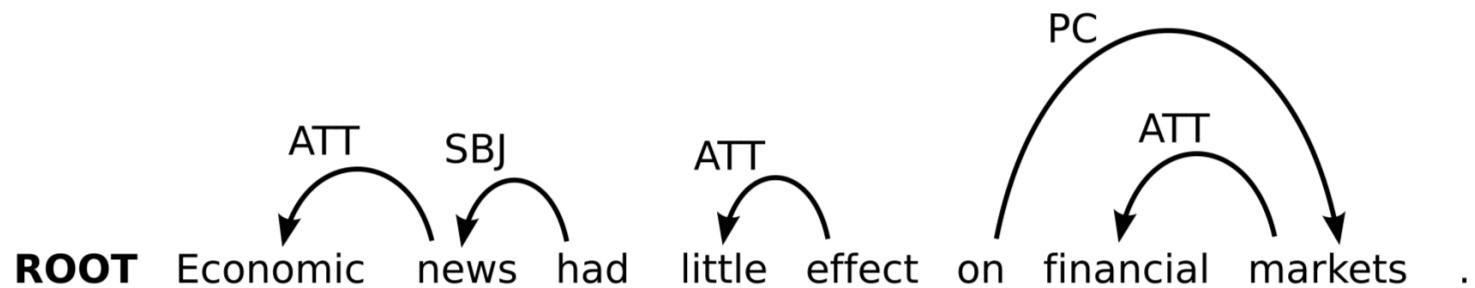
action: Left-Arc(ATT)



# Example Transition Sequence

[ROOT, had, effect, on, markets]<sub>s</sub> [.]<sub>Q</sub>

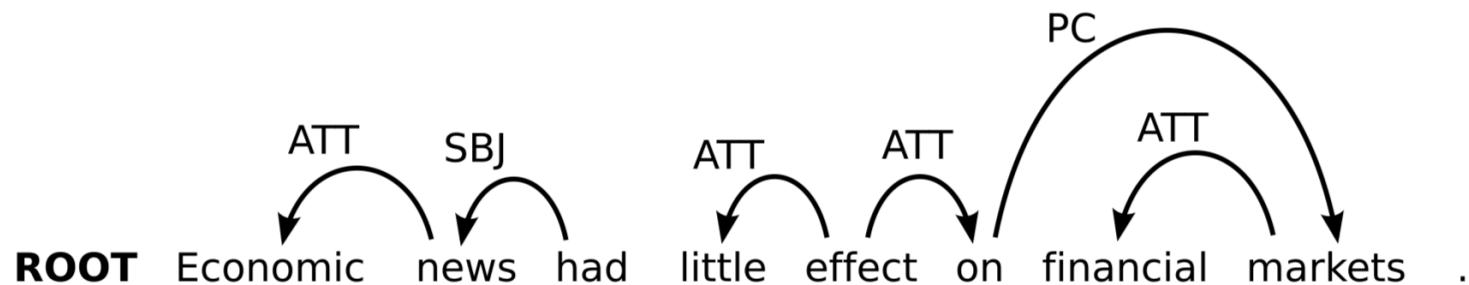
action: Right-Arc(PC)



# Example Transition Sequence

[ROOT, had, effect, on]s [.]<sub>Q</sub>

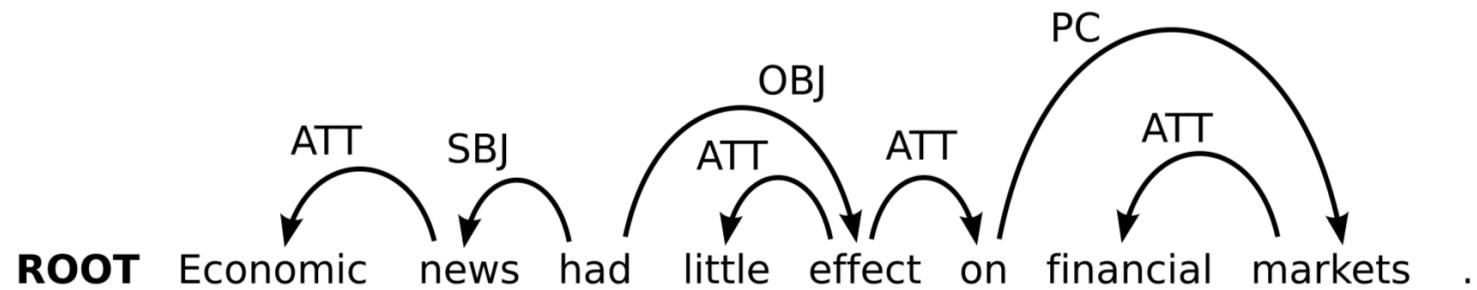
action: Right-Arc(ATT)



# Example Transition Sequence

[ROOT, had, effect]s [.]<sub>Q</sub>

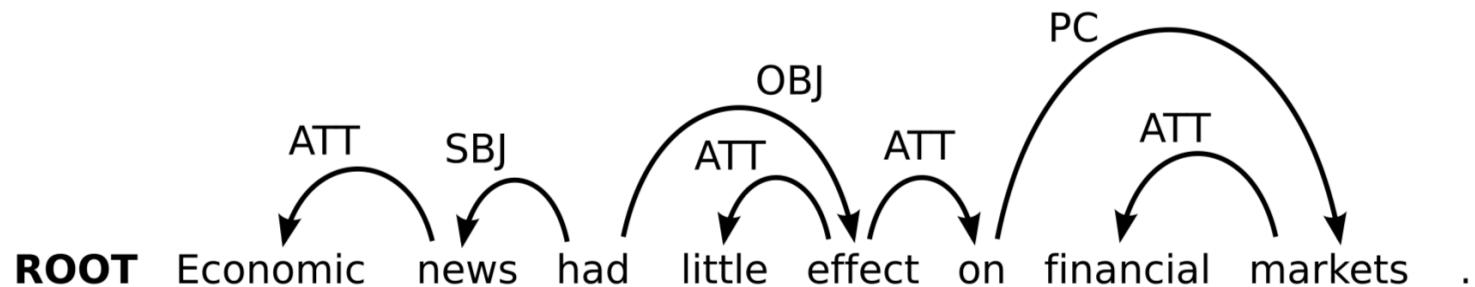
action: Right-Arc(OBJ)



# Example Transition Sequence

[ROOT, had, .]s []<sub>Q</sub>

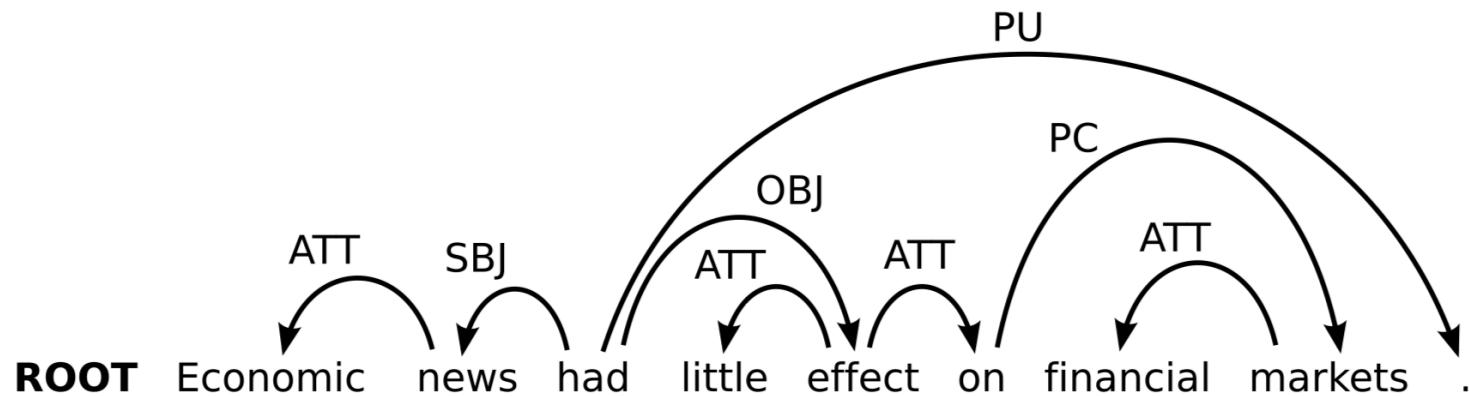
action: Shift



# Example Transition Sequence

[ROOT, had, .]s []<sub>Q</sub>

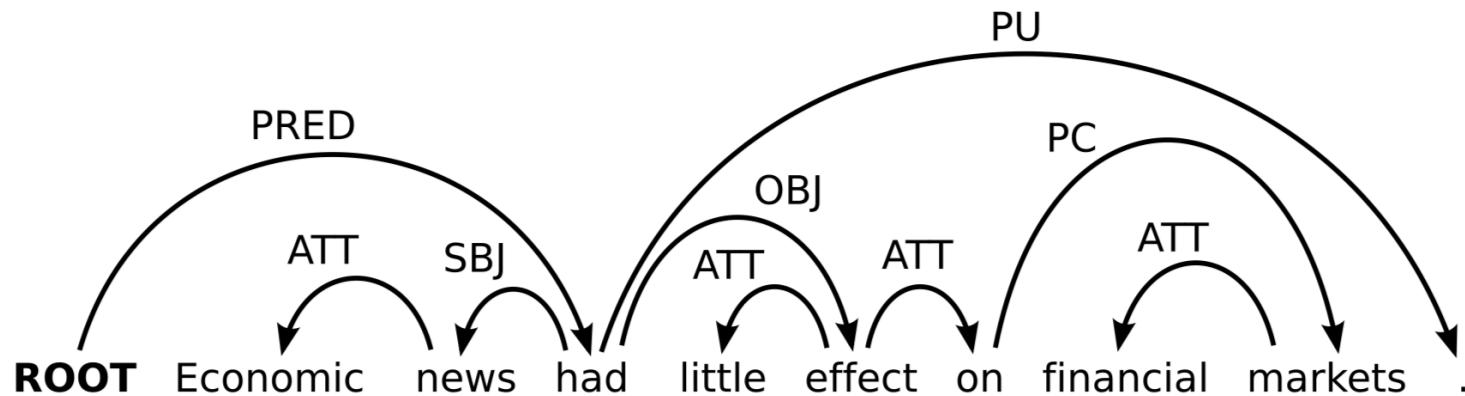
action: Right-Arc(PU)



# Example Transition Sequence

[ROOT, had]<sub>s</sub> []<sub>Q</sub>

action: Right-Arc(PRED)



# Properties of the algorithm

- Every transition sequence outputs a projective dependency tree (**soundness**).
- Every projective dependency tree is output by some transition sequence (**completeness**).
- There are exactly  $2n$  transitions in a sentence with  $n$  words.

# Deterministic parsing

- If we have an **oracle** that correctly predicts the next transition  $o(c)$ , parsing is deterministic:

```
PARSE( $w_1, \dots, w_n$ )
1    $c \leftarrow ([w_0]_S, [w_1, \dots, w_n]_Q, \{ \})$ 
2   while  $Q_c \neq []$  or  $|S_c| > 1$ 
3      $t \leftarrow o(c)$ 
4      $c \leftarrow t(c)$ 
5   return  $T = (\{w_0, w_1, \dots, w_n\}, A_c)$ 
```

# Oracles as classifiers

- An oracle can be approximated by a (linear) **classifier**:  
$$o(c) = \operatorname{argmax}_t \mathbf{w} \cdot \mathbf{f}(c, t)$$
- History-based feature representation  $\mathbf{f}(c, t)$ :
  - Features over input tokens relative to  $S$  and  $Q$
  - Features over the (partial) dependency tree defined by  $A$
  - Features over the (partial) transition sequence
- Weight vector  $\mathbf{w}$  learned from treebank data:
  - Reconstruct oracle transition sequence for each sentence
  - Construct training data set  $D = \{(c, t) \mid o(c) = t\}$
  - Maximise accuracy of local predictions  $o(c) = t$

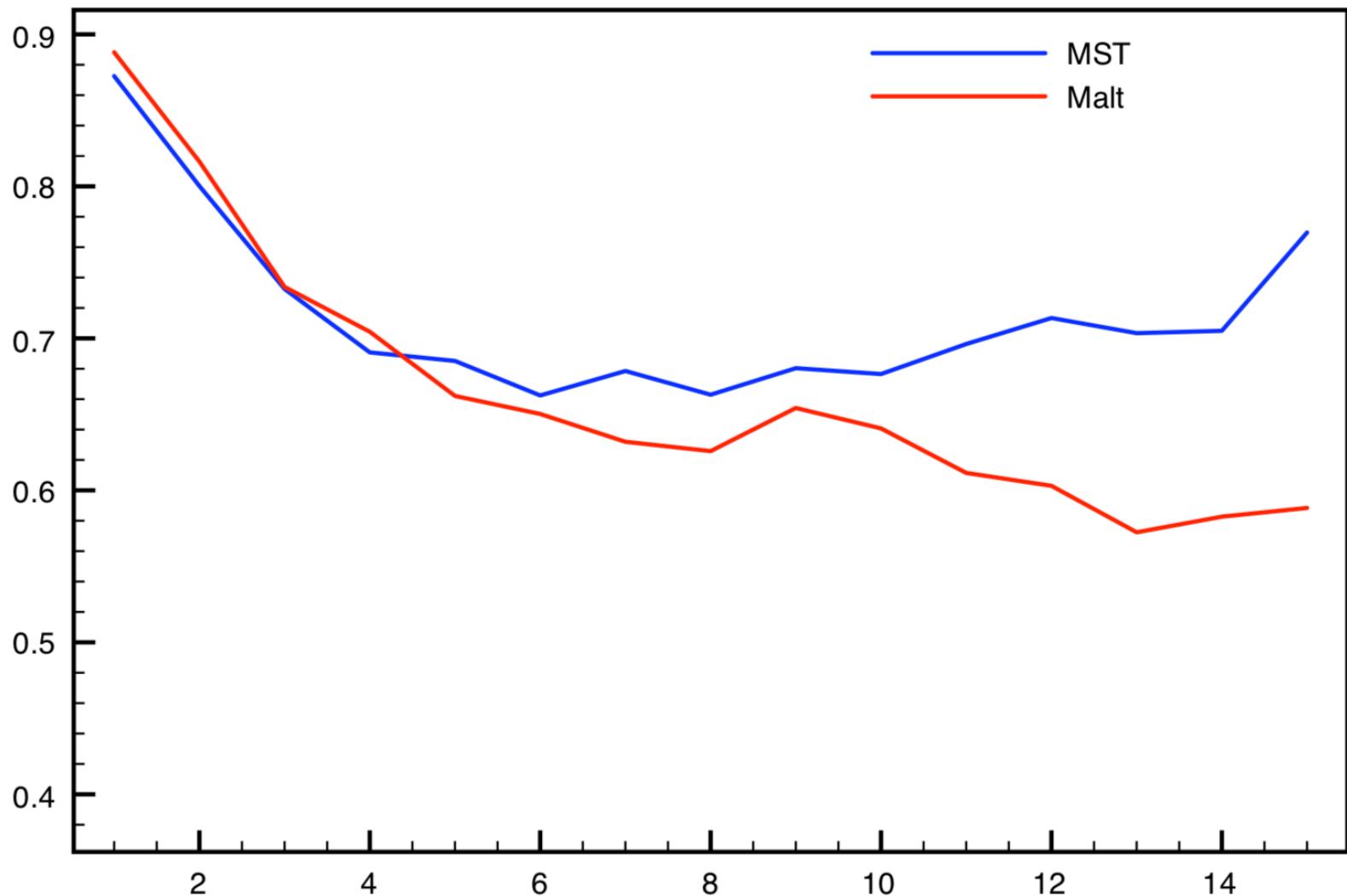
# Deterministic classifier-based parsing

- Advantages:
  - **Highly efficient parsing** – linear time complexity with constant time oracles and transitions
  - **Rich history-based feature representations** – no rigid constraints from inference algorithm
- Drawback:
  - Sensitive to **search errors** and **error propagation** due to deterministic parsing and local learning

# Empirical results: the CoNLL 2006 shared task

- CoNLL 2006 shared task (Buchholz and Marsi 2006):
  - **MaltParser** (Nivre et al. 2006) – transition-based, deterministic, local learning
  - **MSTParser** (McDonald et al. 2006) – graph-based, exact, global learning
    - Same average parsing accuracy over 13 languages
- Comparative **error analysis** (McDonald and Nivre 2007):
  - MaltParser more accurate on short dependencies and disambiguation of core grammatical functions
  - MSTParser more accurate on long dependencies and dependencies near the root of the tree
- Hypothesised **explanation for MaltParser results**:
  - Rich features counteracted by error propagation

# Precision by dependency length



# Beam search and structured prediction



# Beam search

- **Maintain the  $k$  best hypotheses** (Johansson and Nugues 2006):

```
PARSE( $w_1, \dots, w_n$ )
1   BEAM  $\leftarrow \{([w_0]_S, [w_1, \dots, w_n]_Q, \{ \ })\}$ 
2   while  $\exists c \in \text{BEAM} [Q_c \neq [] \text{ or } |S_c| > 1]$ 
3     foreach  $c \in \text{BEAM}$ 
4       foreach  $t$ 
5         ADD( $t(c)$ , NEWBEAM)
6     BEAM  $\leftarrow \text{TOP}(k, \text{NEWBEAM})$ 
7   return  $T = (\{w_0, w_1, \dots, w_n\}, A_{\text{TOP}(1, \text{BEAM})})$ 
```

- **Note:**

- $\text{Score}(c_0, \dots, c_m) = \sum_{i=1}^m \mathbf{w} \cdot \mathbf{f}(c_{j-1}, t_j)$
- Simple combination of locally normalised classifier scores
- Marginal gains in accuracy

# Online question 1

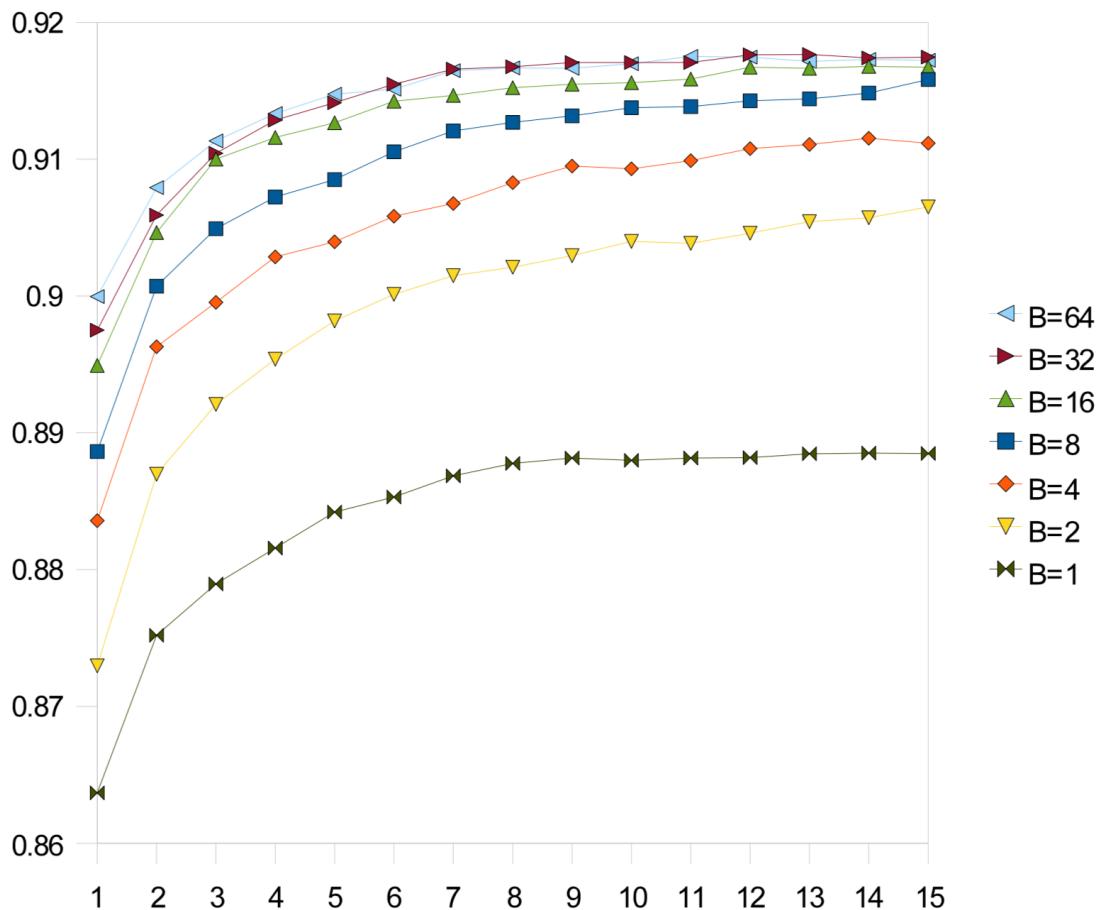
What is the complexity of transition-based parsing extended with beam search?

- 1.  $O(n)$
- 2.  $O(n \log n)$
- 3.  $O(n^2)$

# Structured prediction

- **Parsing as structured prediction** (Zhang and Clark 2008):
  - Minimise loss over entire transition sequence
  - Use beam search to find highest-scoring sequence
- Factored feature representations:
$$\mathbf{f}(c_0, \dots, c_m) = \sum_{i=1}^m \mathbf{f}(c_{i-1}, t_i)$$
- Online learning from oracle transition sequences:
  - Structured perceptron (Collins 2002)
  - Early updates (Collins and Roark 2004)

# Beam size and training iterations

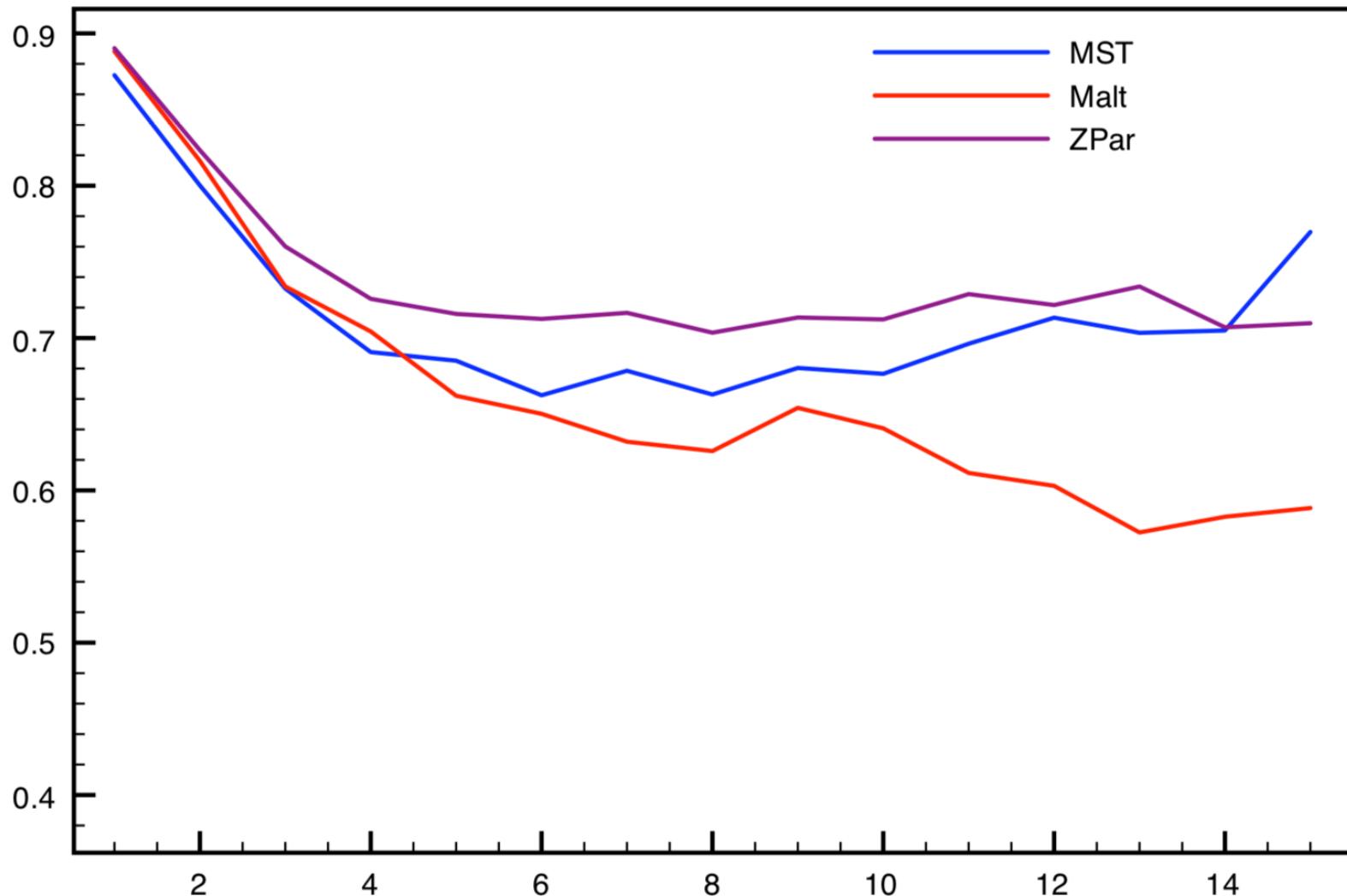


Yue Zhang and Stephen Clark. 2008. A Tale of Two Parsers: Investigating and Combining Graph-Based and Transition-Based Dependency Parsing. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, 562–571.

# The best of two worlds?

- Like graph-based dependency parsing (**MSTParser**):
  - Global learning – minimise loss over entire sentence
  - Non-greedy search – accuracy increases with beam size
- Like deterministic transition-based parsing (**MaltParser**):
  - Highly efficient – complexity still linear for fixed beam size
  - Rich features – no constraints from parsing algorithm
- Example ZPar parser (Zhang and Clark 2011)
  - “Most heavily developed for English and Chinese”

# Precision by dependency length, again



# Even richer feature models

	ZPar	Malt
Baseline	92.18	89.37
+distance	+0.07	-0.14
+valency	+0.24	0.00
+unigrams	+0.40	-0.29
+third-order	+0.18	0.00
+label set	+0.07	+0.06
Extended	93.14	89.00

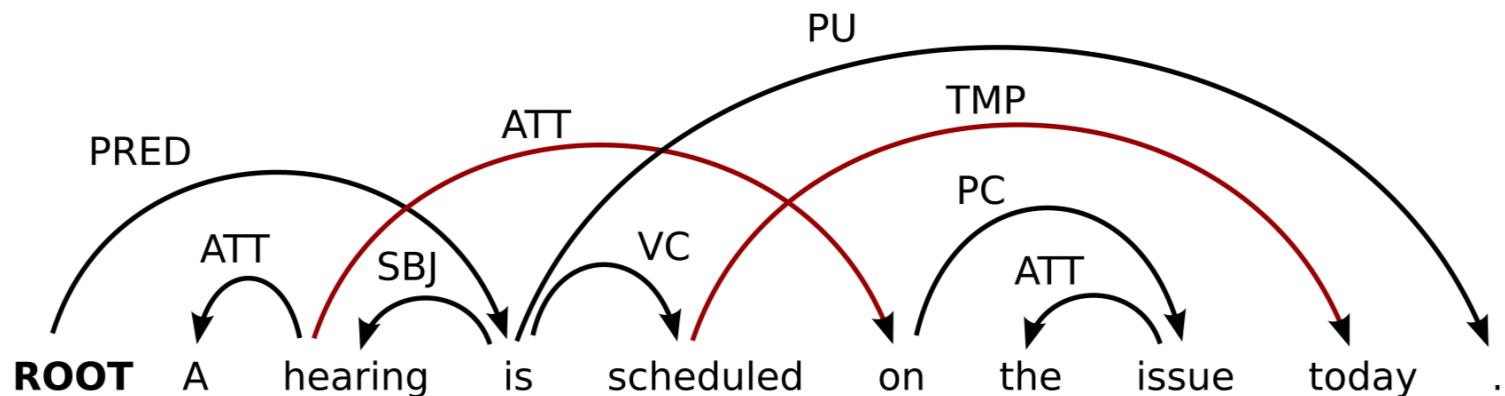
Yue Zhang and Joakim Nivre. 2011. Transition-Based Dependency Parsing with Rich Non-Local Features.  
In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 188–193.

# Online reordering for non-projectivity



# Projectivity

- A dependency arc is **projective** if the head (transitively) dominates all intervening words
- Most dependency grammar theories do not assume projectivity (but many parsers do)

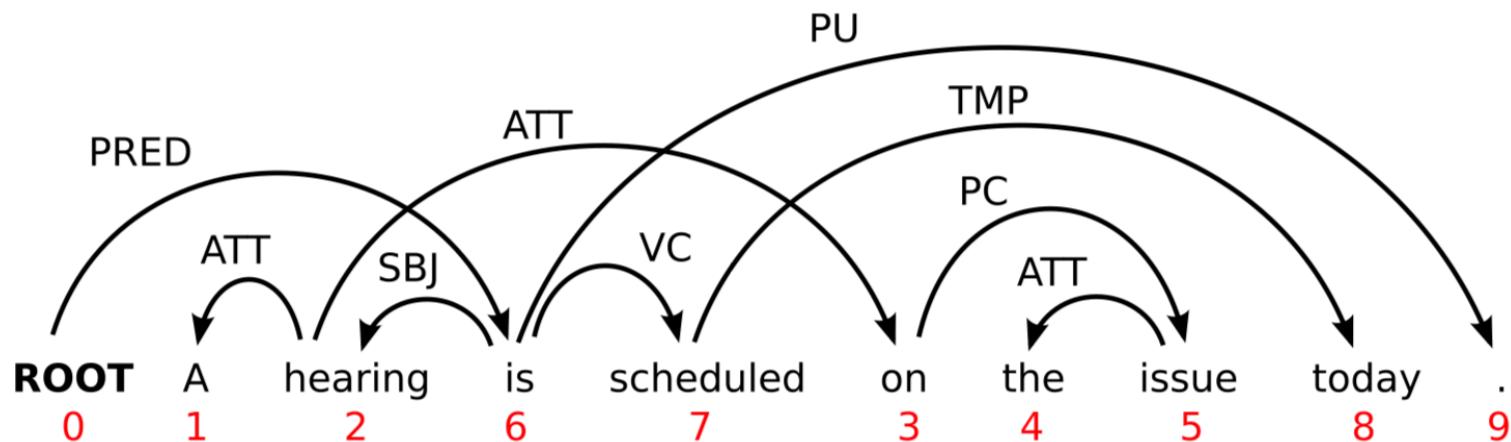


# Non-projectivity in natural languages

Language	Trees	Arcs
Arabic (Hajič et al. 2004)	11,2 %	0,4 %
Basque (Aduriz et al. 2003)	26,2 %	2,9 %
Czech (Hajič et al. 2001)	23,2 %	1,9 %
Danish (Kromann 2003)	15,6 %	1,0 %
Greek (Prokopidis et al. 2005)	20,3 %	1,1 %
Russian (Boguslavsky et al. 2000)	10,6 %	0,9 %
Slovene (Džeroski et al. 2006)	22,2 %	1,9 %
Turkish (Oflazer et al. 2003)	11,6 %	1,5 %

# Projectivity and word order

- Projectivity is a property of a dependency tree only in relation to a particular word order
  - Words can always be reordered to make the tree projective
  - Given a dependency tree  $T = (V, A, <)$ , let the projective order  $<_p$  be the order defined by an inorder traversal of  $T$  with respect to  $<$  (Veselá et al. 2004)



# Parsing with online reordering

- Add transition for reordering words (Nivre 2009):

- **Swap**

$$\frac{([\dots, w_i, w_j]_S, [\dots]_Q, A)}{([\dots, w_j]_S, [w_i, \dots]_Q, A)} \quad [0 < i < j]$$

- Transition-based parsing with two interleaved processes:
  - Sort words into projective order  $<_p$
  - Build dependency tree  $T$  by connecting adjacent subtrees
    - $T$  is always projective with respect to  $<_p$
    - $T$  may be non-projective with respect to  $<$

# Example Transition Sequence

[ROOT]<sub>s</sub> [A, hearing, is, scheduled, on, the, issue, today, .] <sub>Q</sub>

**ROOT** A hearing is scheduled on the issue today .

# Example Transition Sequence

[ROOT, A]s [hearing, is, scheduled, on, the, issue, today, .]q

**action: Shift**

**ROOT A hearing is scheduled on the issue today .**

# Example Transition Sequence

[ROOT, A, hearing]<sub>s</sub> [is, scheduled, on, the, issue, today, .] <sub>Q</sub>

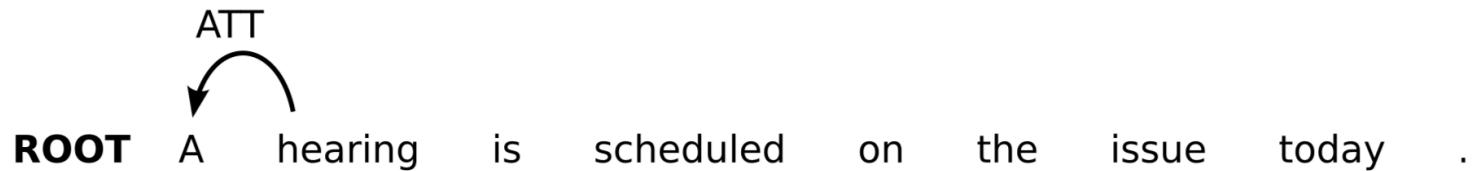
**action: Shift**

**ROOT A hearing is scheduled on the issue today .**

# Example Transition Sequence

[ROOT, A, hearing]<sub>s</sub> [is, scheduled, on, the, issue, today, .] <sub>Q</sub>

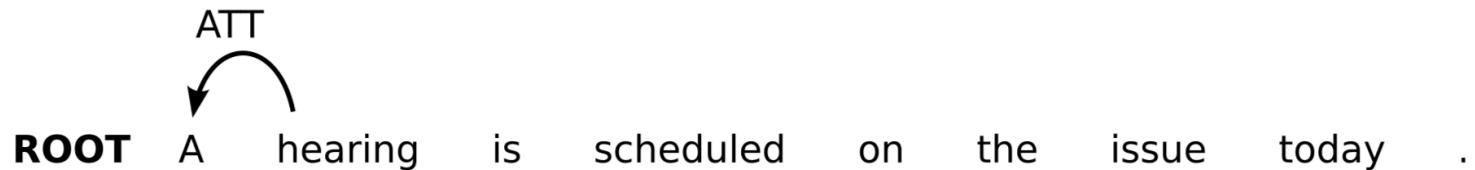
**action: Left-Arc(ATT)**



# Example Transition Sequence

[ROOT, hearing, is]<sub>s</sub> [scheduled, on, the, issue, today, .]<sub>Q</sub>

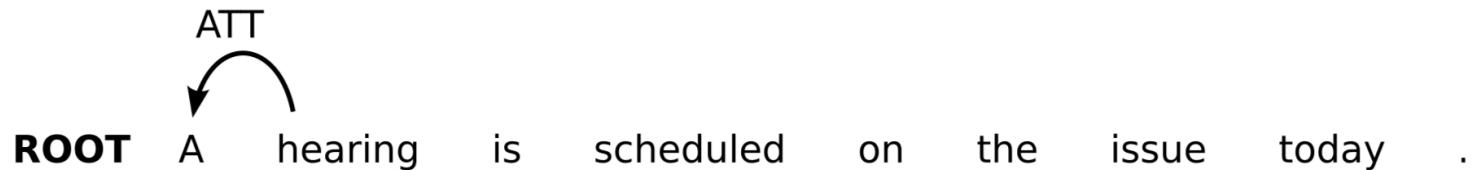
action: Shift



# Example Transition Sequence

[ROOT, hearing, is, scheduled]<sub>S</sub> [on, the, issue, today, .]<sub>Q</sub>

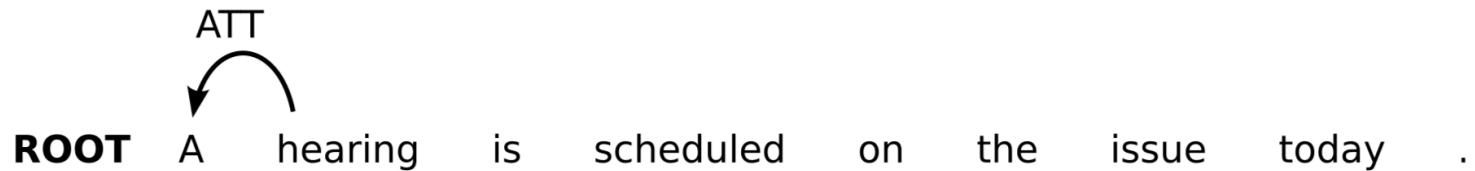
action: Shift



# Example Transition Sequence

[ROOT, hearing, is, scheduled, on]<sub>s</sub> [the, issue, today, .]<sub>Q</sub>

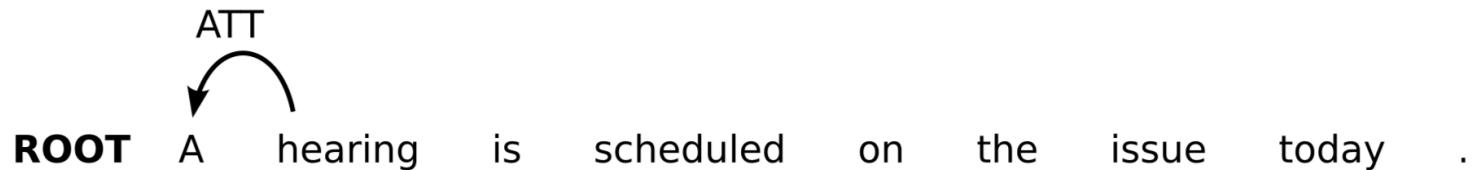
action: Shift



# Example Transition Sequence

[ROOT, hearing, is, scheduled, on, the]<sub>S</sub> [issue, today, .]<sub>Q</sub>

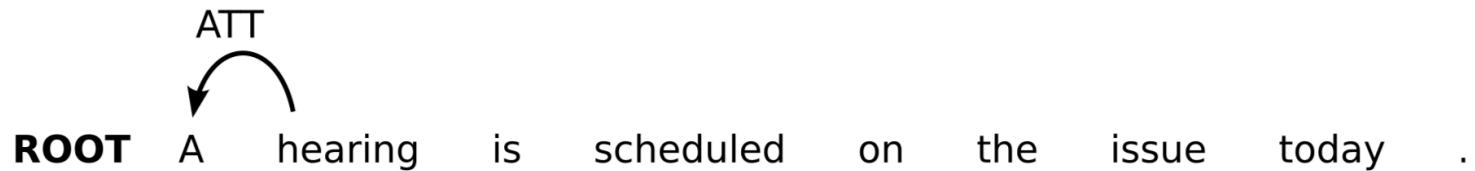
action: Shift



# Example Transition Sequence

[ROOT, hearing, is, scheduled, on, the, issue]<sub>S</sub> [today, .]Q

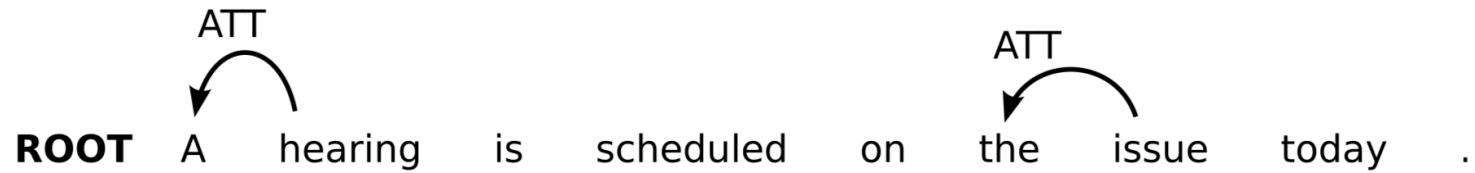
action: Shift



# Example Transition Sequence

[ROOT, hearing, is, scheduled, on, the, issue]<sub>S</sub> [today, .]<sub>Q</sub>

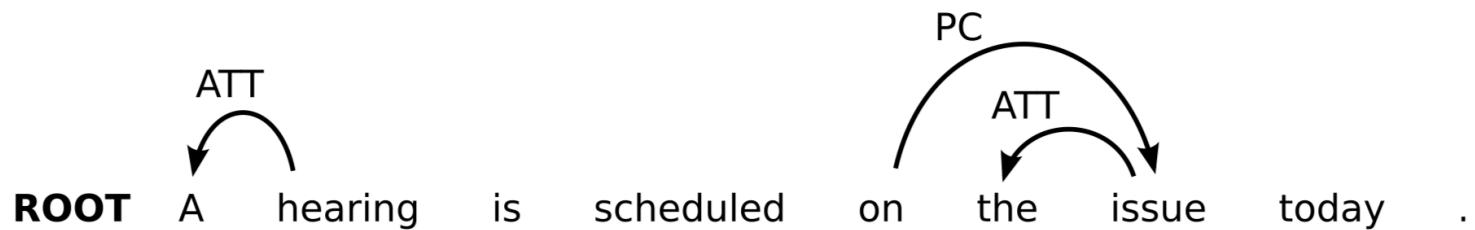
action: Left-Arc(ATT)



# Example Transition Sequence

[ROOT, hearing, is, scheduled, on, issue]s [today, .]Q

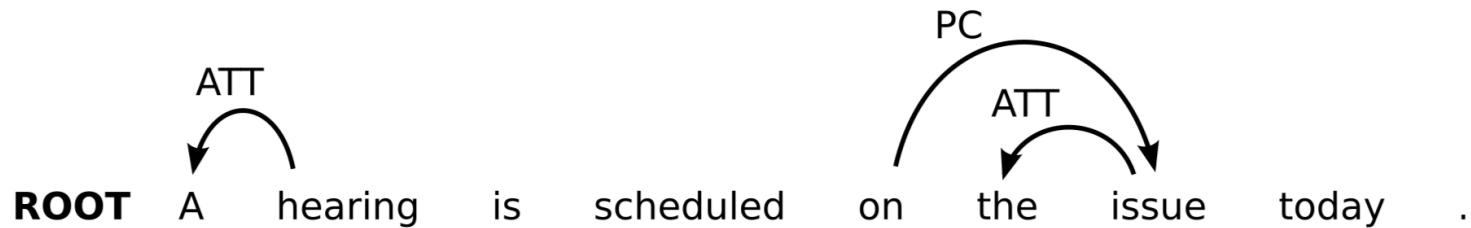
action: Right-Arc(PC)



# Example Transition Sequence

[ROOT, hearing, is, **on**]<sub>s</sub> [scheduled, today, .]<sub>q</sub>

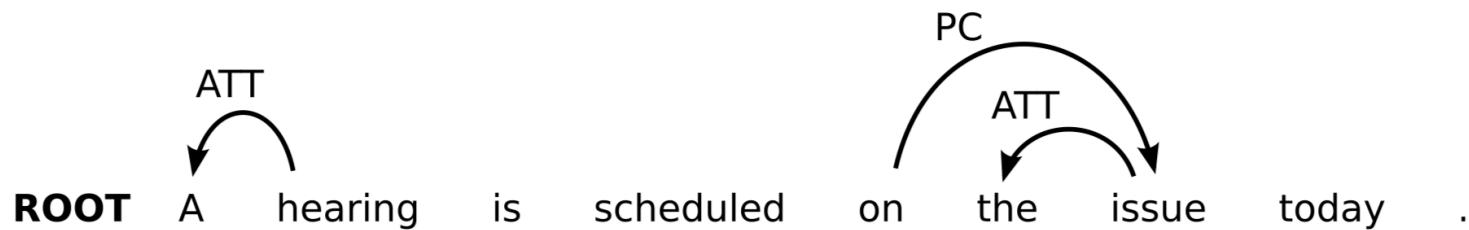
**action: Swap**



# Example Transition Sequence

[ROOT, hearing, **on**]<sub>s</sub> [**is**, scheduled, today, .]<sub>q</sub>

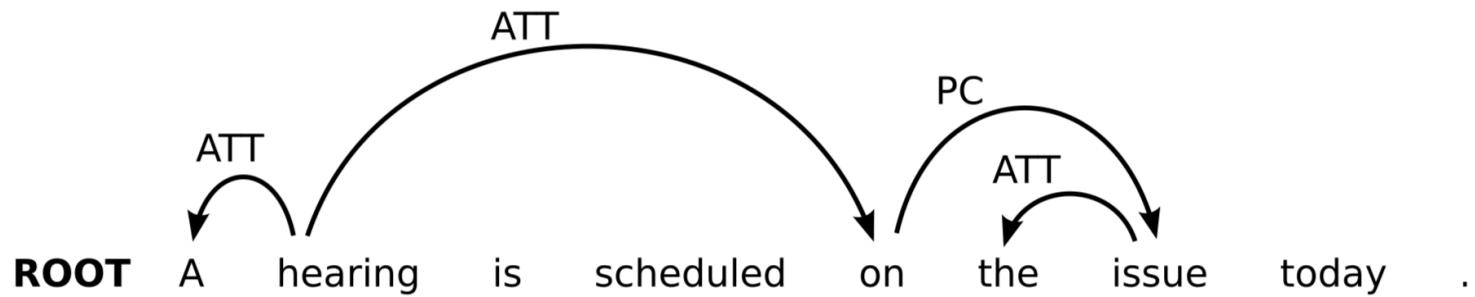
**action: Swap**



# Example Transition Sequence

[ROOT, hearing, **on**]<sub>s</sub> [is, scheduled, today, .]<sub>q</sub>

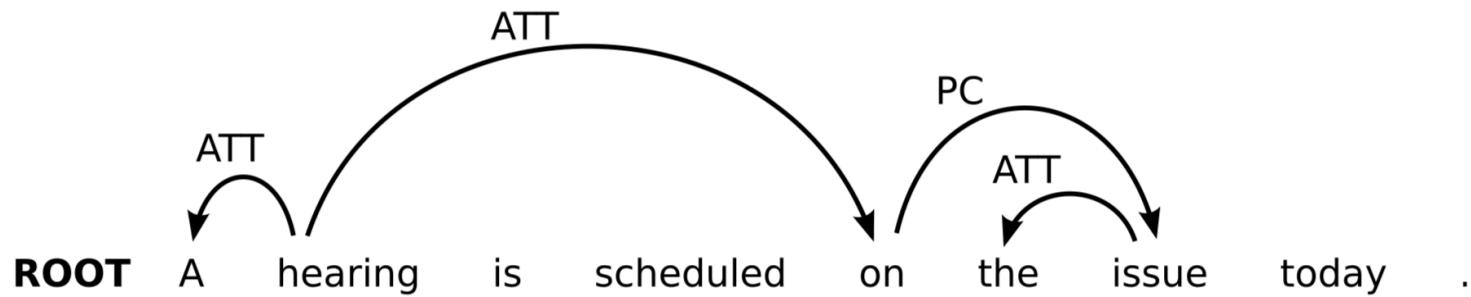
action: Right-Arc(ATT)



# Example Transition Sequence

[ROOT, hearing, is]<sub>s</sub> [scheduled, today, .]<sub>Q</sub>

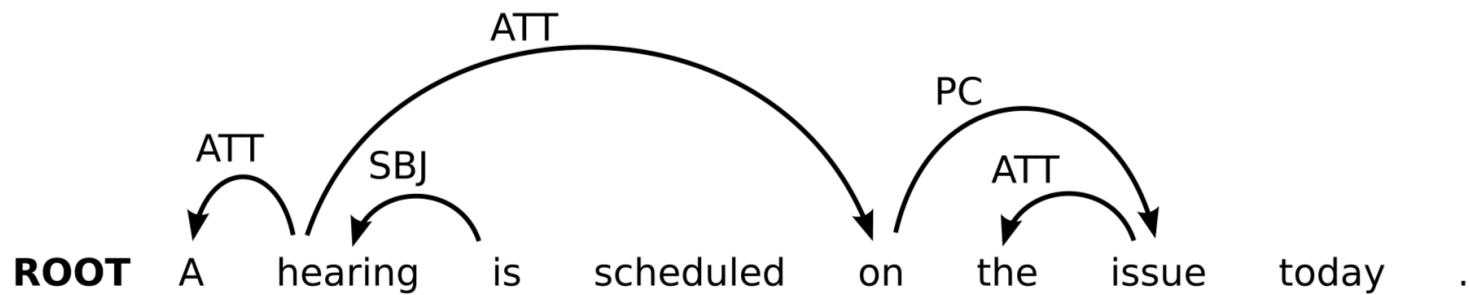
action: Shift



# Example Transition Sequence

[ROOT, hearing, is]<sub>s</sub> [scheduled, today, .]<sub>Q</sub>

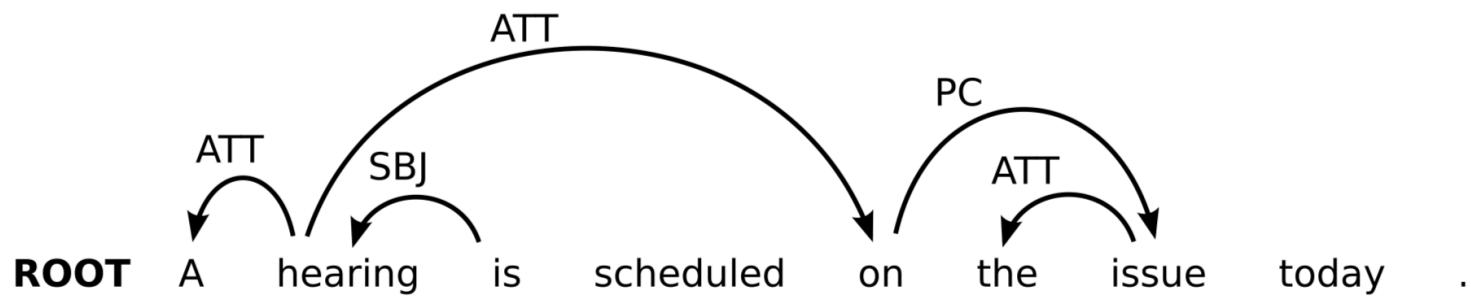
action: Left-Arc(SBJ)



# Example Transition Sequence

[ROOT, is, scheduled]<sub>S</sub> [today, .]<sub>Q</sub>

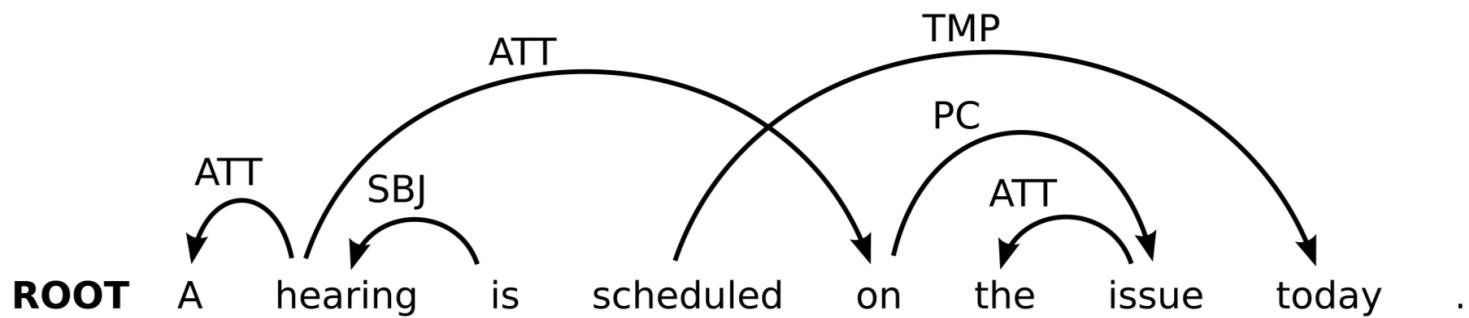
action: Shift



# Example Transition Sequence

[ROOT, is, scheduled, today]s [.]<sub>Q</sub>

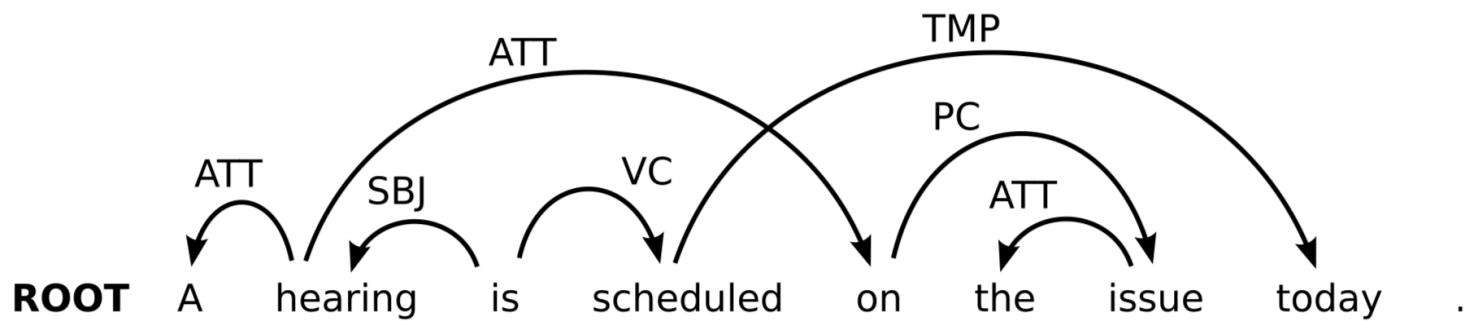
action: Right-Arc(TMP)



# Example Transition Sequence

[ROOT, is, scheduled]<sub>s</sub> [.]<sub>Q</sub>

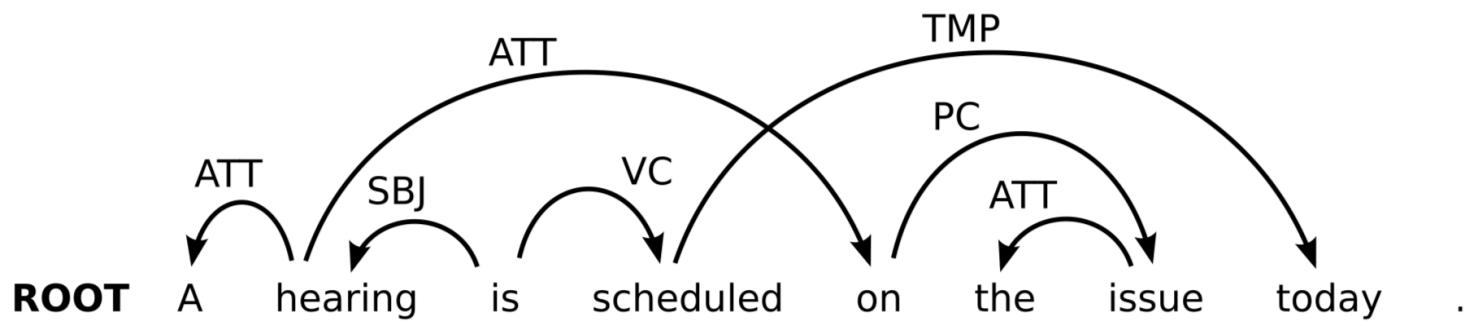
action: Right-Arc(VC)



# Example Transition Sequence

[ROOT, is, .]s []<sub>Q</sub>

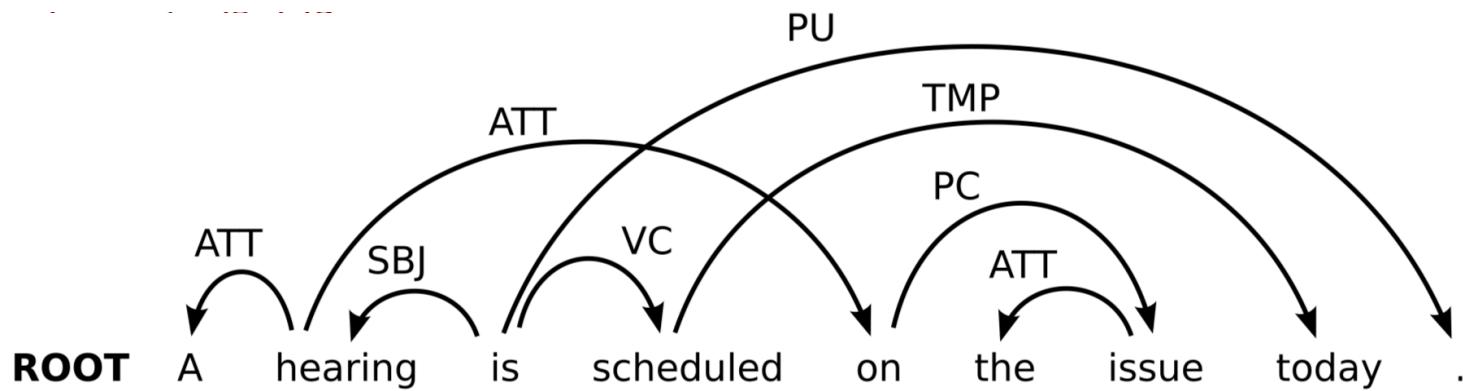
action: Shift



# Example Transition Sequence

[ROOT, is, .]s []<sub>Q</sub>

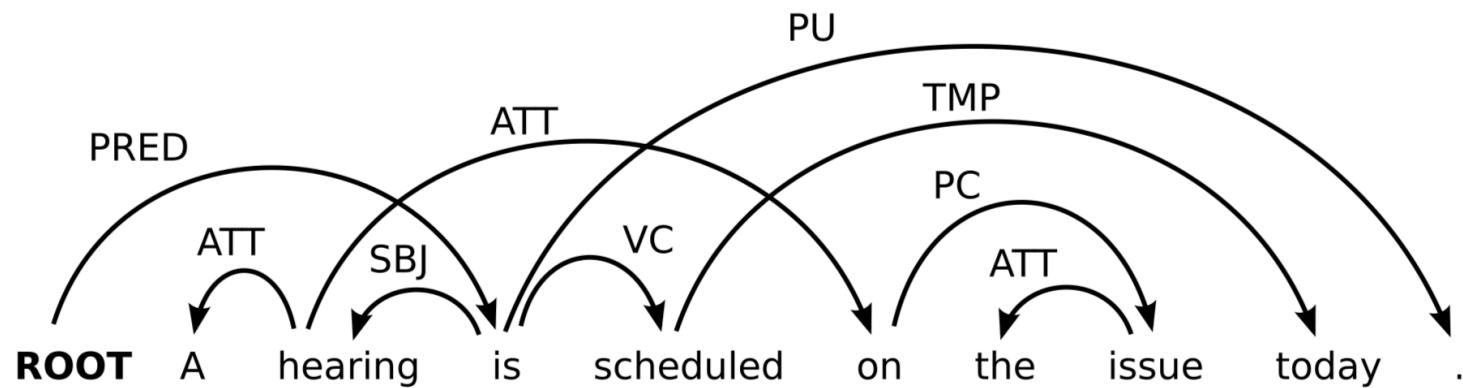
action: Right-Arc(PU)



# Example Transition Sequence

[ROOT, is]<sub>s</sub> []<sub>Q</sub>

action: Right-Arc(PRED)



# Empirical results

- Deterministic transition-based parsing (Nivre 2009):
  - Parsing in linear expected time (quadratic worst-case time)
  - Best results on Czech CoNLL 2006 data sets
- Beam search and structured prediction:
  - Evaluation on CoNLL 2009 data sets (dev sets)

	Czech		German	
	LAS	UAS	LAS	UAS
Projective	80.8	86.3	86.2	88.5
Online reordering	83.9	89.1	88.7	90.9

# Joint PoS tagging and parsing



# Part-of-Speech tagging for parsing

- Part-of-speech tags in dependency parsing:
  - Crucially assumed as input, not predicted by the parser
  - Pipeline approach may lead to error propagation
  - Most PCFG-based parsers instead predict their own tags
- Joint models for tagging and parsing:
  - Richly inflected languages (Lee et al. 2011)
  - Chinese (Li et al. 2011, Hatori et al. 2011)

# Parsing with online reordering

- Redefine Shift transition (Hatori et al. 2011)
  - **ShiftAndTag(p)** where p is a PoS tag  
 $([...]_S, [w_i, ...]_Q, A)$ 

---

 $([\dots, w_i/p]_S, [...]_Q, A)$
- Transition-based parsing now with up to three interleaved processes:
  - Tag wordforms when they are shifted onto the stack
  - Sort words into projective order  $<_p$  (if non-projectivity expected)
  - Build dependency tree T by connecting adjacent subtrees

# Empirical results

- Joint tagging and parsing with online reordering:
  - Baseline = perceptron tagger + parser
  - Up to 3 tags per word allowed in joint model
  - Beam = 40 syntactic hypotheses + 4 tag variants
  - Evaluation on CoNLL 2009 data sets

	Chinese		Czech		English		German	
	LAS	TAcc	LAS	TAcc	LAS	TAcc	LAS	TAcc
Pipeline	77.0	92.8	82.5	99.1	89.3	97.6	88.1	97.2
Joint	77.3	93.3	82.6	99.3	89.7	97.8	88.3	97.8

Bernd Bohnet and Joakim Nivre. 2012. A Transition-Based System for Joint Part-of-Speech Tagging and Labeled Non-Projective Dependency Parsing. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 1455–1465.

# Online question 2

The arc-standard algorithm has several drawbacks. What is this its main drawback?

1. Segmenting the construction of the parse into individual Left-Arc()'s and Right-Arc()'s prevents a correct handling of compound words and idioms
2. The emission of Left-Arc()'s must be performed as soon as possible, which results in decisions being sometimes taken without the appropriate information
3. The emission of Right-Arc()'s is sometimes postponed, which results in decisions being sometimes taken without the appropriate information

# Arc-eager Transition-Based Parsing



# Limitations of the arc-standard algorithm

- The **arc-standard** system considered so far
  - builds a dependency tree strictly bottom-up
  - a dependency arc can only be added between two nodes if the dependent node has already found all its dependents.
  - As a consequence, it is often necessary to postpone the attachment of right dependents.
- This is a problem, as parsing decisions are easier to take when the governor and the governee of a dependency are immediately accessible

# The problem of arc-standard on an example

[ROOT]<sub>s</sub> [La, température, a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, La]<sub>S</sub> [température, a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, La, température]<sub>S</sub> [a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, La, température]<sub>S</sub> [a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Left-Arc()**

La température a un très gros effet sur la concentration.



# The problem of arc-standard on an example

[ROOT, température, a]<sub>s</sub> [un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.



# The problem of arc-standard on an example

[ROOT, température, a]<sub>s</sub> [un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Left-Arc()**



# The problem of arc-standard on an example

[ROOT, a, un]<sub>S</sub> [très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, a, un, très]<sub>S</sub> [gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, a, un, très, gros]<sub>S</sub> [effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, a, un, très, gros]<sub>S</sub> [effet, sur, la, concentration]<sub>Q</sub>

**action: Right-Arc()**



# The problem of arc-standard on an example

[ROOT, a, un, gros, effet]<sub>S</sub> [sur, la, concentration]<sub>Q</sub>

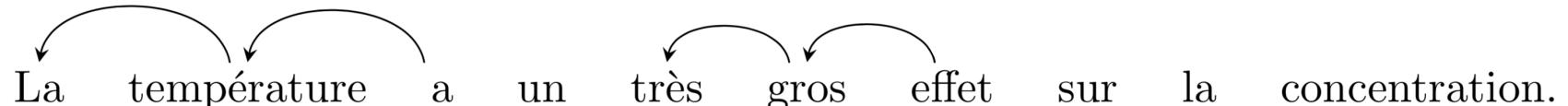
**action: Shift**

La température a un très gros effet sur la concentration.

# The problem of arc-standard on an example

[ROOT, a, un, *gros*, effet]<sub>S</sub> [sur, la, concentration]<sub>Q</sub>

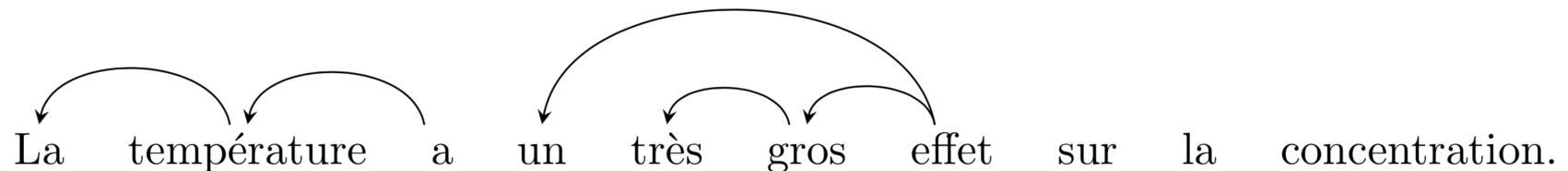
**action: Left-Arc()**



# The problem of arc-standard on an example

[ROOT, a, un, effet]<sub>S</sub> [sur, la, concentration]<sub>Q</sub>

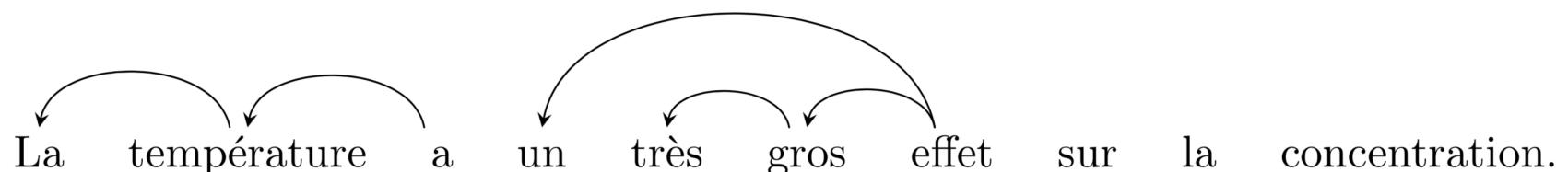
**action: Left-Arc()**



# The problem of arc-standard on an example

[ROOT, a, effet, sur]<sub>S</sub> [la, concentration]<sub>Q</sub>

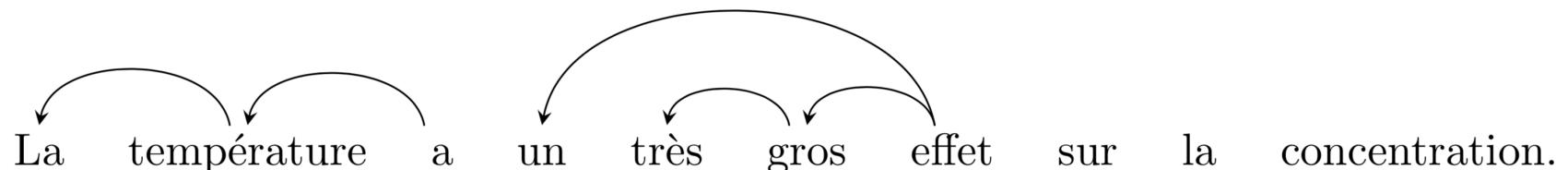
**action: Shift**



# The problem of arc-standard on an example

[ROOT, a, effet, sur, la]<sub>s</sub> [concentration]<sub>Q</sub>

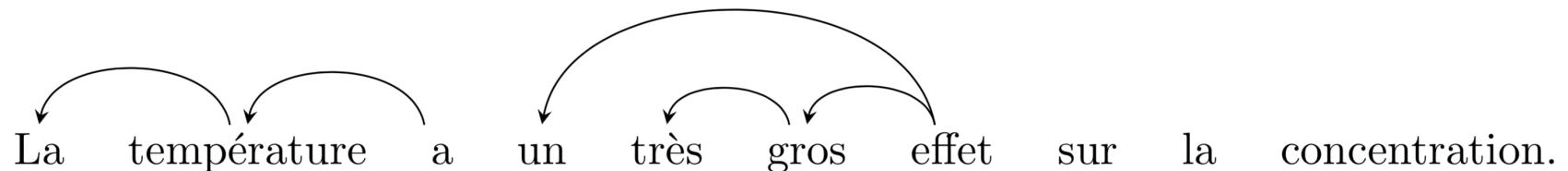
**action: Shift**



# The problem of arc-standard on an example

[ROOT, a, effet, sur, la, concentration]<sub>s</sub> []<sub>Q</sub>

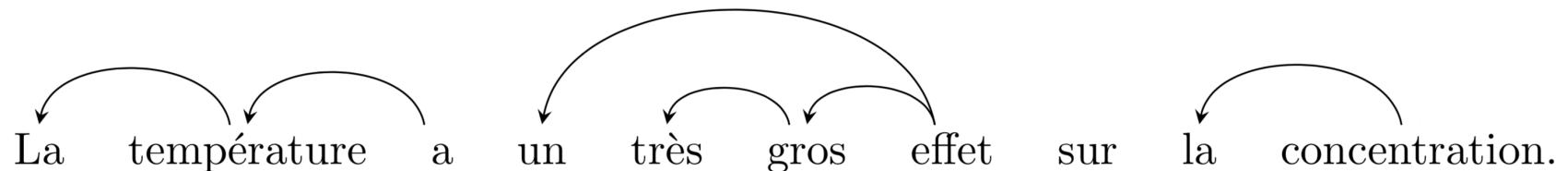
**action: Shift**



# The problem of arc-standard on an example

[ROOT, a, effet, sur, la, concentration]<sub>s</sub> []<sub>Q</sub>

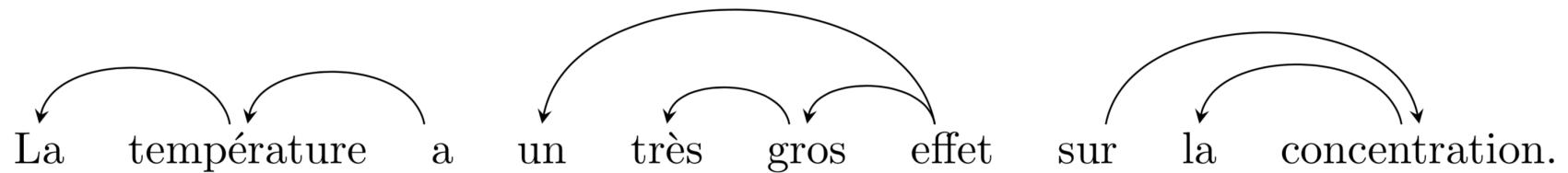
**action: Left-Arc()**



# The problem of arc-standard on an example

[ROOT, a, effet, sur, concentration]<sub>s</sub> [ ]<sub>Q</sub>

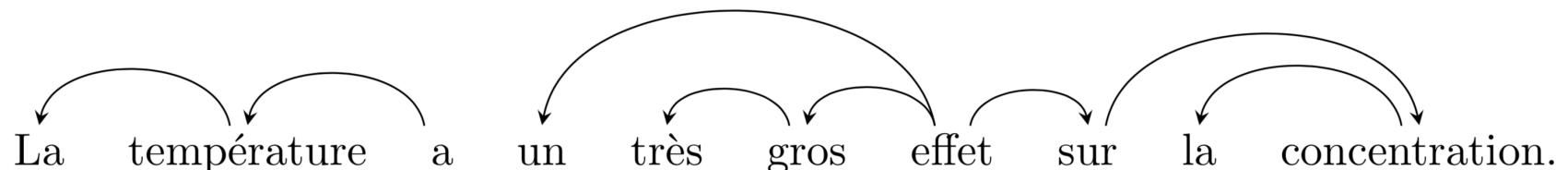
**action: Right-Arc()**



# The problem of arc-standard on an example

[ROOT, a, effet, sur]<sub>s</sub> []<sub>Q</sub>

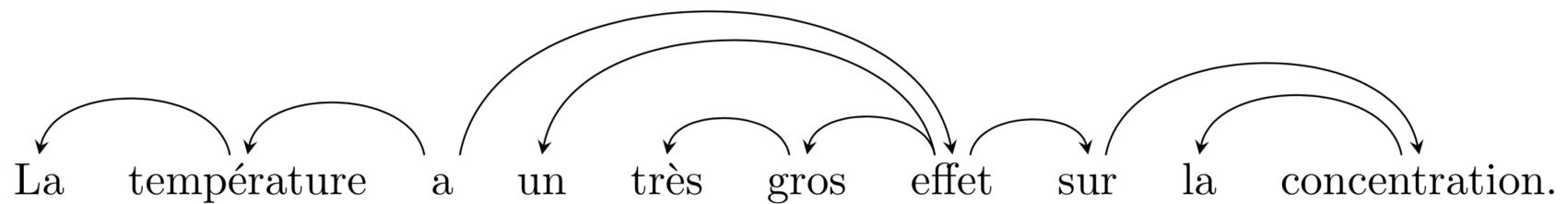
**action: Right-Arc()**



# The problem of arc-standard on an example

[ROOT, a, effet]<sub>s</sub> []<sub>q</sub>

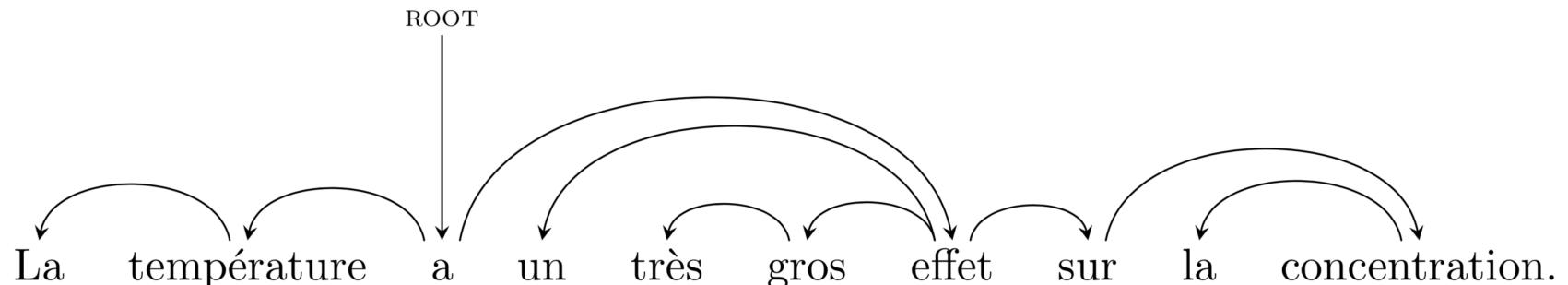
**action: Right-Arc()**



# The problem of arc-standard on an example

[ROOT, a]<sub>s</sub> []<sub>Q</sub>

**action: Right-Arc()**



# The arc-eager system

- We will modify the basic set of actions in order to always add an arc at the earliest possible opportunity:
  - we will now build parts of the tree top-down instead of bottom-up
- Shift remains the same
- **Left-Arc is rewritten and subjected to a stricter condition** (allowed only if the dependent is not the root and has no incoming arcs)

$$([\dots, w_i]_S, [w_j, \dots]_Q, A)$$
$$\frac{}{[i \neq 0 \wedge \nexists(k, l') \mid (k, l', i) \in A]}$$
$$([\dots]_S, [w_j, \dots]_Q, A \cup \{(w_j, l, w_i)\})$$

# The arc-eager system

- **Right-Arc is changed:** it does not discard  $w_i$  anymore:

$$([\dots, w_i]_S, [w_j, \dots]_Q, A)$$

---

$$([\dots, w_i, w_j]_S, [\dots]_Q, A \cup \{(w_i, l, w_j)\})$$

- We postpone the reduction of  $w_i$  to another, new action:
- **Reduction**, only possible if the top of the stack already has a head

$$([\dots, w_i]_S, Q, A)$$

---

only if  $\exists (k, l') \mid (k, l', i) \in A$

$$([\dots]_S, Q, A)$$

# Arc-eager on an example

[ROOT]<sub>s</sub> [La, température, a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

La température a un très gros effet sur la concentration.

# Arc-eager on an example

[ROOT, La]<sub>S</sub> [température, a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.

# Arc-eager on an example

[ROOT, La]<sub>s</sub> [température, a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Left-Arc()**

La température a un très gros effet sur la concentration.

# Arc-eager on an example

[ROOT, température]<sub>s</sub> [a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**

La température a un très gros effet sur la concentration.



# Arc-eager on an example

[ROOT, température]<sub>s</sub> [a, un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Left-Arc()**



# Arc-eager on an example

[ROOT, a]<sub>S</sub> [un, très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Right-Arc()**



# Arc-eager on an example

[ROOT, a, un]<sub>S</sub> [très, gros, effet, sur, la, concentration]<sub>Q</sub>

**action: Shift**



# Arc-eager on an example

[ROOT, a, un, très]<sub>S</sub> [gros, effet, sur, la, concentration]<sub>Q</sub>

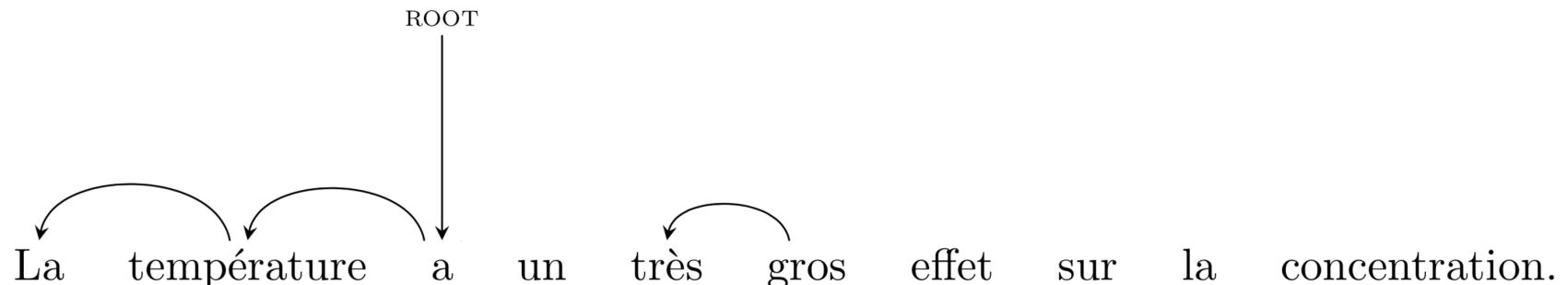
**action: Shift**



# Arc-eager on an example

[ROOT, a, un, très]<sub>s</sub> [gros, effet, sur, la, concentration]<sub>Q</sub>

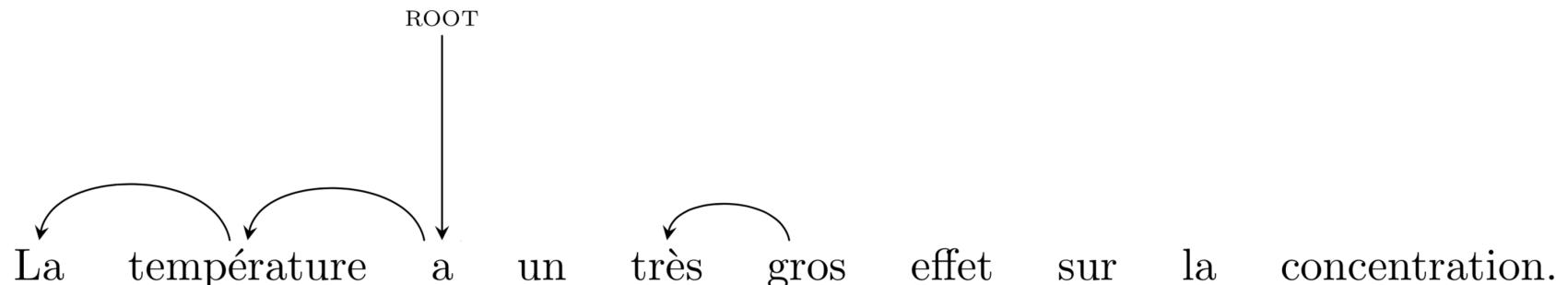
**action: Left-Arc()**



# Arc-eager on an example

[ROOT, a, un, gros]<sub>S</sub> [effet, sur, la, concentration]<sub>Q</sub>

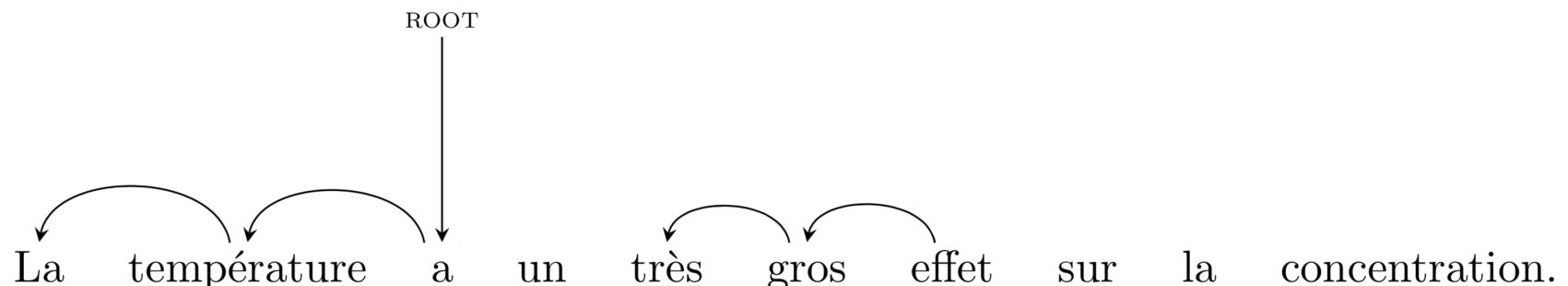
**action: Shift**



# Arc-eager on an example

[ROOT, a, un, *gros*]<sub>S</sub> [effet, sur, la, concentration]<sub>Q</sub>

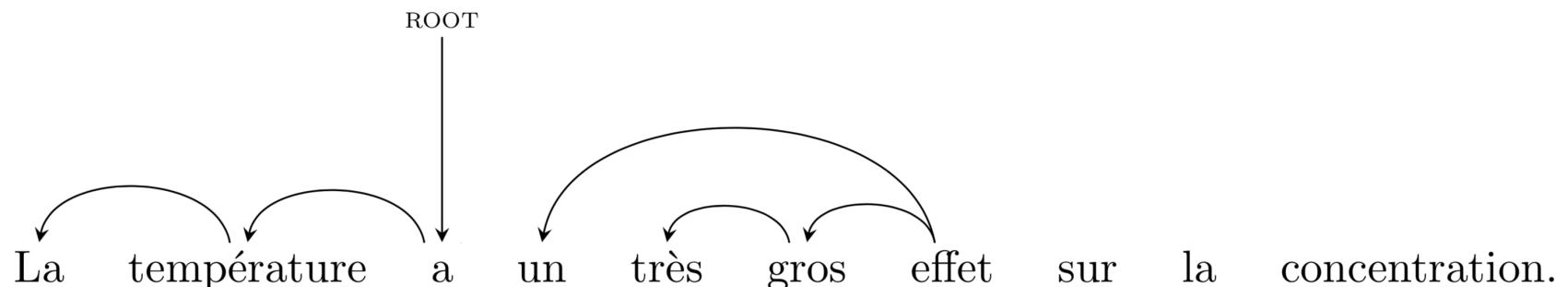
**action: Left-Arc()**



# Arc-eager on an example

[ROOT, a, un]<sub>s</sub> [effet, sur, la, concentration]<sub>Q</sub>

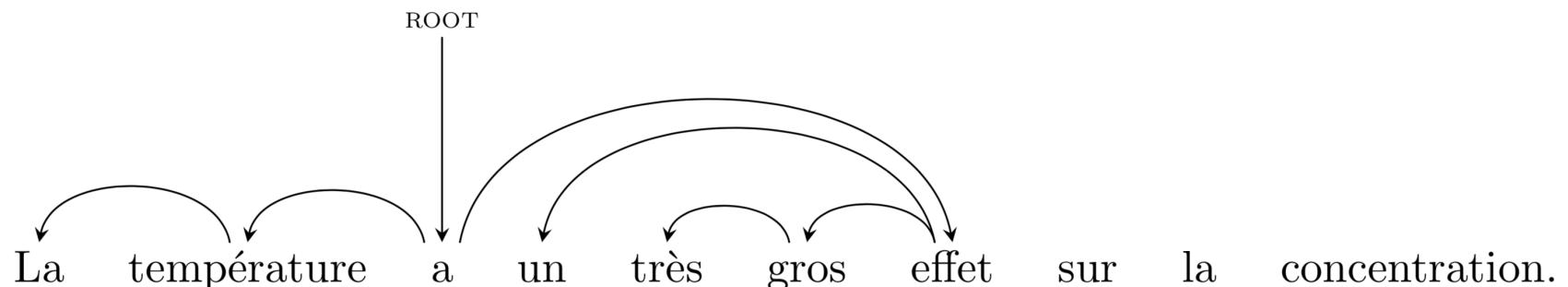
**action: Left-Arc()**



# Arc-eager on an example

[ROOT, a, effet]<sub>S</sub> [sur, la, concentration]<sub>Q</sub>

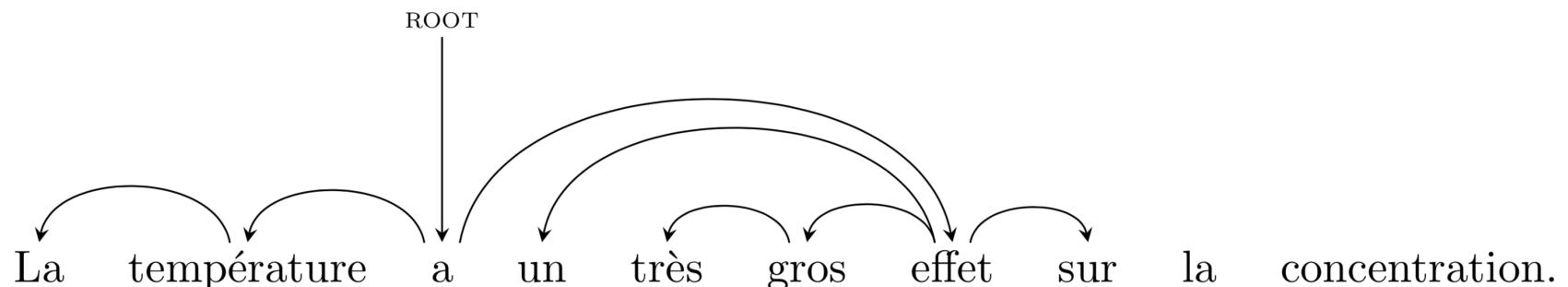
**action: Right-Arc()**



# Arc-eager on an example

[ROOT, a, effet, sur]<sub>S</sub> [la, concentration]<sub>Q</sub>

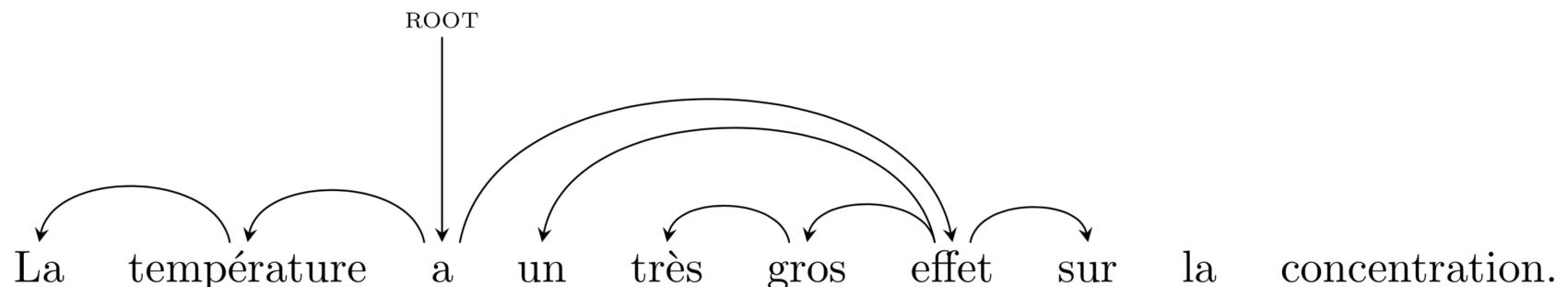
action: Right-Arc()



# Arc-eager on an example

[ROOT, a, effet, sur, la]<sub>s</sub> [concentration]<sub>Q</sub>

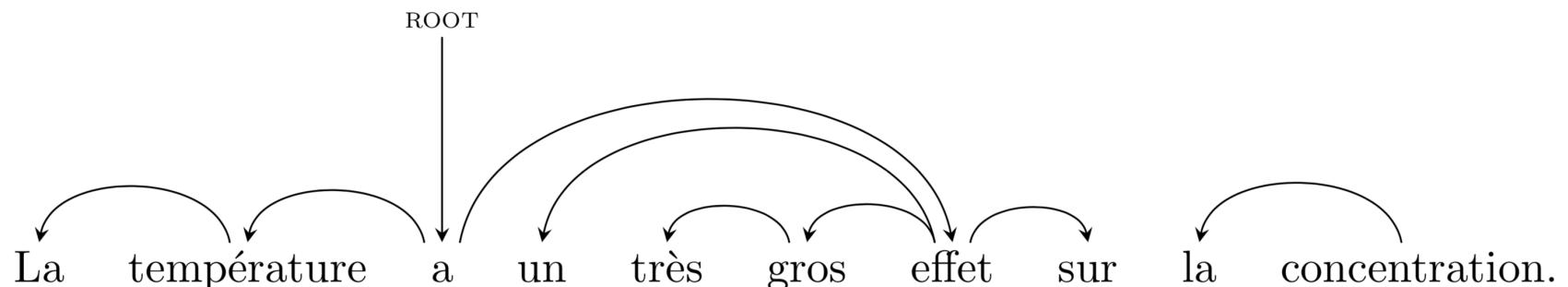
**action: Shift()**



# Arc-eager on an example

[ROOT, a, effet, sur, la]<sub>s</sub> [concentration]<sub>Q</sub>

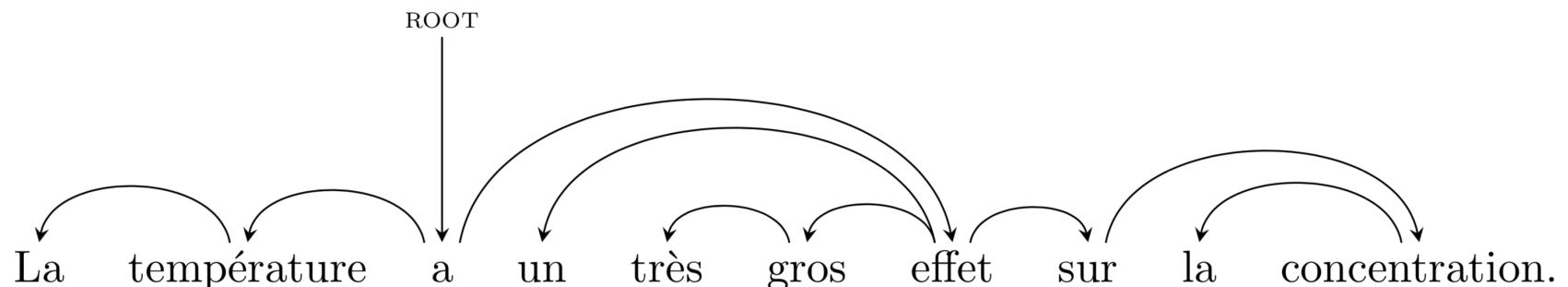
**action: Left-Arc()**



# Arc-eager on an example

[ROOT, a, effet, sur, concentration]<sub>s</sub> []<sub>Q</sub>

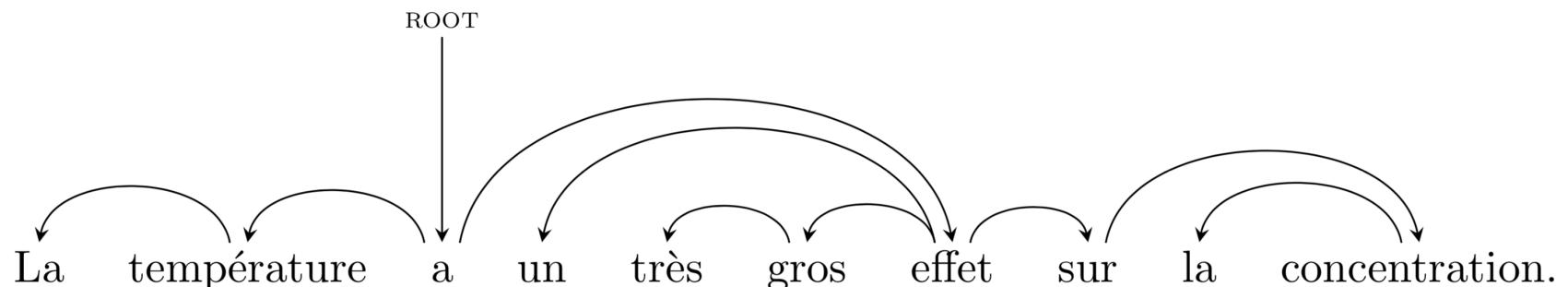
**action: Right-Arc()**



# Arc-eager on an example

[ROOT, a, effet, sur, concentration]<sub>s</sub> []<sub>Q</sub>

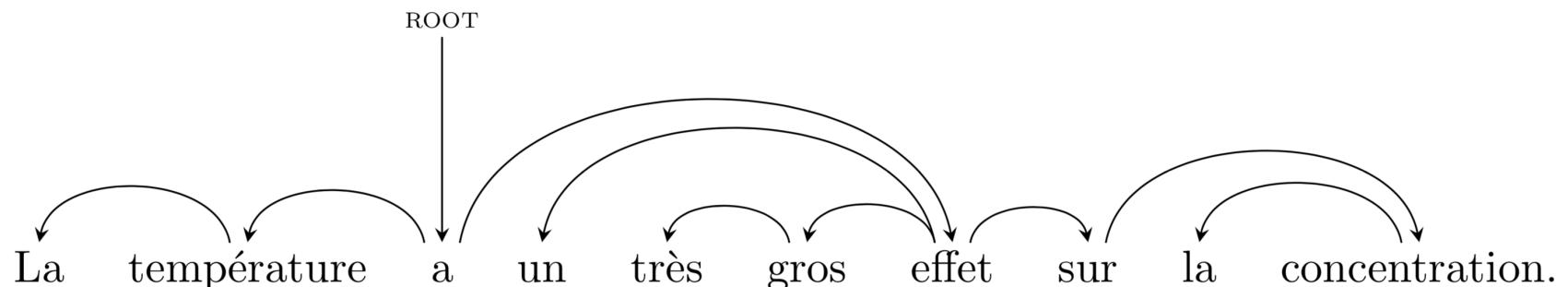
**action: Reduce**



# Arc-eager on an example

[ROOT, a, effet, sur]<sub>s</sub> []<sub>Q</sub>

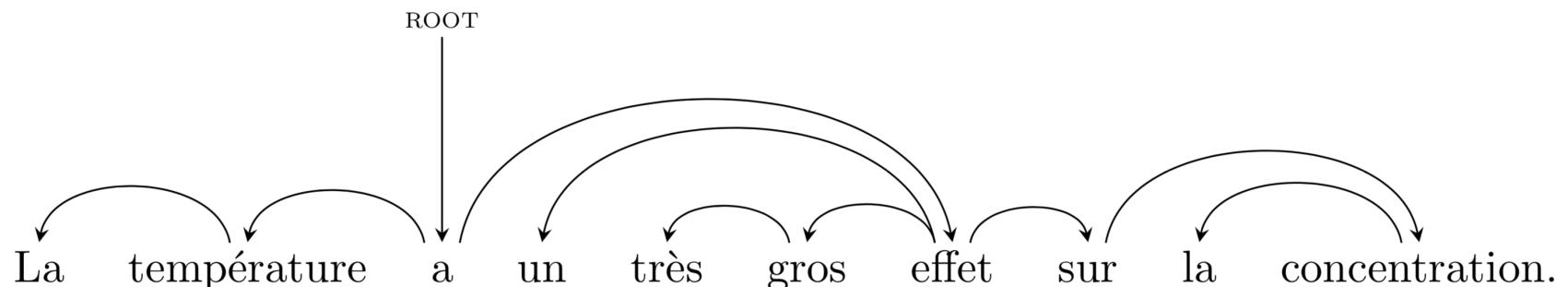
**action: Reduce**



# Arc-eager on an example

[ROOT, a, effet]<sub>s</sub> []<sub>q</sub>

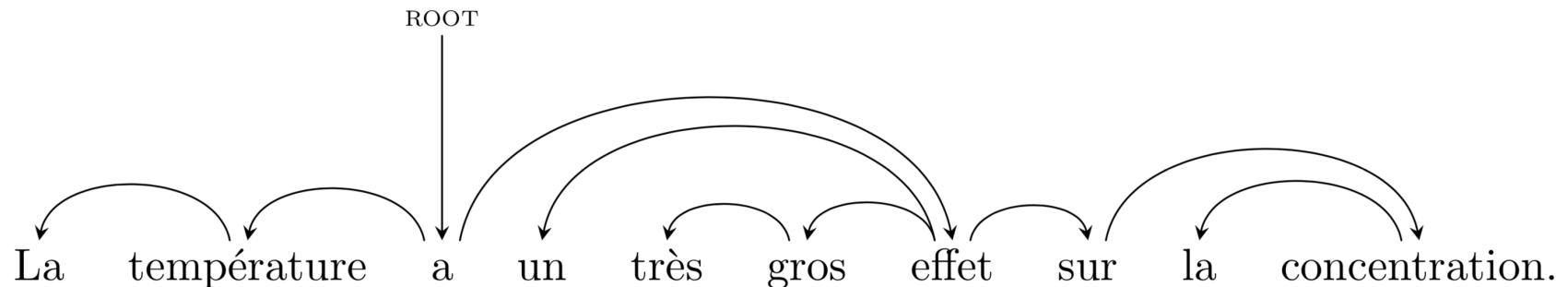
**action: Reduce**



# Arc-eager on an example

[ROOT, a]<sub>s</sub> []<sub>Q</sub>

**action: Reduce**



# Drawbacks of the arc-eager algorithm

- The arc-eager system has a weaker soundness result than the arc-standard system
- It does not guarantee the output to be a dependency tree, only a sequence of (unconnected) trees (a forest).
- In the best case, this is a sequence of length 1, meaning that the tree is in fact a tree.
- In the worst case, this is a sequence of length  $n$ , meaning that each word is its own tree.
- The arc-eager parsers normally have a last step that attaches everything that remains in the stack to the root

# Extending arc-eager: the case of disfluencies



# Drawbacks of the arc-eager algorithm

- Two issues for spoken language processing:
  - ASR errors
  - Speech Disfluencies
    - ~10% of the words in conversational speech are disfluent
    - An extreme case of “noisy” input: <http://www.youtube.com/watch?v=lj3iNxZ8Dww>
- Error propagation from these two errors can wreak havoc on downstream modules such as parsing and semantics

# Drawbacks of the arc-eager algorithm

- Two issues:
  - ASR
  - Speed
    - ~100ms
    - A few seconds
    - <http://www.cs.cmu.edu/~wcohen/>
  - Error rates have been broken and seem

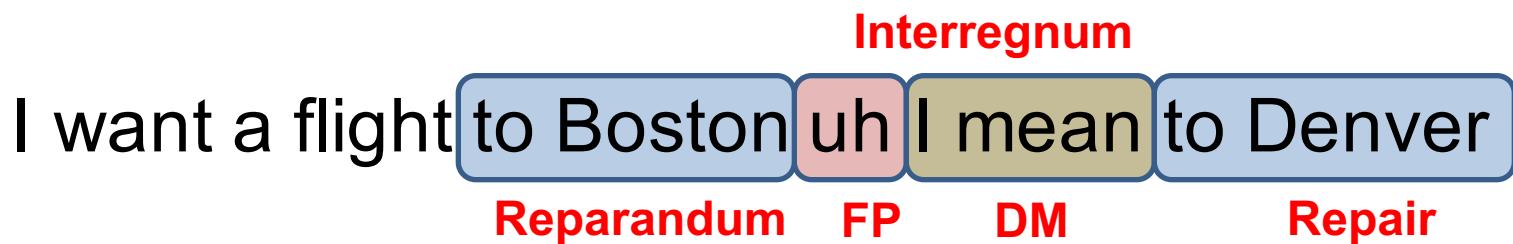


# Drawbacks of the arc-eager algorithm

- Two issues for spoken language processing:
  - ASR errors
  - Speech Disfluencies
    - ~10% of the words in conversational speech are disfluent
    - An extreme case of “noisy” input: <http://www.youtube.com/watch?v=lj3iNxZ8Dww>
- Error propagation from these two errors can wreak havoc on downstream modules such as parsing and semantics

# Disfluencies

- Three types
  - Filled pauses: e.g. uh, um
  - Discourse markers and parentheticals: e.g., I mean, you know
  - Reparandum (edited phrase)



# Processing disfluent sentences

- Most approaches deal solely with disfluency detection as a pre-processing step before parsing
- Serialised method of disfluency detection and then parsing can be slow...
- Why not parse disfluent sentences at the same time as detecting disfluencies?
  - Advantage: speed-up processing, especially for dialogue systems
- **Joint dependency parsing and disfluency detection** with high accuracy and processing speed (Rasooli and Tetrault 2013)
  - Source: *I want a flight to Boston uh I mean to Denver*
  - Output: *I want a flight to Denver*

# Rasooli and Tetrault' (2014) proposal

- Three new actions:
  - **IJ( $k$ )**: remove the first  $k$  words from the queue and tag them as interjection
  - **DM( $k$ )**: remove the first  $k$  words from the queue and tag them as discourse marker
  - **RP( $i,j$ )**: remove  $i$  to  $j$  in the stack and tag them as reparandum

## On an example

[ROOT]<sub>s</sub> [I, want, a, flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

I want a flight to Boston uh I mean to Denver

# On an example

[ROOT, I]<sub>S</sub> [want, a, flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

action: Shift()

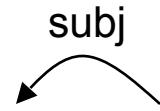
I want a flight to Boston uh I mean to Denver

# On an example

[ROOT, I]<sub>s</sub> [want, a, flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

action: Left-Arc(subj)

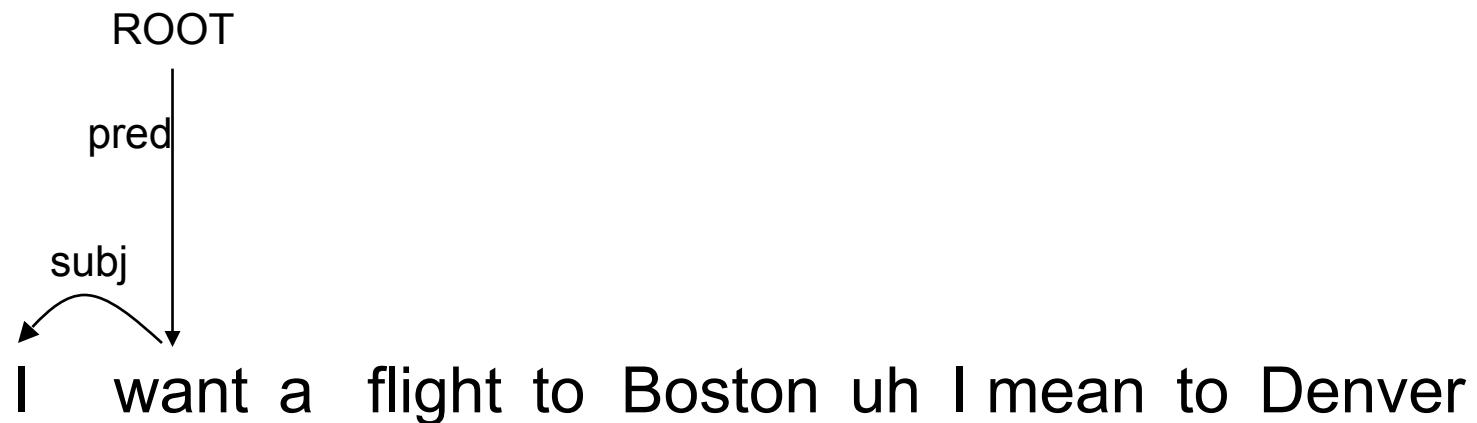
I want a flight to Boston uh I mean to Denver



# On an example

[ROOT, want]<sub>S</sub> [a, flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

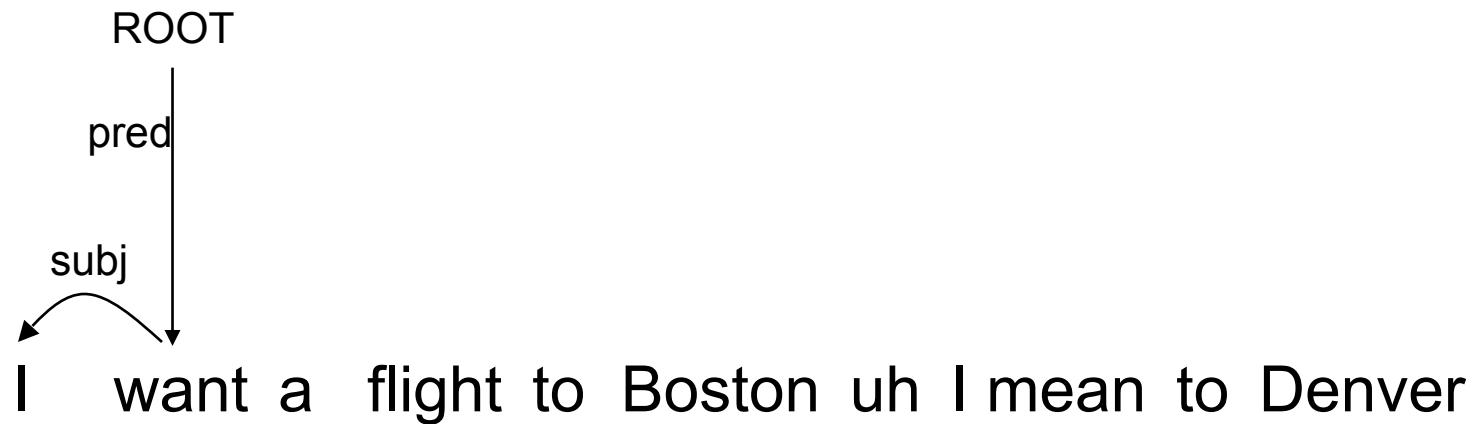
action: Right-Arc(PRED)



# On an example

[ROOT, want, a]<sub>S</sub> [flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

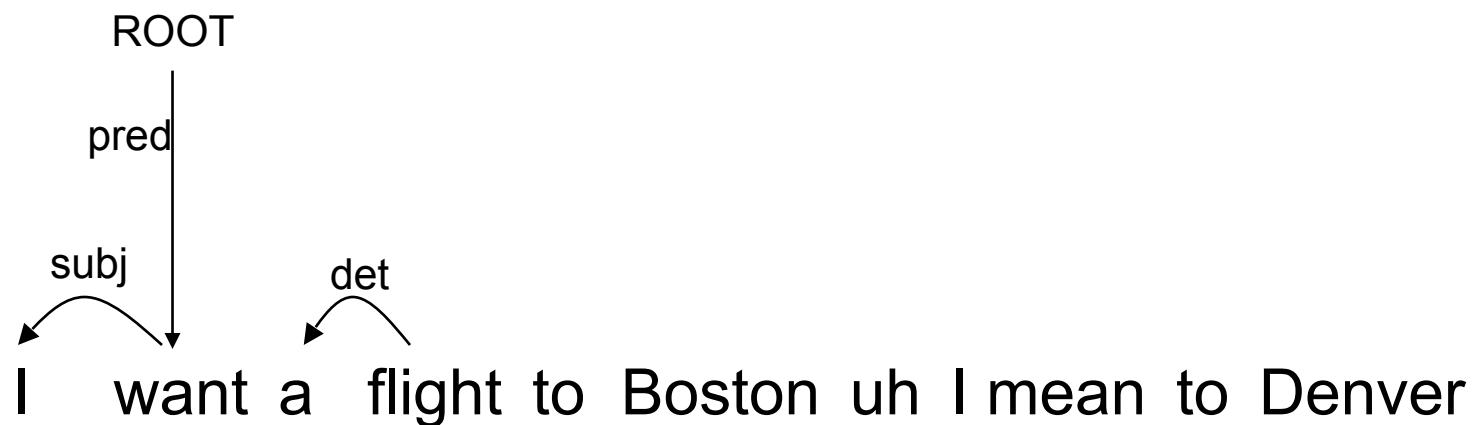
action: Shift()



# On an example

[ROOT, want, a]<sub>s</sub> [flight, to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

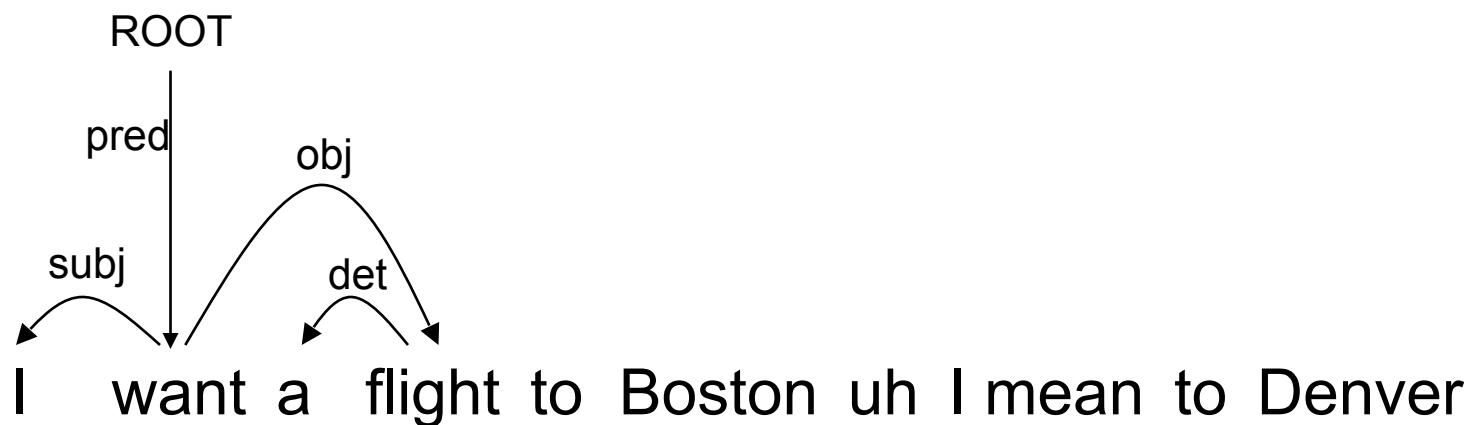
action: Left-Arc(DET)



# On an example

[ROOT, want, flight]<sub>S</sub> [to, Boston, uh, I, mean, to, Denver]<sub>Q</sub>

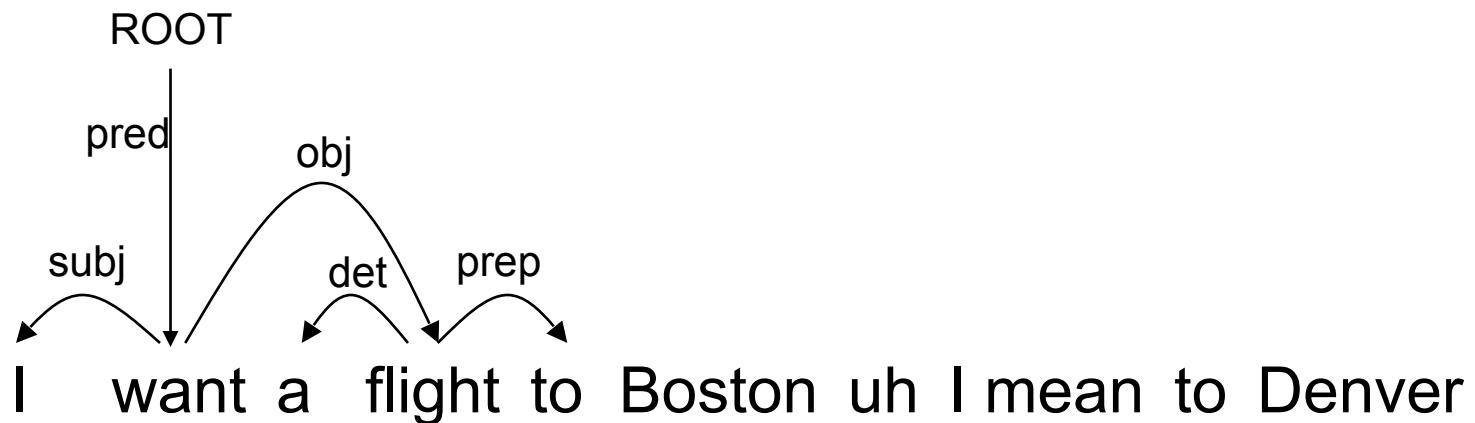
action: Right-Arc(OBJ)



# On an example

[ROOT, want, flight, to]<sub>s</sub> [Boston, uh, I, mean, to, Denver]<sub>Q</sub>

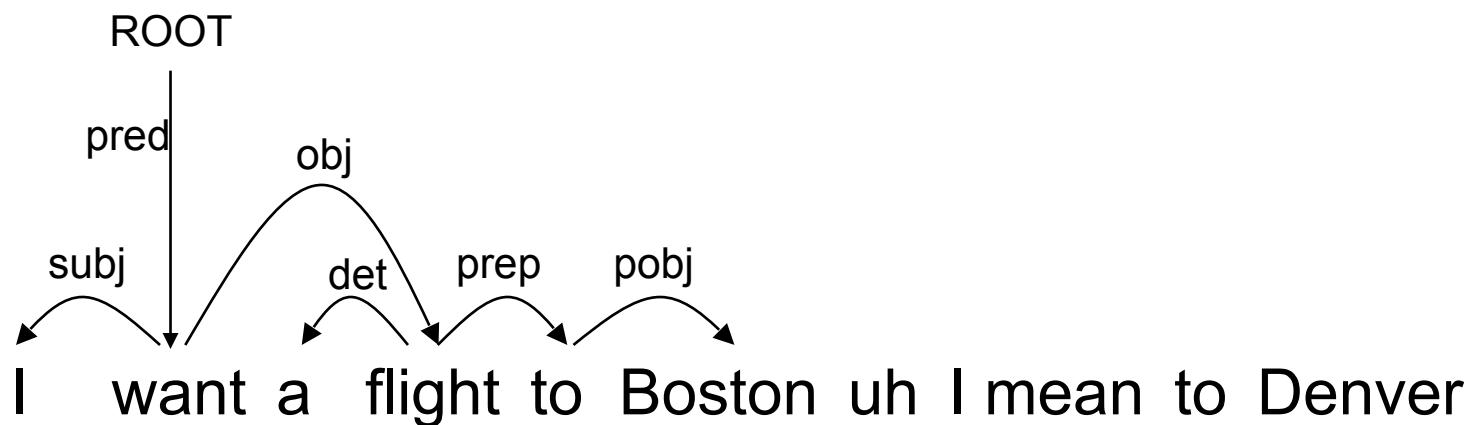
action: Right-Arc(PREP)



# On an example

[ROOT, want, flight, to, Boston]<sub>S</sub> [uh, I, mean, to, Denver]<sub>Q</sub>

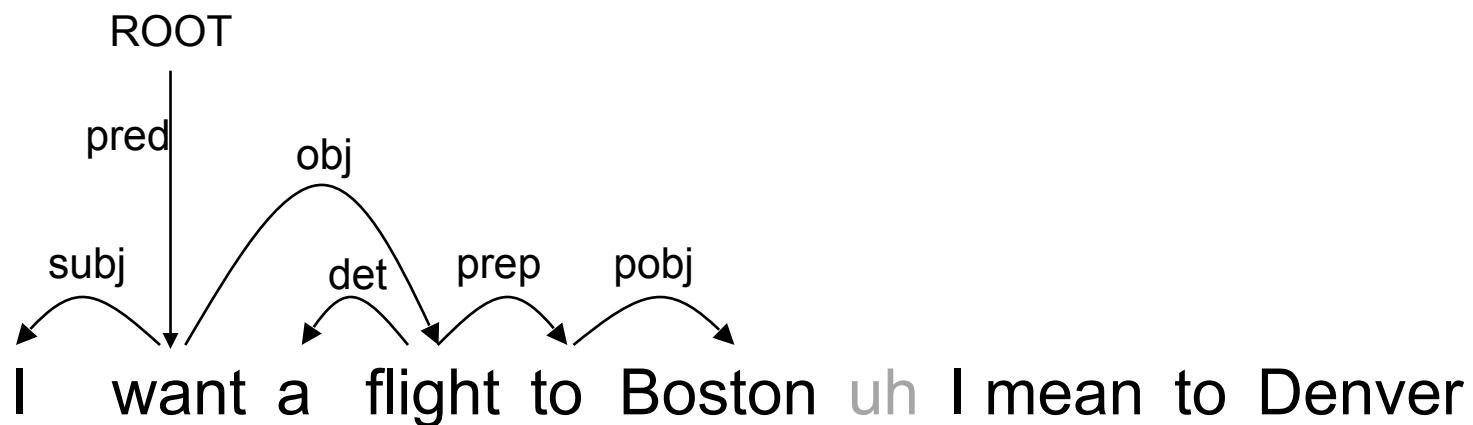
action: Right-Arc(POBJ)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [uh, I, mean, to, Denver]<sub>Q</sub>

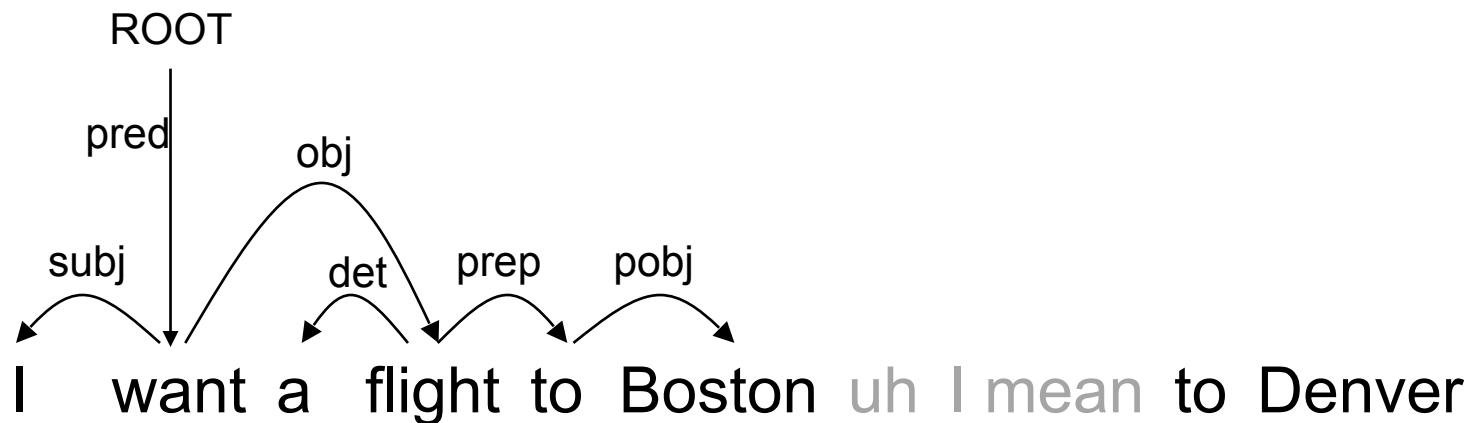
action: IJ



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [I, mean, to, Denver]<sub>Q</sub>

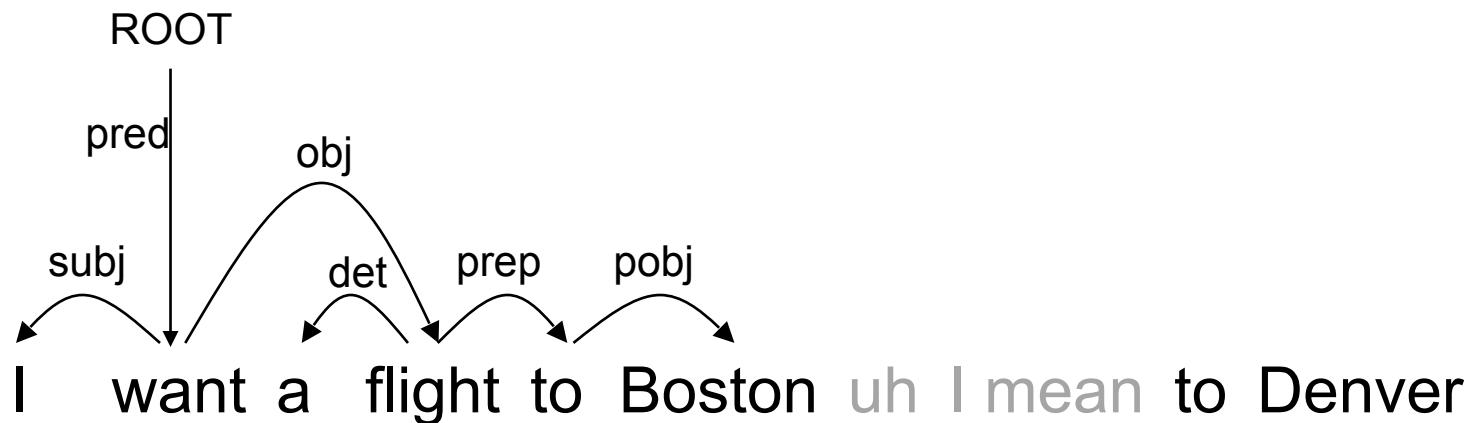
action: DM(2)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [to, Denver]<sub>Q</sub>

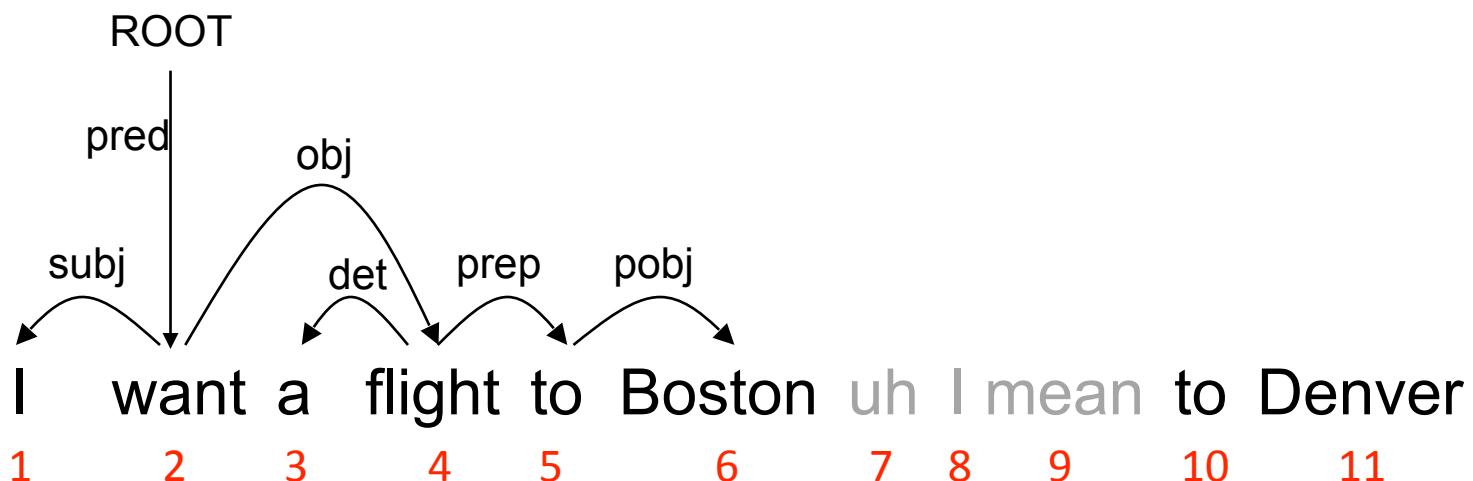
action: RP(5..6)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [to, Denver]<sub>Q</sub>

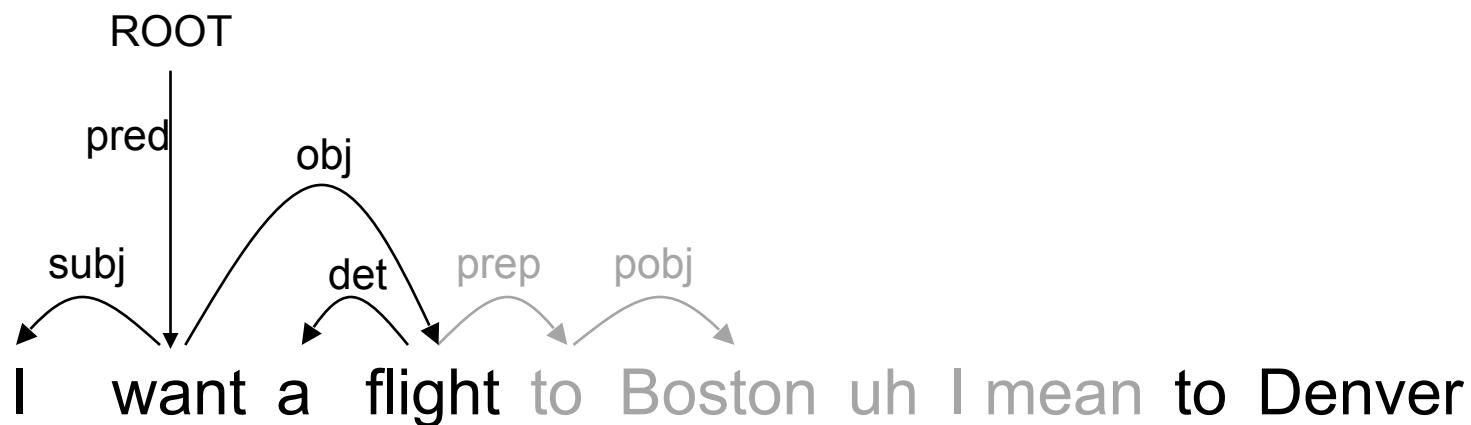
action: RP(5..6)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [to, Denver]<sub>Q</sub>

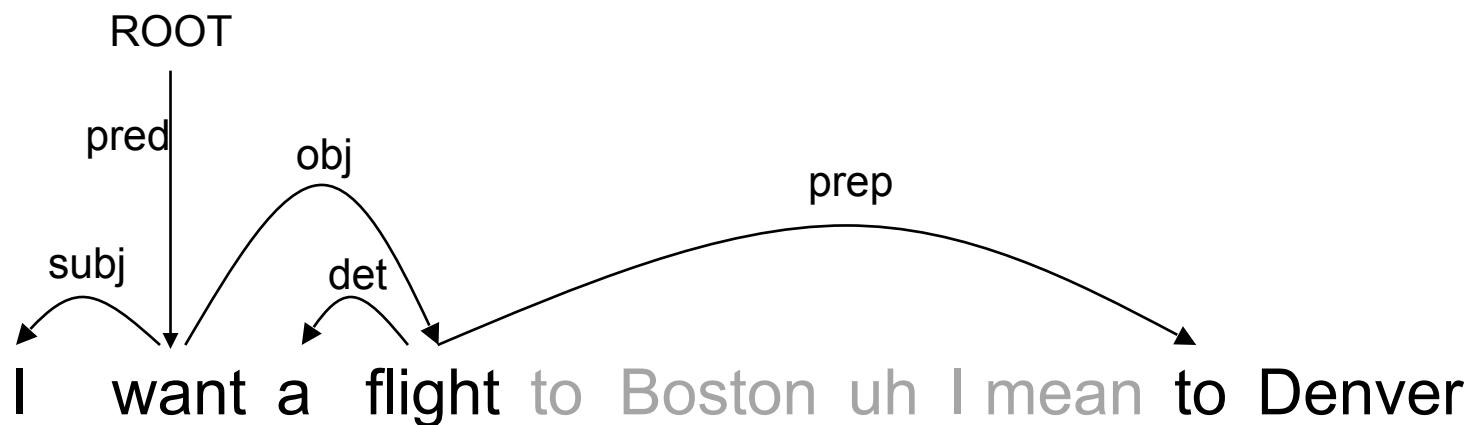
action: RP(5..6)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [to, Denver]<sub>Q</sub>

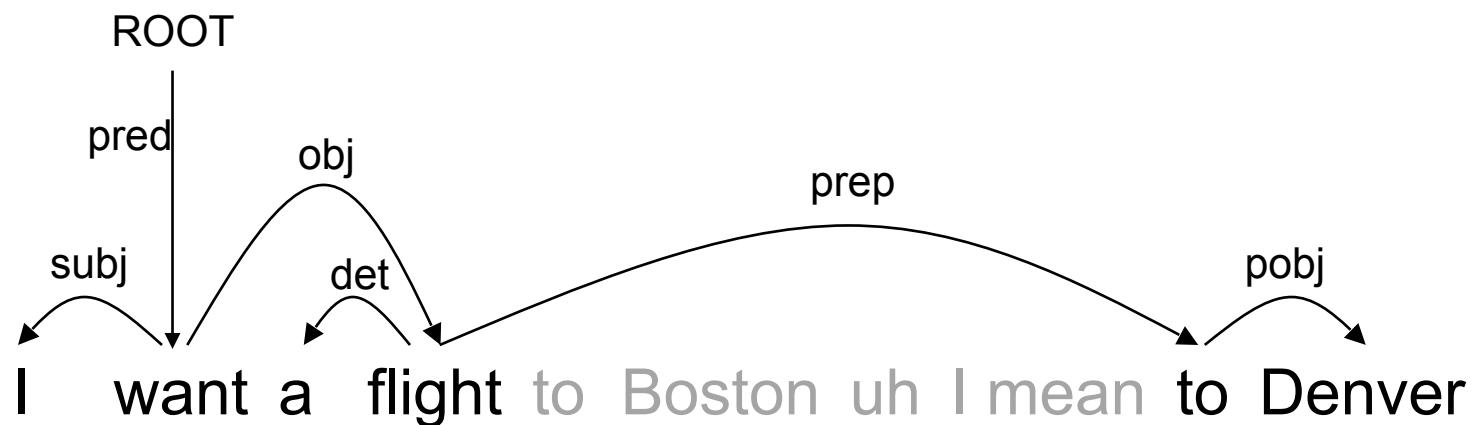
action: Right-Arc(PREP)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> [Denver]<sub>Q</sub>

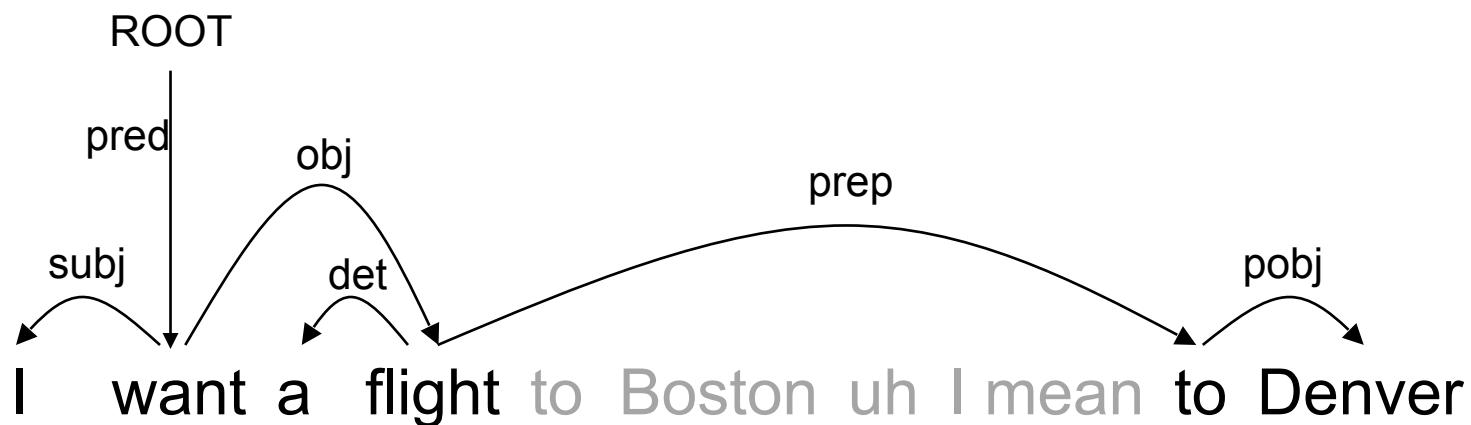
action: Right-Arc(POBJ)



# On an example

[ROOT, want, flight, to, Boston]<sub>s</sub> []<sub>Q</sub>

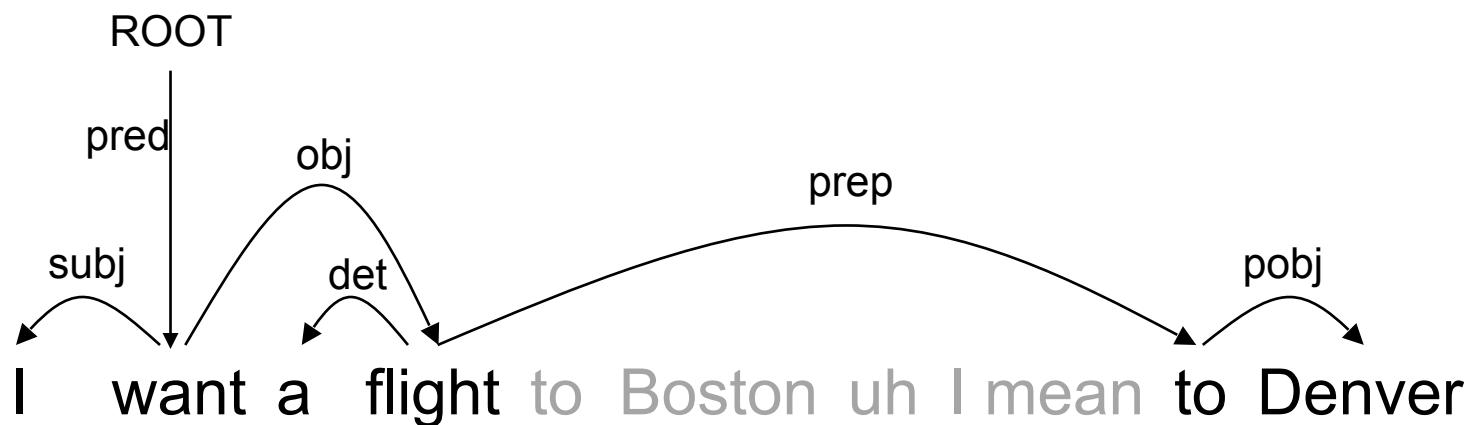
**action: Reduce**



# On an example

[ROOT, want, flight, to]s []<sub>Q</sub>

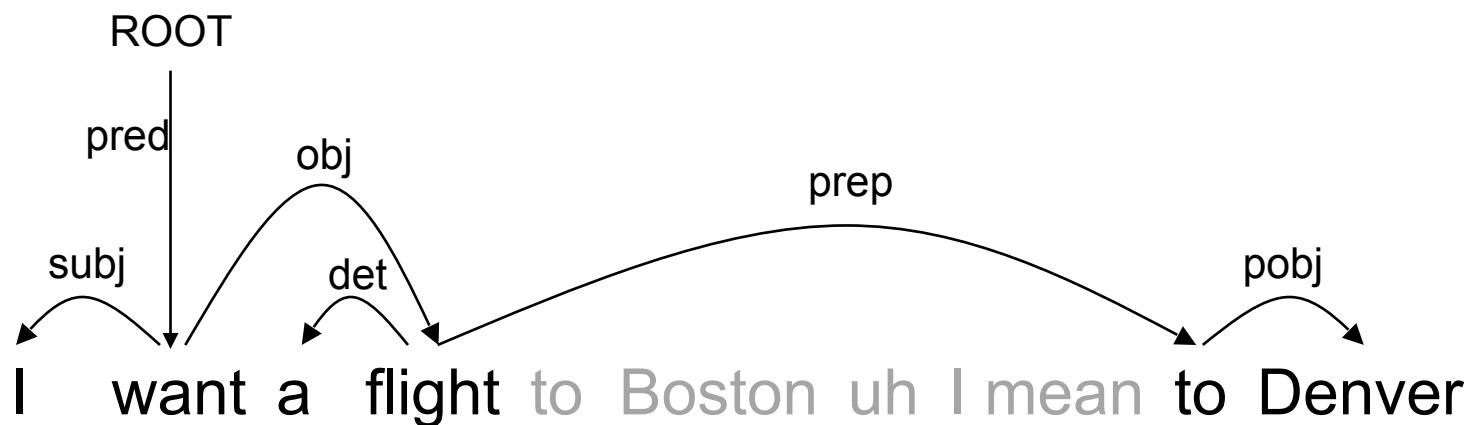
action: Reduce



# On an example

[ROOT, want, flight]<sub>s</sub> []<sub>Q</sub>

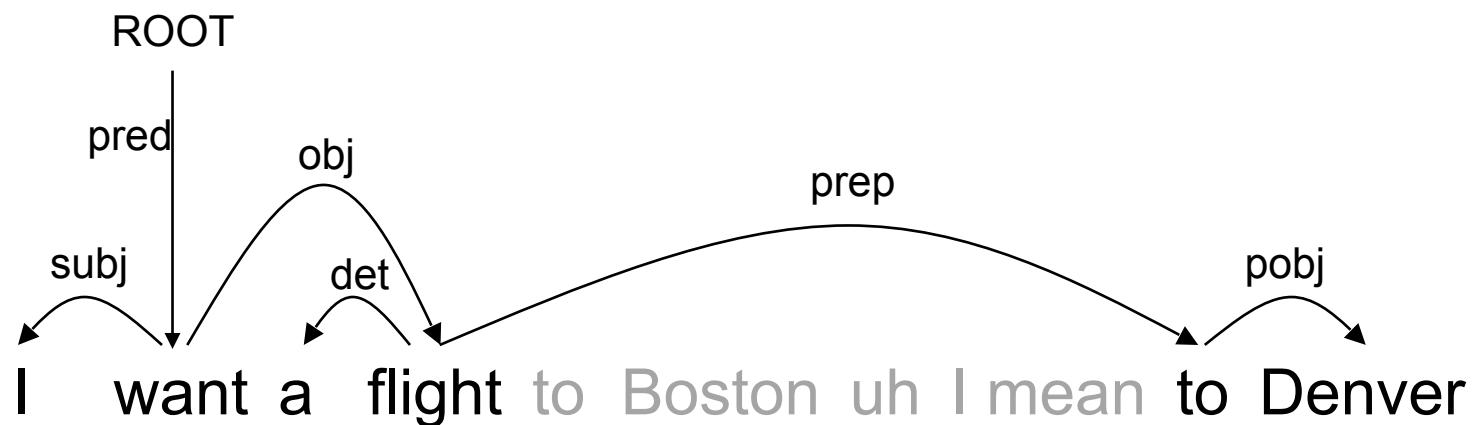
**action: Reduce**



# On an example

[ROOT, want]<sub>s</sub> []<sub>Q</sub>

**action: Reduce**



# Evaluation

Model	Description	F-score
Miller and Schuler (2008)	Joint + PCFG Parsing	30.6
Lease and Johnson (2006)	Joint + PCFG Parsing	62.4
Kahn et al. (2005)	TAG + LM rerank	78.2
Qian and Lui, (2013)	IOB tagging	82.5
<b>Rasooli and Tetrault (2013)</b>	<b>Arc-Eager Parsing</b>	<b>81.4</b>
Rasooli and Tetrault (2014)	Arc-Eager Parsing	82.6

# Neural Transition-Based Parsing



# Example of a neural arc-standard algorithm

- Chen and Manning (2014)
- Replace the feature-based action selection module by a neural network

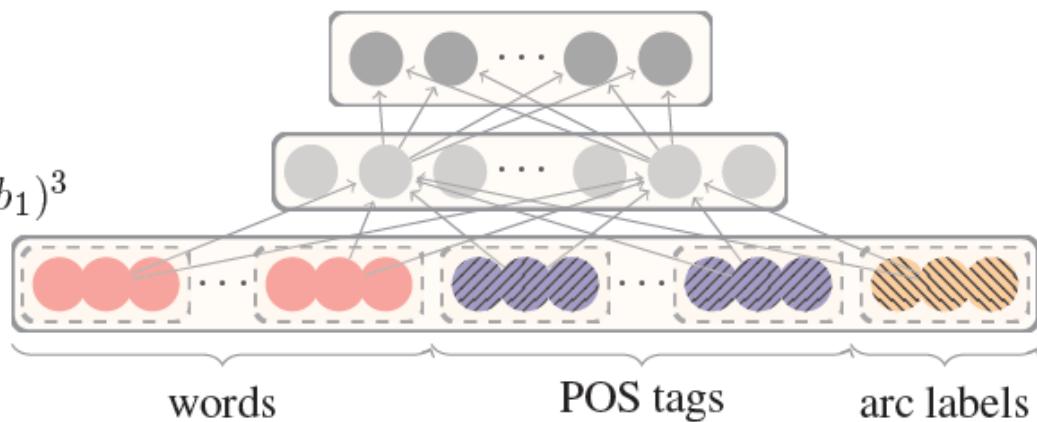
**Softmax layer:**

$$p = \text{softmax}(W_2 h)$$

**Hidden layer:**

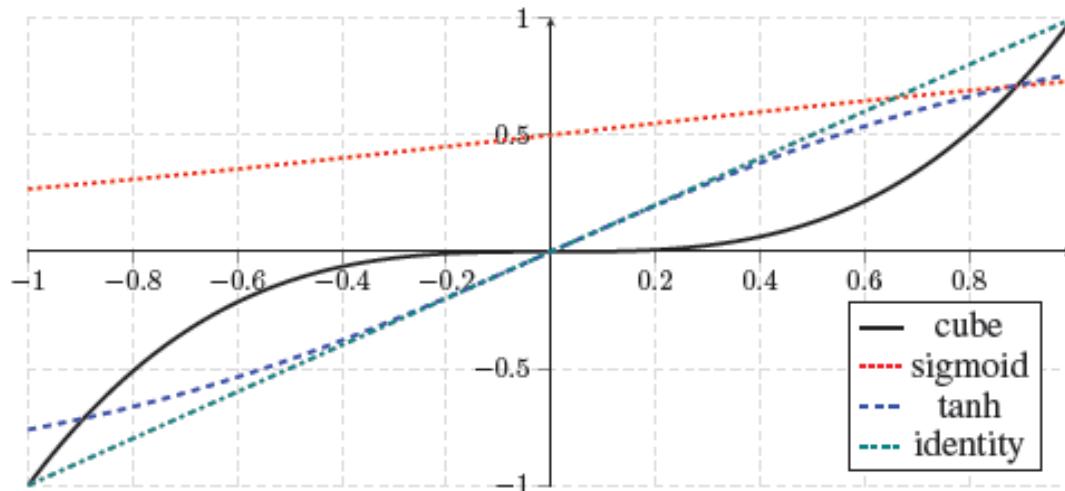
$$h = (W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1)^3$$

**Input layer:**  $[x^w, x^t, x^l]$



# Example of a neural arc-standard algorithm

- Chen and Manning (2014)
- Replace the feature-based action selection module by a neural network
- Cube activation function



$$h = (W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1)^3$$

# Example of a neural arc-standard algorithm

- Chen and Manning (2014)
- Replace the feature-based action selection module by a neural network
- Cube activation function
- POS tags and arc labels are discrete sets
  - Normally represented as one-hot vectors
  - Just like words, there should be similarities
    - NN (singular noun) should be similar to NNP (plural noun)
  - Dense embedding layer for POS tags and arc labels should capture relationships

# Example of a neural arc-standard algorithm: Experimental Results

Parser	Dev		Test		Speed (sent/s)
	UAS	LAS	UAS	LAS	
standard	89.9	88.7	89.7	88.3	51
eager	90.3	89.2	89.9	88.6	63
Malt:sp	90.0	88.8	89.9	88.5	560
Malt:eager	90.1	88.9	90.1	88.7	535
MSTParser	92.1	90.8	<b>92.0</b>	90.5	12
Our parser	<b>92.2</b>	<b>91.0</b>	<b>92.0</b>	<b>90.7</b>	<b>1013</b>

Table 4: Accuracy and parsing speed on PTB + CoNLL dependencies.

Parser	Dev		Test		Speed (sent/s)
	UAS	LAS	UAS	LAS	
standard	90.2	87.8	89.4	87.3	26
eager	89.8	87.4	89.6	87.4	34
Malt:sp	89.8	87.2	89.3	86.9	469
Malt:eager	89.6	86.9	89.4	86.8	448
MSTParser	91.4	88.1	90.7	87.6	10
Our parser	<b>92.0</b>	<b>89.7</b>	<b>91.8</b>	<b>89.6</b>	<b>654</b>

Table 5: Accuracy and parsing speed on PTB + Stanford dependencies.

Parser	Dev		Test		Speed (sent/s)
	UAS	LAS	UAS	LAS	
standard	82.4	80.9	82.7	81.2	72
eager	81.1	79.7	80.3	78.7	80
Malt:sp	82.4	80.5	82.4	80.6	420
Malt:eager	81.2	79.3	80.2	78.4	393
MSTParser	<b>84.0</b>	82.1	83.0	81.2	6
Our parser	<b>84.0</b>	<b>82.4</b>	<b>83.9</b>	<b>82.4</b>	<b>936</b>

Table 6: Accuracy and parsing speed on CTB.

# Example of a neural arc-standard algorithm: Model comparison

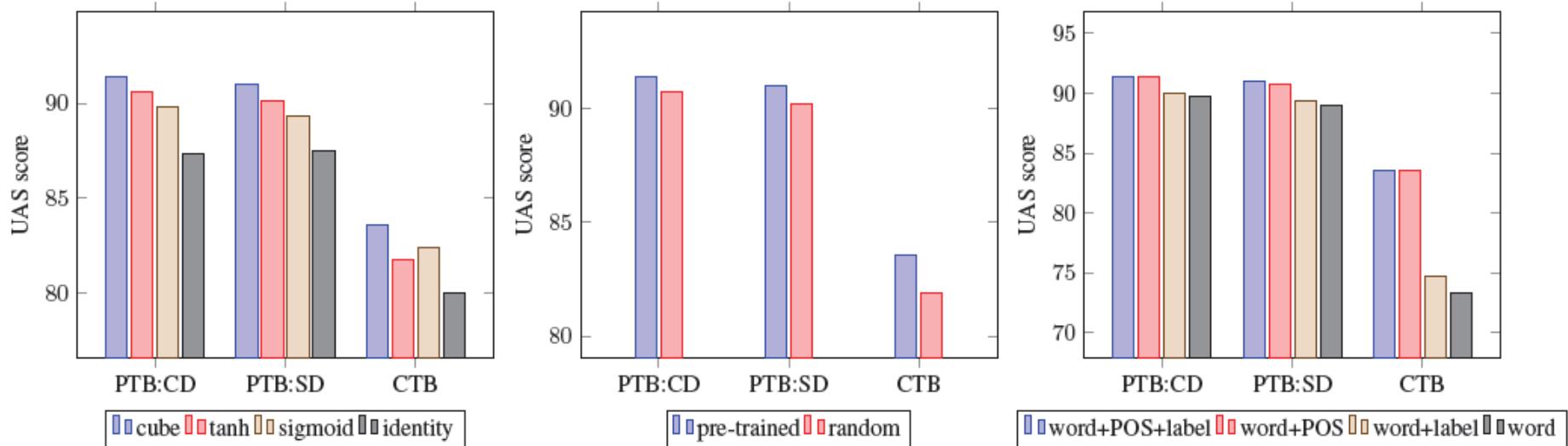


Figure 4: Effects of different parser components. Left: comparison of different activation functions. Middle: comparison of pre-trained word vectors and random initialization. Right: effects of POS and label embeddings.

# Practical assignment 5

(will count 4 times less than other assignments)



# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation

FRIED RICE		Fri
Tom Moy Special .....	7.00	Ca
Pork.....	3.50	Ca
→ Chicken .....	3.70	M
Shrimp.....	5.00	Ch
BBQ Pork .....	6.75	Eg
Beef.....	6.75	1/
→ Real Chicken.....	6.75	D
Meatless.....	3.70	1/

# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation



Trønder-Avisa  
@trondneravisa

Suivre

▼

OL-leiren bestilte 1500 egg gjennom å oversette via Google Translate. Men det slo feil. 15.000 ble levert på døra. Vi ønsker lykke til og håper at de norske gullhåpene er glade – veldig glade – i egg: 😊



# Goal

- Understand the current state of the art in machine translation



Trønder-Avisa  
@tronderavisa

Suivre

OL-leiren bestilte 1500 egg gjennom å oversette via Google Translate. Men det slo feil. 15.000 ble levert på døra. Vi ønsker lykke til og håper at de norske gullhåpene er glade – veldig glade – i egg: 😊



09:30 - 3 févr. 2018

09:30 - 3 févr. 2018

# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation



Trønder-Avisa  
@trondneravisa

Suivre

OL-leiren bestilte 1500 egg gjennom å oversette via Google Translate. Men det slo feil. 15.000 ble levert på døra. Vi ønsker lykke til og håper at de norske gullhåpene er glade – veldig glade – i egg: 😊

**Google Translate translation:**  
The OL camp ordered 1500 eggs through translating through Google Translate. But it was wrong. 15,000 was delivered to the door. We wish good luck and hope that the Norwegian golden hopes are happy - very happy - in eggs

# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation



**Google Translate translation:**  
I'm not in the office right now.  
Send any work to translate it.

# Goal

- Understand the current state of the art in machine translation



# Goal

- Understand the current state of the art in machine translation
- Use whichever (freely available) machine translation system to investigate the following questions:
  - Which syntactic constructions are correctly translated? Which ones are not? Can you guess why?
  - What about questions, long and/or complex sentences, sentences in the first person, idioms (*kick the bucket* vs. *casser sa pipe*), named entities (*London*)...
  - What is the influence of the language pair? Cf. English<->French vs. English<->Chinese or English<->(a rare language)
  - What is the influence of the context. E.g., how are the following discourse translated (English <-> French):
    - *A mouse appeared. It looked hungry.*
    - *Il aime bien le mouton. Surtout les côtelettes de mouton.*

# Assignment deliverable

- A 1- to 2-page report discussing these questions, with examples

[https://github.com/edupoux/MVA\\_2018\\_SL/tree/master/TD\\_%235](https://github.com/edupoux/MVA_2018_SL/tree/master/TD_%235)

A detailed reproduction of Pieter Bruegel the Elder's painting "The Tower of Babel". The scene depicts a massive, multi-tiered tower under construction, rising from a rocky base. The tower is built of light-colored stone and features numerous arched windows and doorways. In the foreground, a group of people in period clothing, including a man with a long white beard, stand on a rocky shore. The background shows a vast landscape with distant hills and a cloudy sky.

**That's all for today!**