

“Training Deterministic Parsers with Non-Deterministic Oracles”

by Yoav Goldberg and Joakim Nivre, 2013

Seminarvortrag

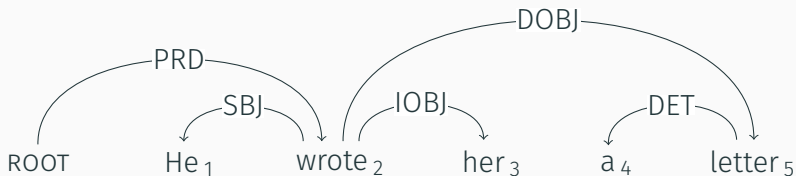
Pius Meinert

July 13, 2018

Training **Deterministic Parsers** with Non-Deterministic Oracles

He₁ wrote₂ her₃ a₄ letter₅

Training **Deterministic Parsers** with Non-Deterministic Oracles



Definition (Transition System)

A transition system for dependency parsing is a quadruple $S = (C, T, c_s, C_t)$, where

1. C is a set (configurations),
2. T is a set of transitions, each of which is a (partial) function $t : C \rightarrow C$,
3. c_s is an initialization function, mapping sentence $w = w_1w_2\dots w_n$ to a configuration $c \in C$,
4. $C_t \subseteq C$ (terminal configurations).

Training **Deterministic Parsers** with Non-Deterministic Oracles

He₁ wrote₂ her₃ a₄ letter₅

Training **Deterministic Parsers** with Non-Deterministic Oracles

ROOT He₁ wrote₂ her₃ a₄ letter₅

$\top_{C_S(w)}$

[ROOT], [He₁, wrote₂, her₃, a₄, letter₅], {}

Training **Deterministic Parsers** with Non-Deterministic Oracles

ROOT He₁ wrote₂ her₃ a₄ letter₅

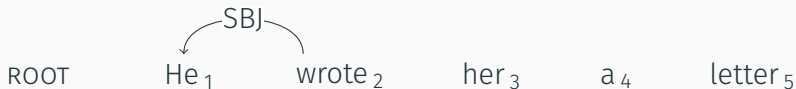
[ROOT], [He₁, wrote₂, her₃, a₄, letter₅]

┌
| SHIFT

[ROOT, He₁], [wrote₂, her₃, a₄, letter₅]

Training **Deterministic Parsers** with Non-Deterministic Oracles

ROOT He₁ wrote₂ her₃ a₄ letter₅



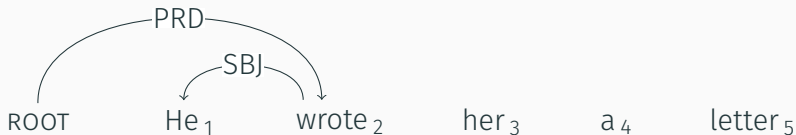
```
graph TD
    ROOT --- He1[He1]
    ROOT --- wrote2[wrote2]
    wrote2 --- her3[her3]
    wrote2 --- a4[a4]
    wrote2 --- letter5[letter5]
    He1 -- SBJ --> wrote2
```

[ROOT, He₁], [wrote₂, her₃, a₄, letter₅]

┌ LEFT_{SBJ}

[ROOT], [wrote₂, her₃, a₄, letter₅]

Training **Deterministic Parsers** with Non-Deterministic Oracles

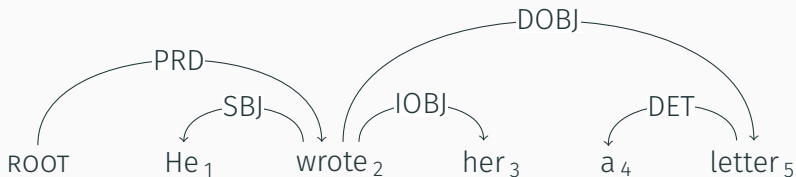


[ROOT], [wrote₂, her₃, a₄, letter₅]

┌
RIGHT_{PRD}

[ROOT, wrote₂], [her₃, a₄, letter₅]

Training **Deterministic Parsers** with Non-Deterministic Oracles



[ROOT, wrote₂], [her₃, a₄, letter₅]

┌ RIGHT_{IOBJ}, SHIFT, LEFT_{DET}, REDUCE, RIGHT_{DOBJ}

[ROOT, wrote₂, letter₅], [] ∈ C_t

Training Deterministic Parsers with Non-Deterministic Oracles

```
1 if  $c = (\sigma|i,j|\beta, A)$  and  $(j, i) \in T$  then
2    $t \leftarrow$  LEFT
3 else if  $c = (\sigma|i,j|\beta, A)$  and  $(i, j) \in T$  then
4    $t \leftarrow$  RIGHT
5 else if  $c = (\sigma|i,j|\beta, A)$  and  $\exists k[k < i \wedge [(k, j) \in T \vee (j, k) \in T]]$ 
   then
6    $t \leftarrow$  REDUCE
7 else
8    $t \leftarrow$  SHIFT
9 return  $t$ 
```

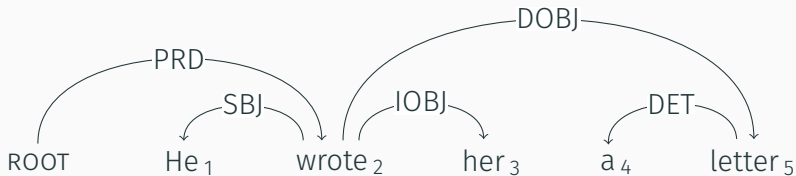
Greedy Classifier-based Parsing

```
1  
2    $c \leftarrow C_S(w)$   
3   while  $c \notin C_t$  do  
4      $t_p \leftarrow \arg \max_{t \in \text{LEGAL}(c)} \mathbf{w} \cdot \phi(c, t)$   
5  
6  
7  
8  
9  
10  
11      $c \leftarrow t_p(c)$   
12 return  $A_c$ 
```

Training Deterministic Parsers with Non-Deterministic Oracles

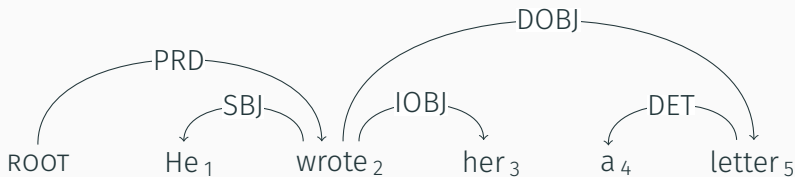
```
1 for  $(w, T) \in d$  do
2    $c \leftarrow C_S(w)$ 
3   while  $c \notin C_t$  do
4      $t_p \leftarrow \arg \max_{t \in \text{LEGAL}(c)} \mathbf{w} \cdot \phi(c, t)$ 
5      $\text{CORRECT}(c) \leftarrow \{t \mid o(t; c, T) = \text{true}\}$ 
6      $t_o \leftarrow \arg \max_{t \in \text{CORRECT}(c)} \mathbf{w} \cdot \phi(c, t)$ 
7     if  $t_p \notin \text{CORRECT}(c)$  then
8        $\text{UPDATE}(\mathbf{w}, \phi(c, t_o), \phi(c, t_p))$ 
9        $c \leftarrow t_o(c)$ 
10    else
11       $c \leftarrow t_p(c)$ 
12 return  $\mathbf{w}$ 
```

Training Deterministic Parsers with **Non-Deterministic** Oracles



SH, LA_{SBJ}, RA_{PRD}, RA_{IOBJ}, SH, LA_{DET}, RE, RA_{DOBJ}

Training Deterministic Parsers with **Non-Deterministic** Oracles

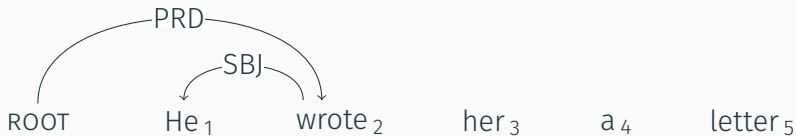


SH, LA_{SBJ}, RA_{PRD}, RA_{IOBJ}, SH, LA_{DET}, RE, RA_{DOBJ}

SH, LA_{SBJ}, RA_{PRD}, RA_{IOBJ}, RE, SH, LA_{DET}, RA_{DOBJ}

→ spurious ambiguity requires non-deterministic oracle instead of static oracle

... with Non-Deterministic and Complete Oracles

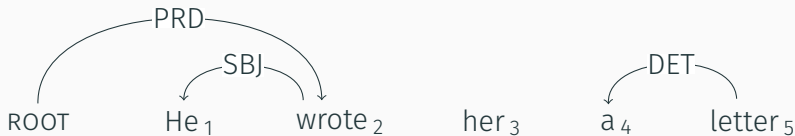


[ROOT], [He₁, wrote₂, her₃, a₄, letter₅]

┌ SH, LA_{SBJ}, RA_{PRD}, SH

[ROOT, wrote₂, her₃], [a₄, letter₅]

... with Non-Deterministic and Complete Oracles



[ROOT, wrote₂, her₃], [a₄, letter₅]

⌈ SH, LA_{DET}, SH

[ROOT, wrote₂, her₃, letter₅], [] ∈ C_t

→ error propagation can be mitigated by complete oracle

→ **dynamic oracle**: non-deterministic + complete

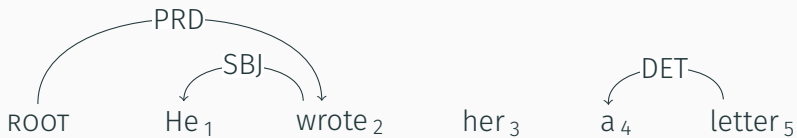
Training (Standard)

```
1 for  $(w, T) \in d$  do
2    $c \leftarrow C_S(w)$ 
3   while  $c \notin C_t$  do
4      $t_p \leftarrow \arg \max_{t \in \text{LEGAL}(c)} \mathbf{w} \cdot \phi(c, t)$ 
5      $\text{CORRECT}(c) \leftarrow \{t \mid o(t; c, T) = \text{true}\}$ 
6      $t_o \leftarrow \arg \max_{t \in \text{CORRECT}(c)} \mathbf{w} \cdot \phi(c, t)$ 
7     if  $t_p \notin \text{CORRECT}(c)$  then
8        $\text{UPDATE}(\mathbf{w}, \phi(c, t_o), \phi(c, t_p))$ 
9        $c \leftarrow t_o(c)$ 
10    else
11       $c \leftarrow t_p(c)$ 
12 return  $w$ 
```

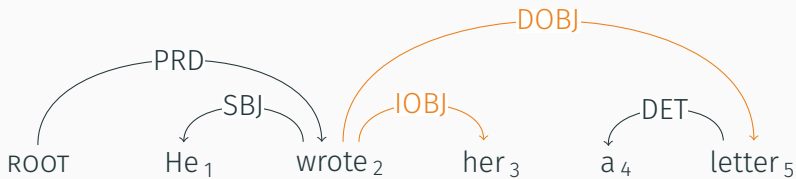
Training with Exploration

```
1 for  $(w, T) \in d$  do
2    $c \leftarrow C_S(w)$ 
3   while  $c \notin C_t$  do
4      $t_p \leftarrow \arg \max_{t \in \text{LEGAL}(c)} w \cdot \phi(c, t)$ 
5      $\text{OPTIMAL}(c) \leftarrow \{t \mid o(t; c, T) = \text{true}\}$ 
6      $t_o \leftarrow \arg \max_{t \in \text{OPTIMAL}(c)} w \cdot \phi(c, t)$ 
7     if  $t_p \notin \text{OPTIMAL}(c)$  then
8       UPDATE( $w, \phi(c, t_o), \phi(c, t_p)$ )
9        $c \leftarrow \text{EXPLORE}(c, t_o, t_p)$ 
10    else
11       $c \leftarrow t_p(c)$ 
12 return  $w$ 
```

Optimality / Transition Costs

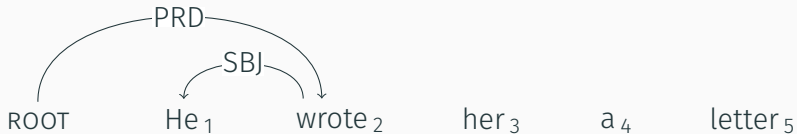


Optimality / Transition Costs



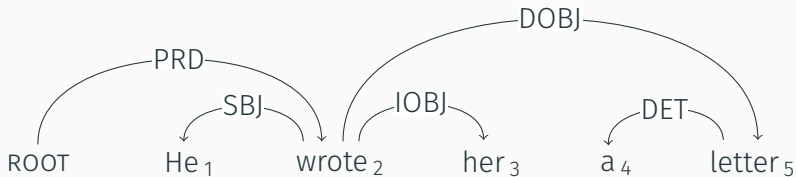
$$C(A, T) = 2$$

Optimality / Transition Costs



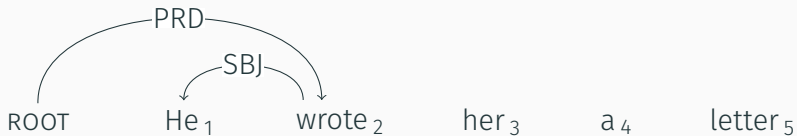
[ROOT, wrote₂, her₃], [a₄, letter₅]

Optimality / Transition Costs



$$\min_{A: C \rightsquigarrow A} C(A, T) = 0$$

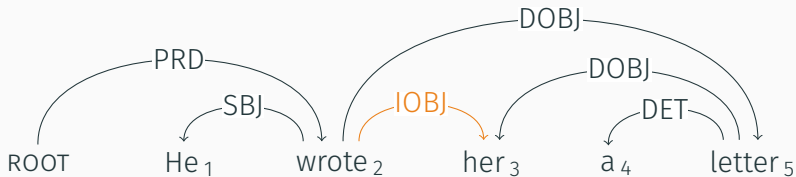
Optimality / Transition Costs



[ROOT, wrote₂, her₃], [a₄, letter₅]

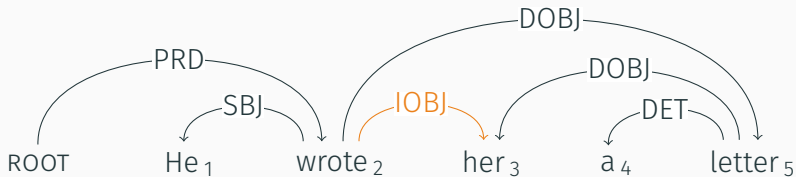
⌈ SH, ...

Optimality / Transition Costs



$$C(\text{SHIFT}; c, T) = \min_{A:t(c)\rightsquigarrow A} C(A, T) - \min_{A:c\rightsquigarrow A} C(A, T) = 1$$

Optimality / Transition Costs



$$\mathcal{C}(\text{SHIFT}; c, T) = \min_{A:t(c) \rightsquigarrow A} \mathcal{C}(A, T) - \min_{A:c \rightsquigarrow A} \mathcal{C}(A, T) = 1$$

$$o_d(c, T) = \{t \mid \mathcal{C}(t; c, T) = 0\}$$

Arc Decomposition - Definition

Definition (Tree Consistency)

A set of arcs A is said to be *tree consistent* if there exists a projective dependency tree T such that $A \subseteq T$.

Definition (Arc Decomposition)

A transition system is said to be *arc decomposable* if, for every tree consistent arc set A and configuration c , $c \rightsquigarrow A$ is entailed by $c \rightsquigarrow (h, d)$ for every arc $(h, d) \in A$.

Arc Decomposition - Arc-Standard Counterexample

$$c = ([a, b, c], \beta)$$



Arc-Standard Transitions

$$\text{LEFT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_0, \beta, A \cup \{(s_0, s_1)\})$$

$$\text{RIGHT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_1, \beta, A \cup \{(s_1, s_0)\})$$

$$\text{SHIFT}[(\sigma, b|\beta, A)] = (\sigma|b, \beta, A)$$

Arc Decomposition - Arc-Standard Counterexample

$$c = ([a, b, c], \beta) \stackrel{\text{LEFT}}{\vdash} ([a, c], \beta)$$



Arc-Standard Transitions

$$\text{LEFT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_0, \beta, A \cup \{(s_0, s_1)\})$$

$$\text{RIGHT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_1, \beta, A \cup \{(s_1, s_0)\})$$

$$\text{SHIFT}[(\sigma, b|\beta, A)] = (\sigma|b, \beta, A)$$

Arc Decomposition - Arc-Standard Counterexample

$$c = ([a, b, c], \beta) \stackrel{\text{RIGHT}}{\vdash} ([a, b], \beta) \stackrel{\text{LEFT}}{\vdash} ([b], \beta)$$



Arc-Standard Transitions

$$\text{LEFT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_0, \beta, A \cup \{(s_0, s_1)\})$$

$$\text{RIGHT}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_1, \beta, A \cup \{(s_1, s_0)\})$$

$$\text{SHIFT}[(\sigma, b|\beta, A)] = (\sigma|b, \beta, A)$$

Arc Decomposition - Arc-Eager Proof Sketch

Given: arbitrary configuration $c = (\sigma, \beta, A)$ and tree consistent arc set A' such that all arc are reachable from c .

To show: $c \rightsquigarrow A'$

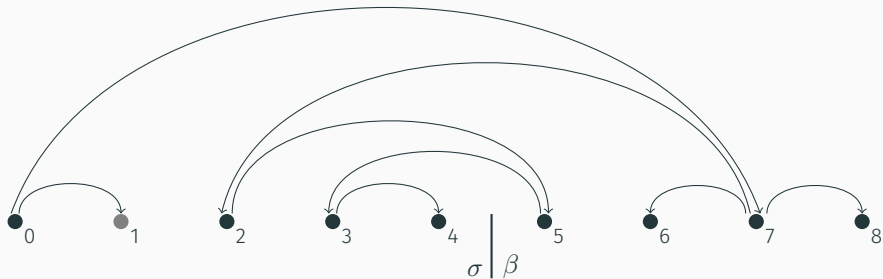
$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

Arc Decomposition - Arc-Eager Proof Sketch



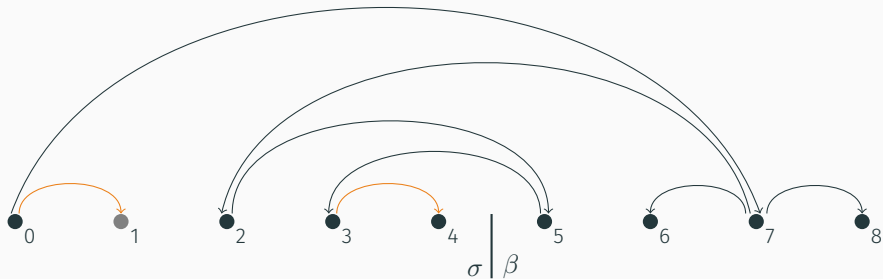
$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

Arc Decomposition - Arc-Eager Proof Sketch



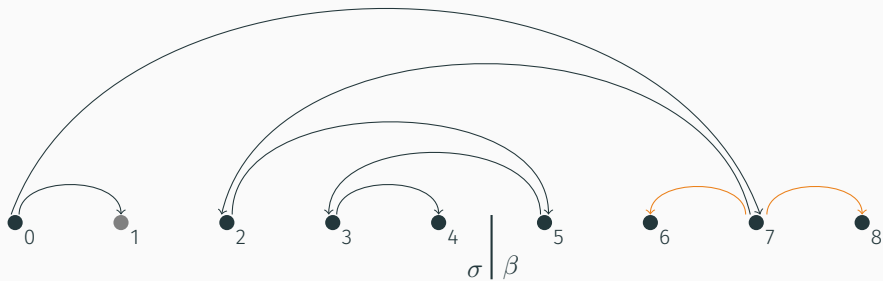
$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

Arc Decomposition - Arc-Eager Proof Sketch



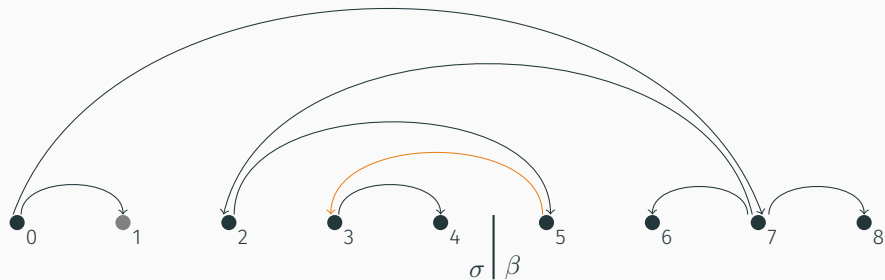
$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

Arc Decomposition - Arc-Eager Proof Sketch



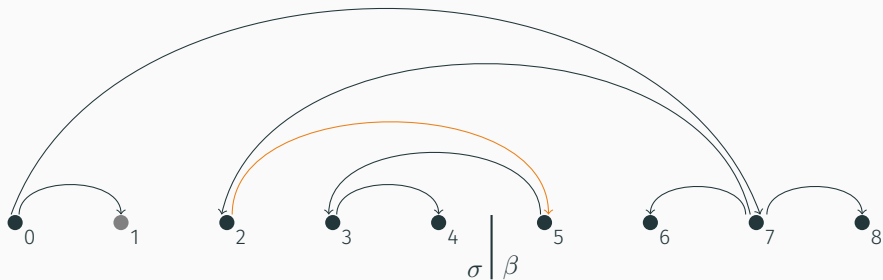
$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

Arc Decomposition - Arc-Eager Proof Sketch



$$\bar{\mathcal{B}} = \{(h, d) \mid h, d \notin \beta\}$$

$$\mathcal{B} = \{(h, d) \mid h, d \in \beta\}$$

$$\mathcal{B}_h = \{(h, d) \mid h \in \beta, d \in \sigma\}$$

$$\mathcal{B}_d = \{(h, d) \mid d \in \beta, h \in \sigma\}$$

$$o_d(c, T) = \{t \mid \mathcal{C}(t; c, T) = 0\}$$

$$\mathcal{C}(t; c, T) = \min_{A:t(c)\rightsquigarrow A} \mathcal{C}(A, T) - \min_{A:c\rightsquigarrow A} \mathcal{C}(A, T)$$

Efficiently compute transition costs:

1. Intersect set of individually reachable arcs with goal arc set.
2. Gain set of individually reachable goal arcs and thusly, reachable goal arc set.
3. See how a given transition affects this set of reachable arcs.

✓ Arc-Eager

- Nivre 2003
- Goldberg and Nivre 2012

✗ Arc-Standard

- Nivre 2004
- Goldberg, Sartorio, and Satta 2014

✓ Hybrid

- Kuhlmann, Gómez-Rodríguez, and Satta 2011

✓ Easy-First

- Goldberg and Elhadad 2010

Conclusion

- spurious ambiguity
static → non-deterministic oracle
- error propagation
incomplete → complete oracle
- **dynamic oracle** (non-deterministic + complete)
arc decomposability
- good runtime
optimization during training
- experiments show improved accuracy



Yoav Goldberg and Michael Elhadad. “An Efficient Algorithm for Easy-First Non-Directional Dependency Parsing”. In: *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 2-4, 2010, Los Angeles, California, USA*. 2010, pp. 742–750. URL: <http://www.aclweb.org/anthology/N10-1115>.



Yoav Goldberg and Joakim Nivre. “A Dynamic Oracle for Arc-Eager Dependency Parsing”. In: *COLING 2012, 24th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, 8-15 December 2012, Mumbai, India*. 2012, pp. 959–976. URL: <http://aclweb.org/anthology/C/C12/C12-1059.pdf>.



Yoav Goldberg and Joakim Nivre. “Training Deterministic Parsers with Non-Deterministic Oracles”. In: *TACL 1 (2013)*, pp. 403–414. URL: <https://tacl2013.cs.columbia.edu/ojs/index.php/tacl/article/view/145>.



Yoav Goldberg, Francesco Sartorio, and Giorgio Satta.
“A Tabular Method for Dynamic Oracles in
Transition-Based Parsing”. In: *TACL 2* (2014),
pp. 119–130. URL:
[https://tacl2013.cs.columbia.edu/ojs/
index.php/tacl/article/view/302](https://tacl2013.cs.columbia.edu/ojs/index.php/tacl/article/view/302).



Marco Kuhlmann, Carlos Gómez-Rodríguez, and Giorgio Satta. “Dynamic Programming Algorithms for Transition-Based Dependency Parsers”. In: *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA*. 2011, pp. 673–682. URL: <http://www.aclweb.org/anthology/P11-1068>.



Joakim Nivre. “An Efficient Algorithm for Projective Dependency Parsing”. In: *Proceedings of the 8th International Workshop on Parsing Technologies (IWPT)*. 2003, pp. 149–160.

References v



Joakim Nivre. “Incrementality in Deterministic Dependency Parsing”. In: *Proceedings of the Workshop on Incremental Parsing: Bringing Engineering and Cognition Together*. 2004. URL: <http://www.aclweb.org/anthology/W04-0308>.



Joakim Nivre. “Algorithms for Deterministic Incremental Dependency Parsing”. In: *Computational Linguistics* 34.4 (2008), pp. 513–553. DOI: [10.1162/coli.07-056-R1-07-027](https://doi.org/10.1162/coli.07-056-R1-07-027). URL: <https://doi.org/10.1162/coli.07-056-R1-07-027>.

Transition Costs - Arc-Eager

$$\begin{aligned}\mathcal{C}(\text{SHIFT}; (\sigma, b | \beta, A), T) &= |\{(k, b) \in T \mid k \in \sigma\} \\ &\quad \cup \{(b, k) \in T \mid k \in \sigma \wedge \forall x \in V : (x, k) \notin A\}| \\ \mathcal{C}(\text{RIGHT}; (\sigma | s, b | \beta, A), T) &= |\{(k, b) \in T \mid k \in \sigma \cup \beta\} \\ &\quad \cup \{(b, k) \in T \mid k \in \sigma \wedge \forall x \in V : (x, k) \notin A\}| \\ \mathcal{C}(\text{LEFT}; (\sigma | s, b | \beta, A), T) &= |\{(k, s) \in T \mid k \in \beta\} \cup \{(s, k) \in T \mid k \in \beta\}| \\ \mathcal{C}(\text{REDUCE}; (\sigma | s, \beta, A), T) &= |\{(s, k) \in T \mid k \in \beta\}| \end{aligned}$$

Transition Costs - Hybrid

$$\mathcal{C}(\text{SHIFT}; (\sigma|s_1|s_0, b|\beta, A), T) = |\{(b, k) \in T \mid k \in \{s_0, s_1\} \cup \sigma\} \\ \cup \{(k, b) \in T \mid k \in \{s_1\} \cup \sigma\}|$$

$$\mathcal{C}(\text{RIGHT}; (\sigma|s_1|s_0, \beta, A), T) = |\{(s_0, k) \in T \mid k \in \beta\} \\ \cup \{(k, s_0) \in T \mid k \in \beta\}|$$

$$\mathcal{C}(\text{LEFT}; (\sigma|s_1|s_0, b|\beta, A), T) = |\{(s_0, k) \in T \mid k \in \{b\} \cup \beta\} \\ \cup \{(k, s_0) \in T \mid k \in \{s_1\} \cup \beta\}|$$

Transition Costs - Easy-First

$$\mathcal{C}(\text{TR}; (\lambda, A), T) = |\{(h', d) \in T \mid h' \in \lambda \wedge h' \neq h\} \cup \{(d, d') \in T \mid d' \in \lambda\}|$$

$\text{TR} \in \{\text{LEFT}_{lb}^i \mid 1 < i \leq |\lambda|\} \cup \{\text{RIGHT}_{lb}^i \mid 1 \leq i < |\lambda|\}$ and (h, d) added by TR