

Natural Language Processing

Yue Zhang
Westlake University



Chapter 11

Transition-based Methods for Structured Prediction

Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Contents

- **11.1 Transition-based Structured Prediction**
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

A Transition System

- A different view from graph-based models for structured prediction
 - Encodes rich unbounded features without complexity constraints
 - A general framework that is easy to adapt to different tasks
- Maps output building process into a state transducer
 - State — S_i
 - Corresponds to partial results during decoding
 - Action — a_i
 - The operations that can be applied for state transition
 - Construct output incrementally

A Transition System

- **Automata:**

- State — S_i

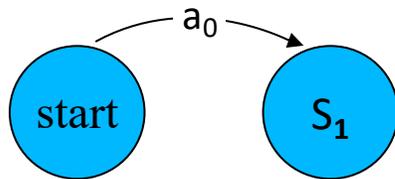
- Action — a_i



A Transition System

- **Automata:**

- State — S_i
- Action — a_i

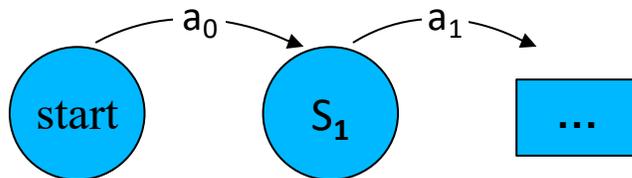


A Transition System

- **Automata:**

- State — S_i

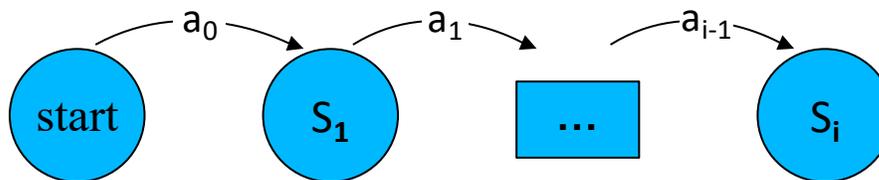
- Action — a_i



A Transition System

- **Automata:**

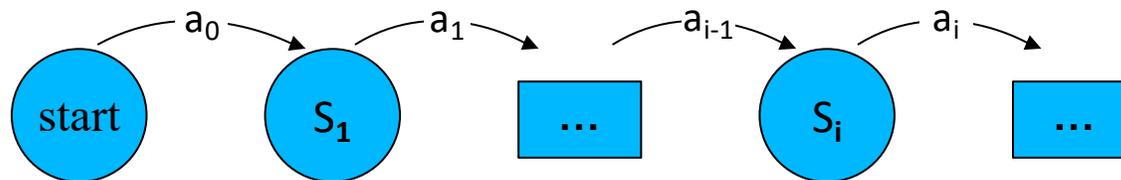
- State — S_i
- Action — a_i



A Transition System

- **Automata:**

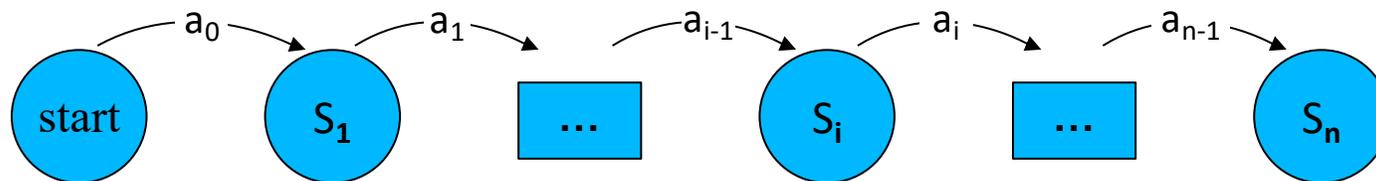
- State — S_i
- Action — a_i



A Transition System

- **Automata:**

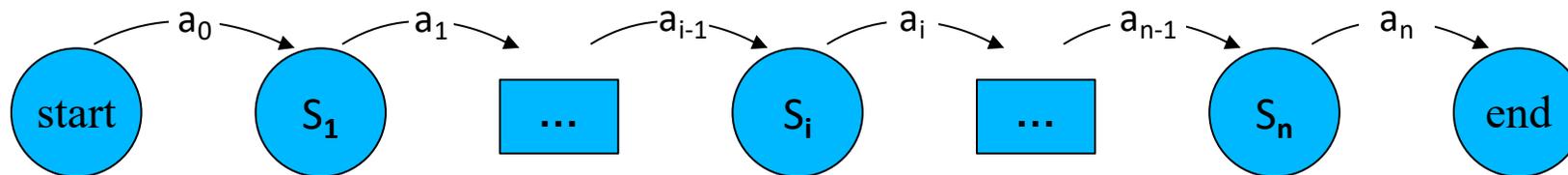
- State — S_i
- Action — a_i



A Transition System

- **Automata:**

- State — S_i
- Action — a_i



Transition-based structured prediction

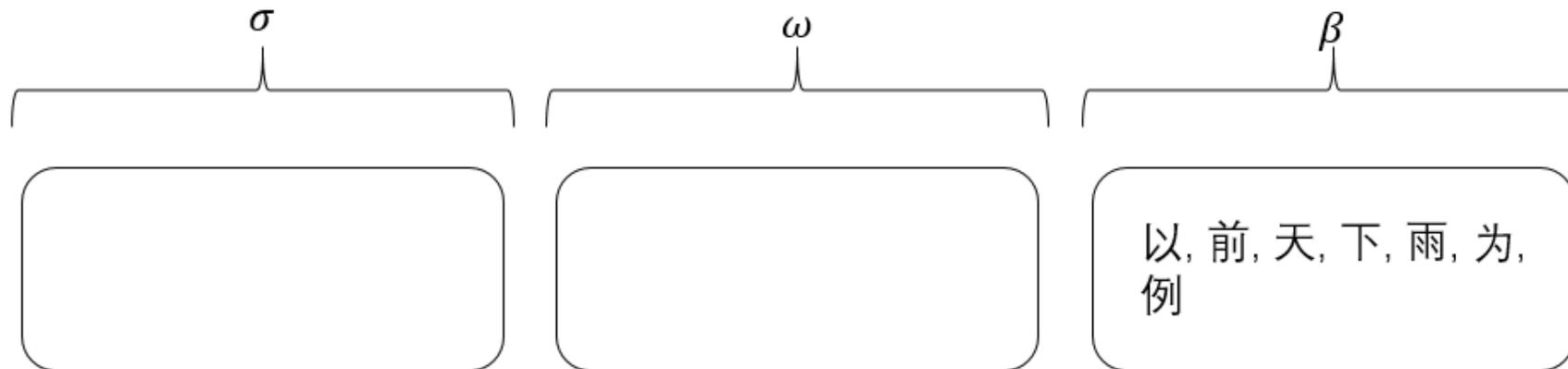
Transition-based structured prediction is a **state transition** process.

Automata:

- State
 - Start state — an empty structure
 - End state — the output structure
 - Intermediate state — partially constructed structures
- Transition actions
 - Incremental steps that build output structures, change one state to another

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

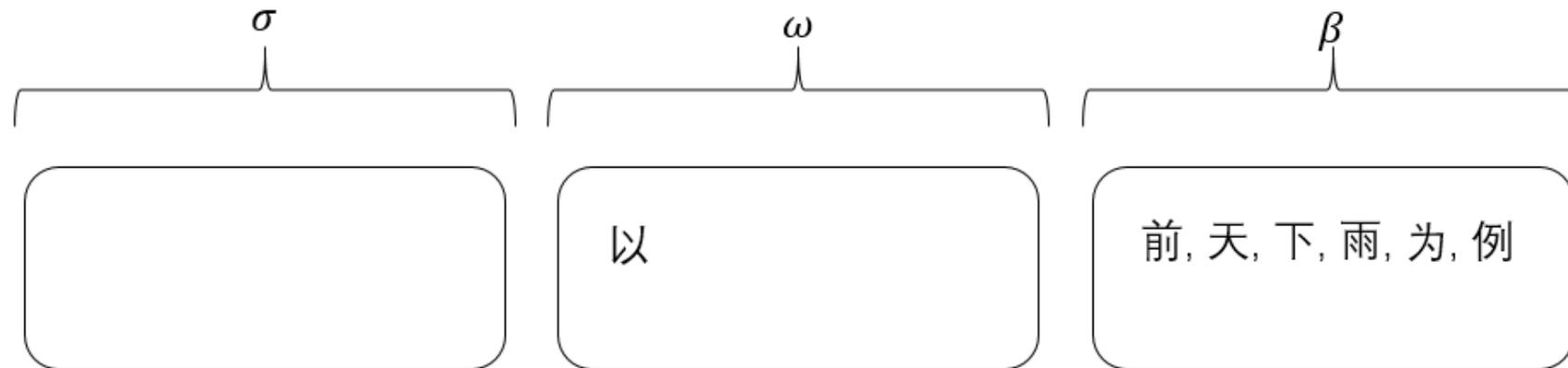


Next Action: *SEP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

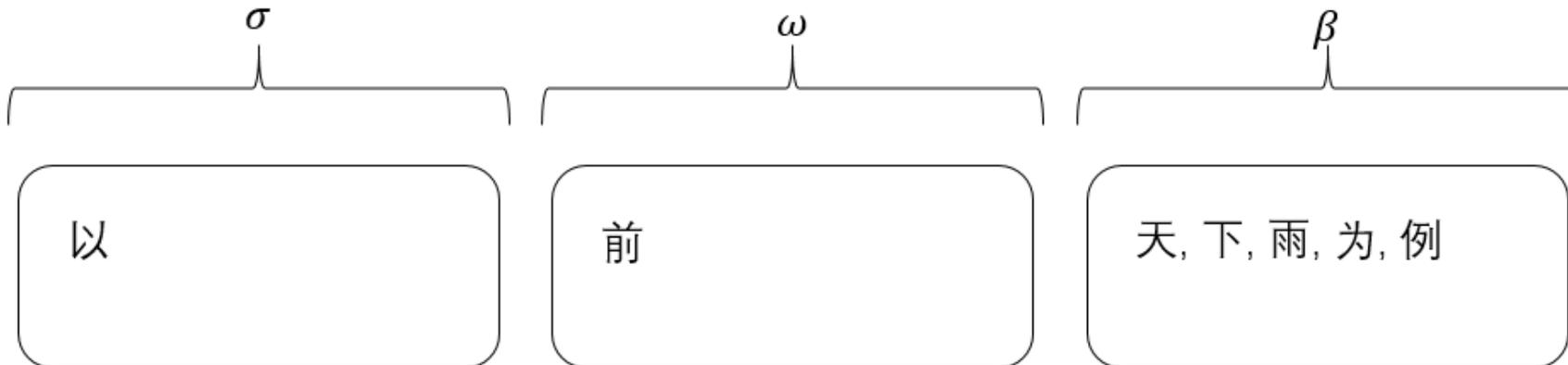


Next Action: *SEP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

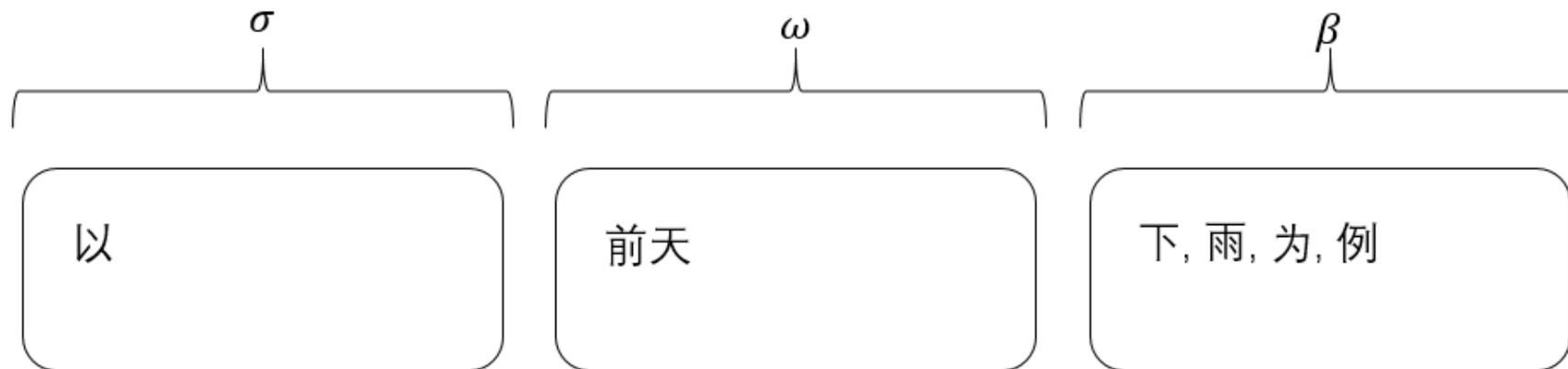


Next Action: *APP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

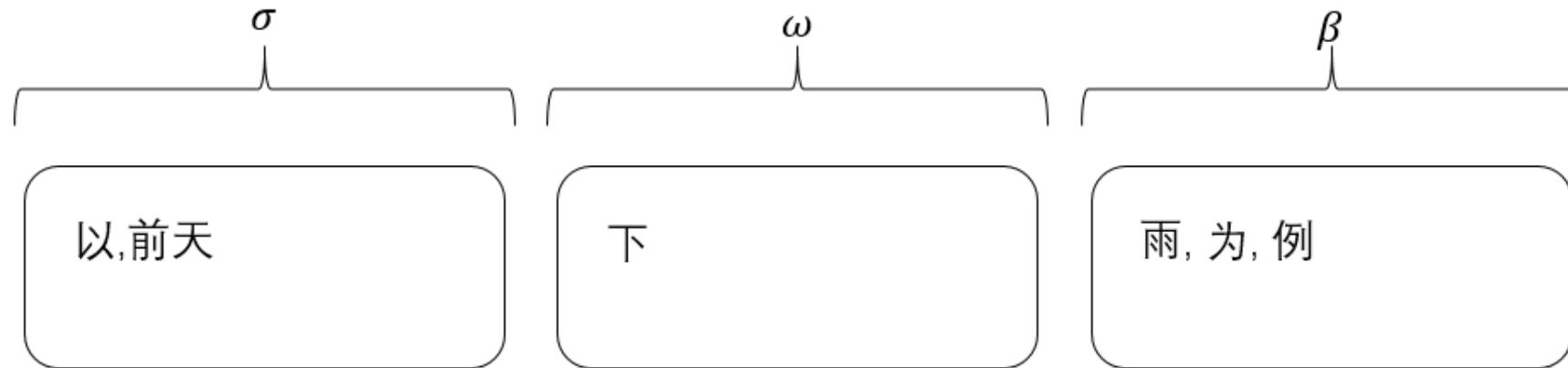


Next Action: *SEP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

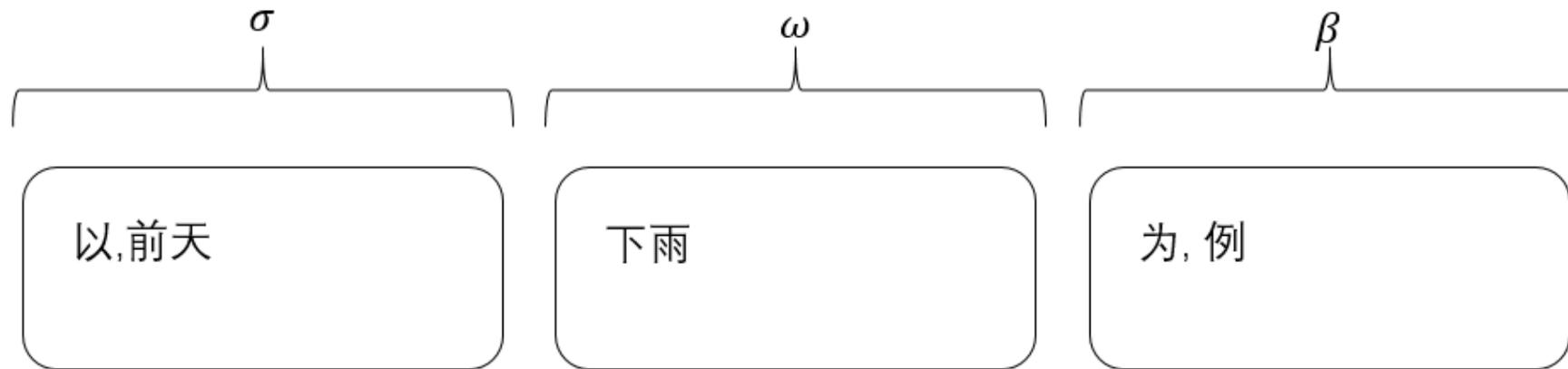


Next Action: *APP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

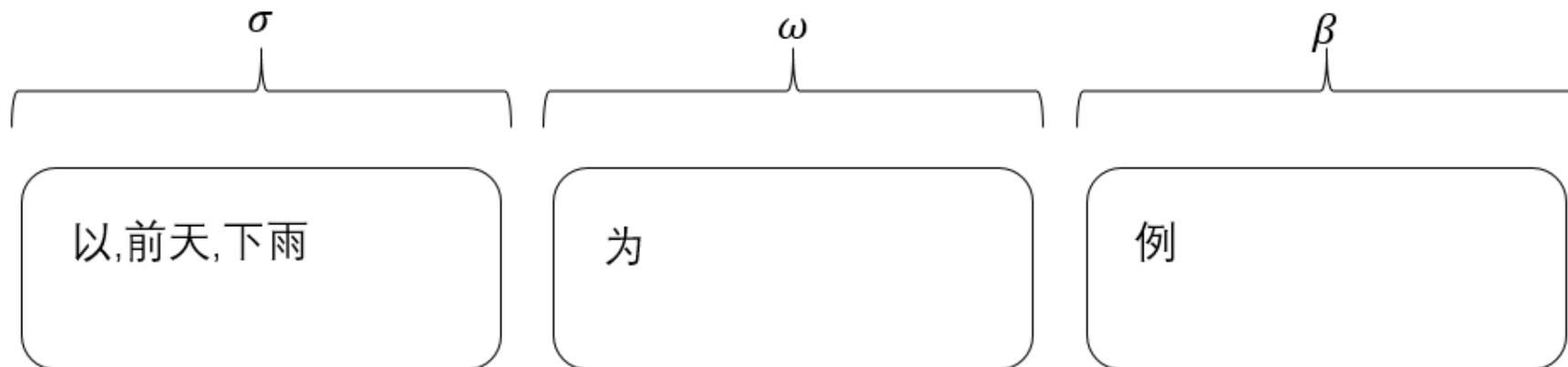


Next Action: *SEP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

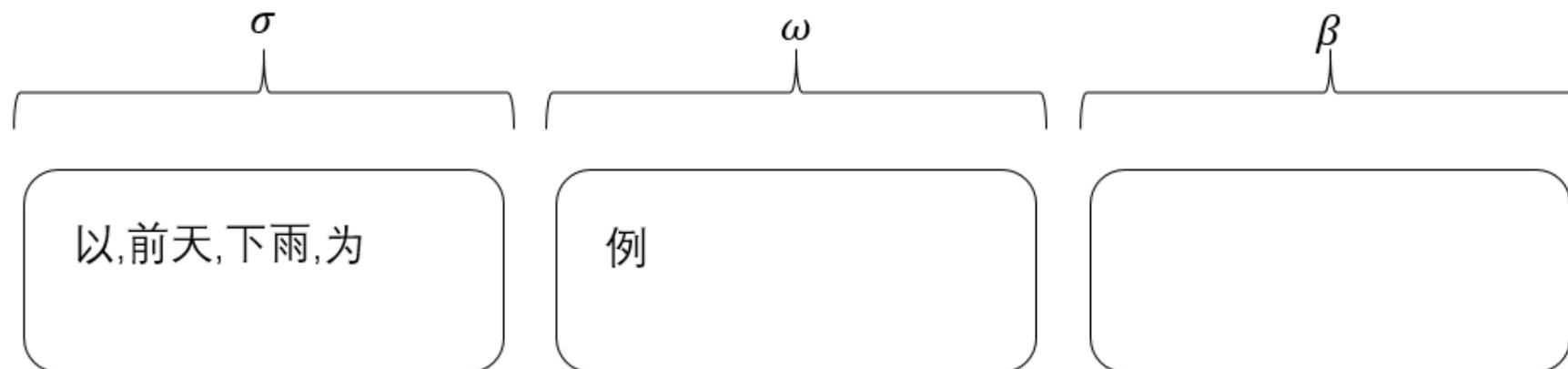


Next Action: *SEP*

σ : partial output w : current partial word β : list of next incoming characters

Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

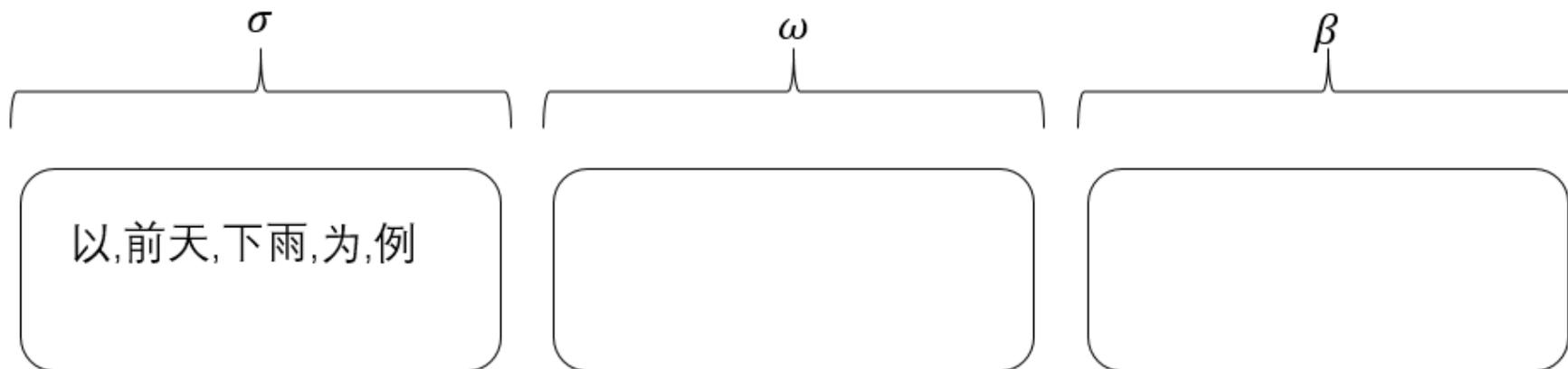


Next Action: *FIN*

σ : partial output w : current partial word β : list of next incoming characters

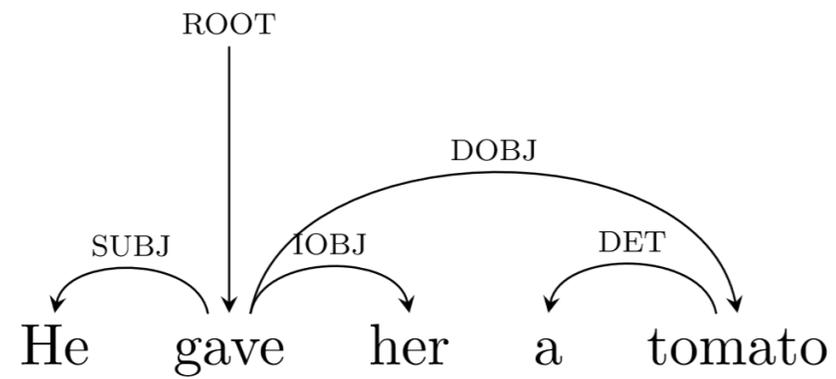
Example (Word Segmentation)

- Input: "以(take) 前(before) 天(day) 下(fall) 雨(rain) 为(be) 例(example)"

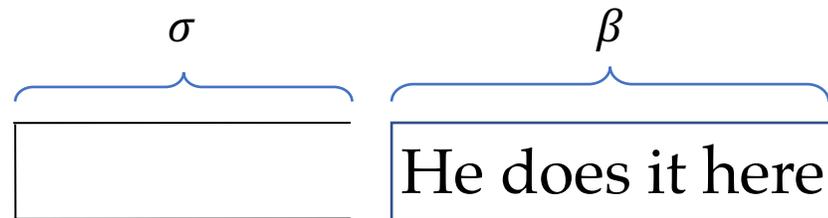


σ : partial output w : current partial word β : list of next incoming characters

Example (Dependency Parsing)



Example (Dependency Parsing)

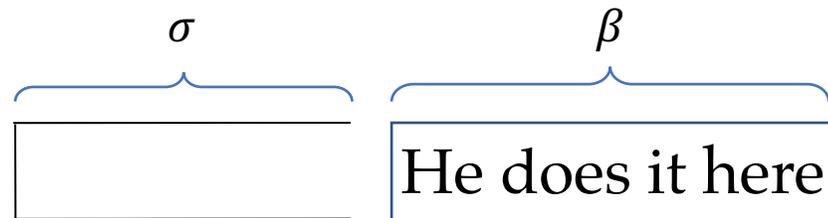


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Shift*

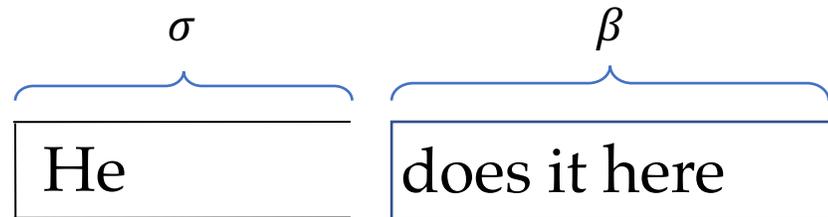


σ : stack

β : buffer

Example (Dependency Parsing)

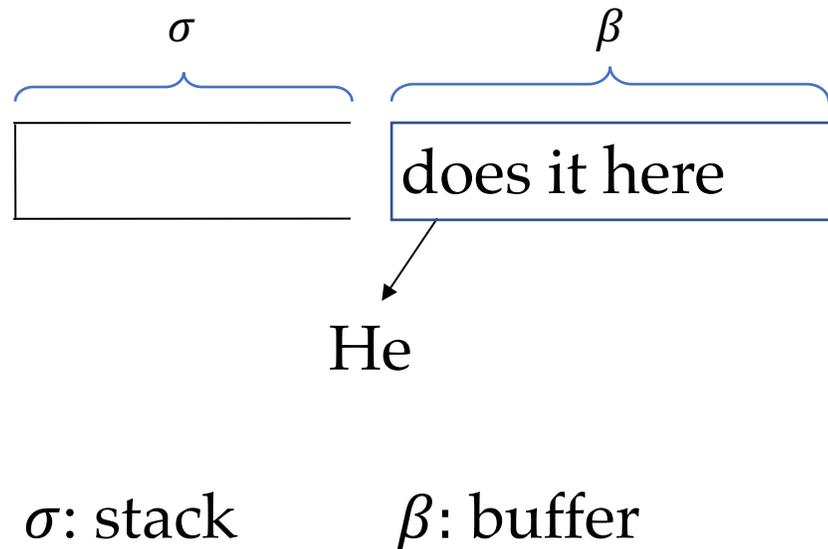
- NEXT ACTION: *Left – Arc*



σ : stack β : buffer

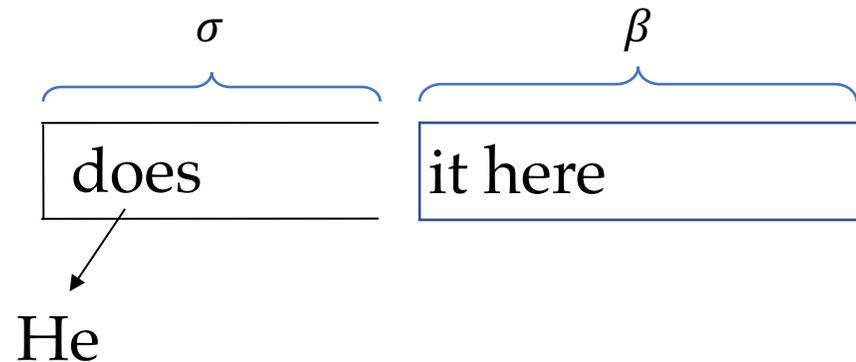
Example (Dependency Parsing)

- NEXT ACTION: *Shift*



Example (Dependency Parsing)

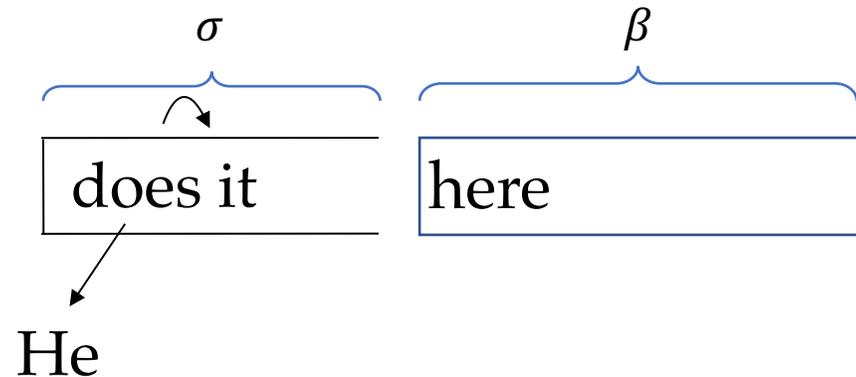
- NEXT ACTION: *Right – Arc*



σ : stack β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Reduce*

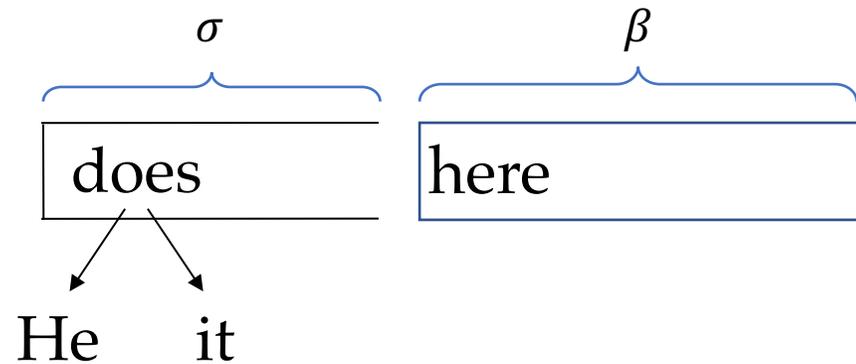


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Right – Arc*

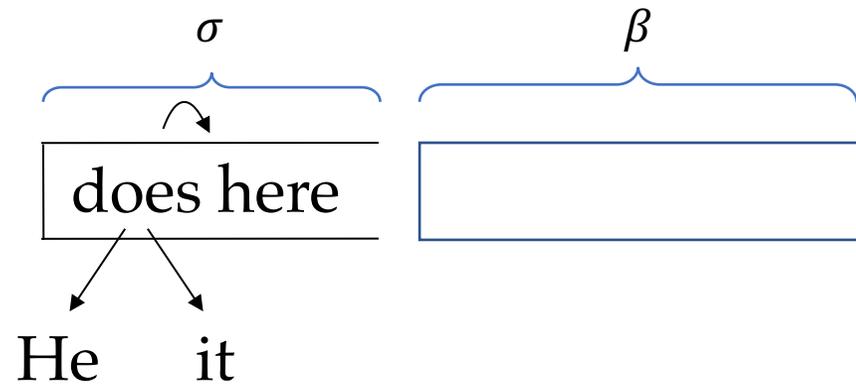


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Reduce*

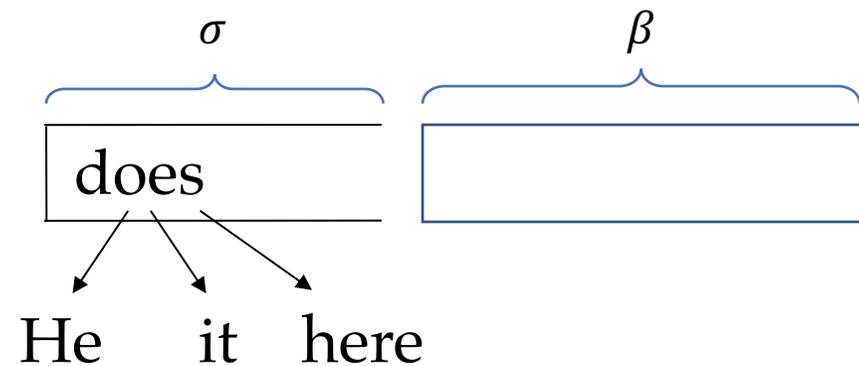


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Finish*



σ : stack β : buffer

Contents

- **11.1 Transition-based Structured Prediction**
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Transition-based Modeling

- Given a state s_{i-1} , our goal is to disambiguate all possible actions $a_i \in \text{POSSIBLEACTIONS}(s_{i-1})$ by using a discriminative model to score transition actions:

$$\text{score}(a_i | s_{i-1}) = \vec{\theta} \cdot \vec{\phi}(s_{i-1}, a_i)$$

- $\vec{\theta}$: model parameter vector
- $\vec{\phi}(s_{i-1}, a_i)$: feature vector on the input-output pair (s_{i-1}, a_i)

Greedy Local Method

- Training

- Training dataset $D = \{(X_i, Y_i)\}_{i=1}^N$

↓ break down

- A sequence of gold transitions: $(s_{j-1}^{(i)}, a_j^{(i)})$

↓ merge all the state-action pairs

- A training set for the discriminative model

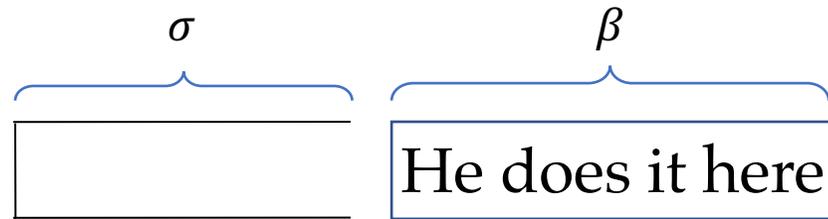
- Testing

- Start from initial state $s_0(X)$

- Repeatedly find $\hat{a}_i = \operatorname{argmax}_{\alpha} \vec{\theta} \cdot \vec{\phi}(s_{i-1}, \alpha)$

Greedy Local Method

- NEXT ACTION: *Shift*

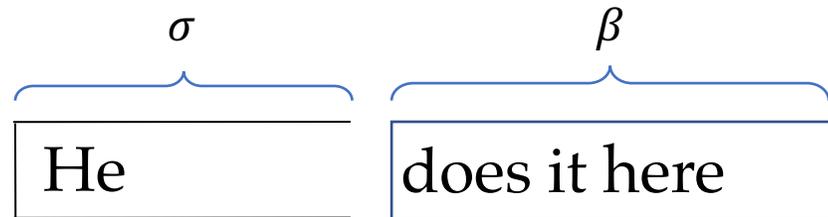


σ : stack

β : buffer

Greedy Local Method

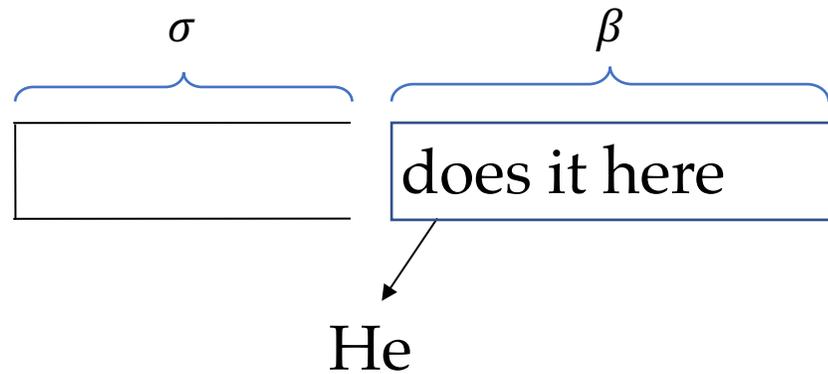
- NEXT ACTION: *Left – Arc*



σ : stack β : buffer

Greedy Local Method

- NEXT ACTION: *Shift*

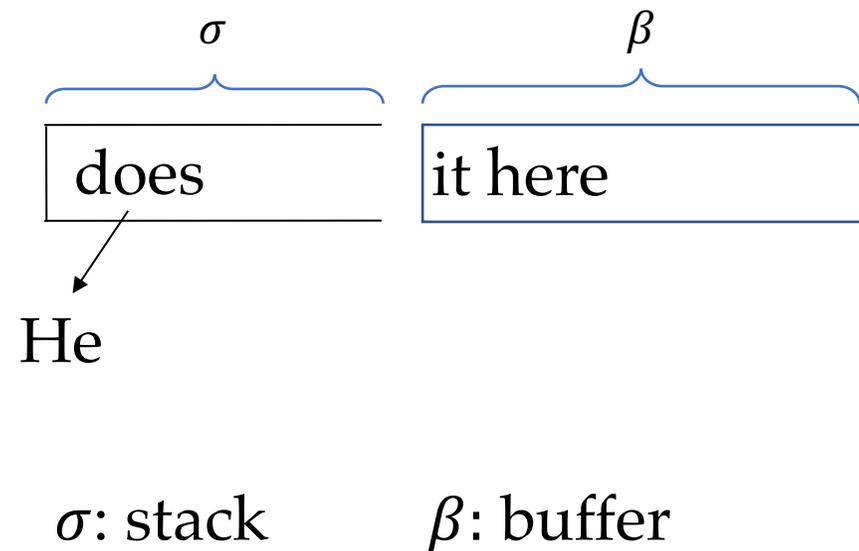


σ : stack

β : buffer

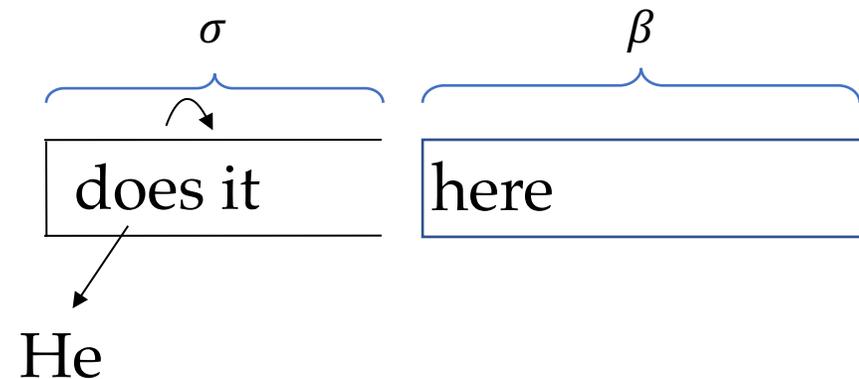
Greedy Local Method

- NEXT ACTION: *Right – Arc*



Greedy Local Method

- NEXT ACTION: *Reduce*

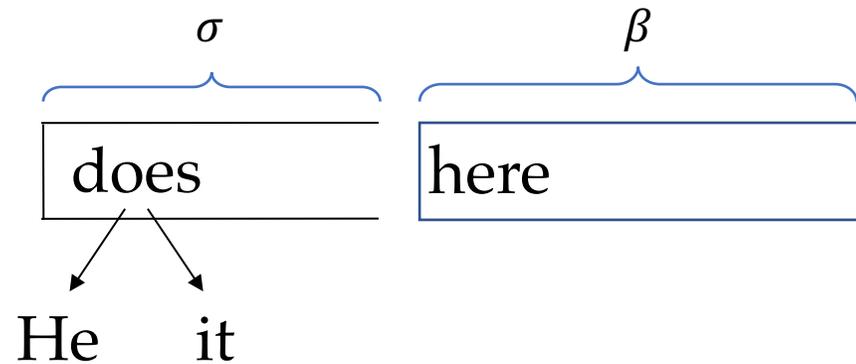


σ : stack

β : buffer

Greedy Local Method

- NEXT ACTION: *Right – Arc*

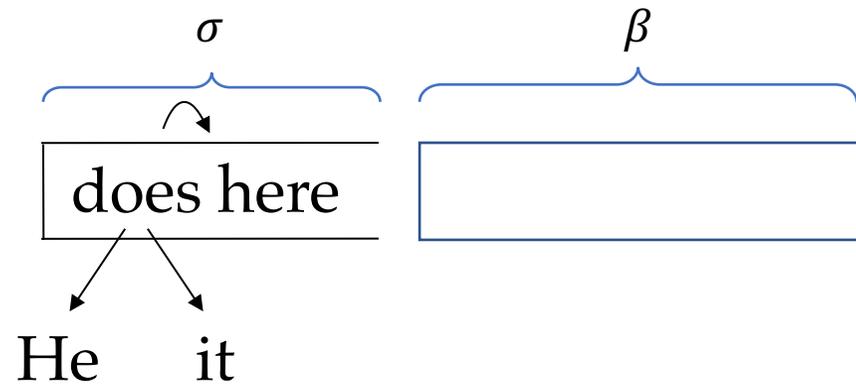


σ : stack

β : buffer

Greedy Local Method

- NEXT ACTION: *Reduce*

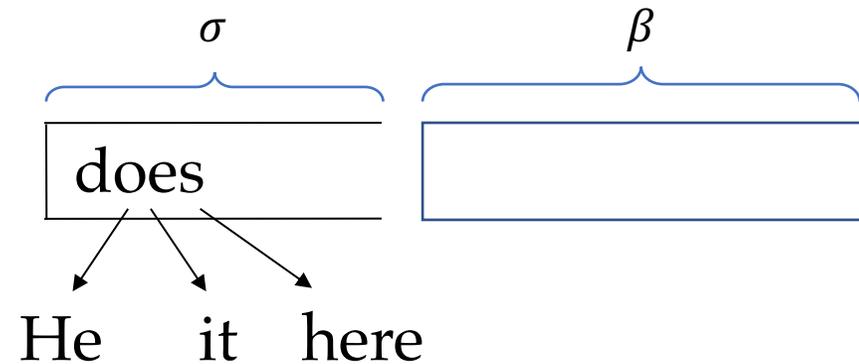


σ : stack

β : buffer

Greedy Local Method

- NEXT ACTION: *Finish*



σ : stack β : buffer

Greedy Local Method

- An Example

He does it here

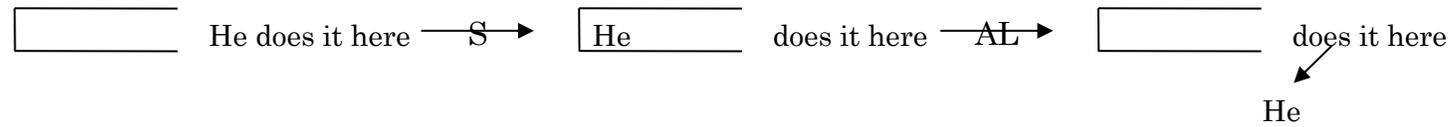
Greedy Local Method

- An Example

He does it here \xrightarrow{S} He does it here

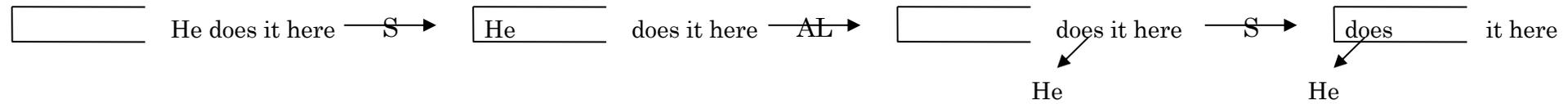
Greedy Local Method

- An Example



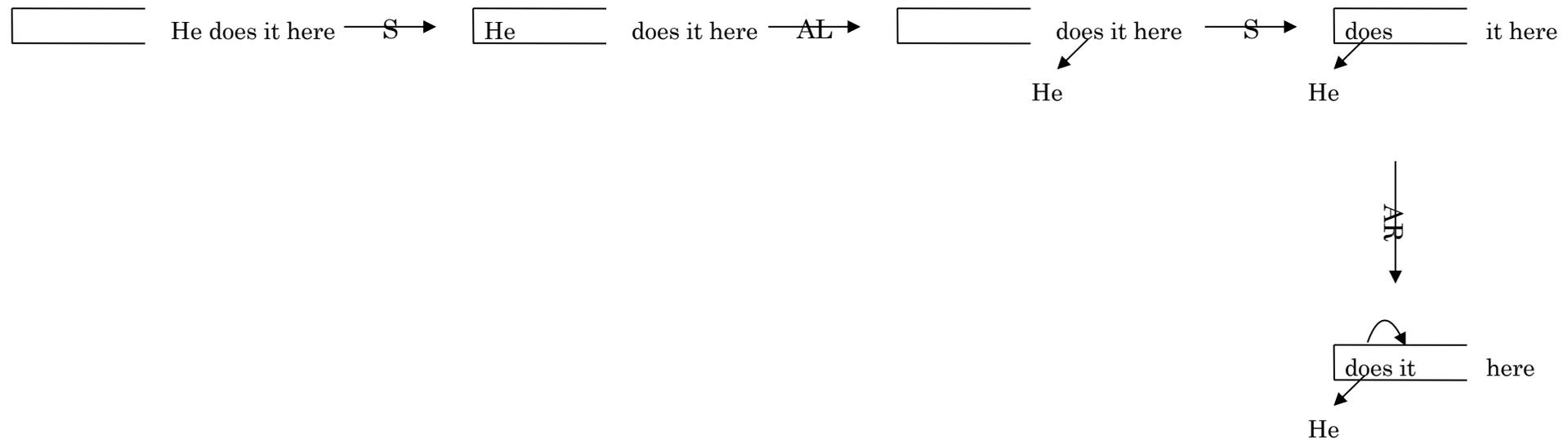
Greedy Local Method

- An Example



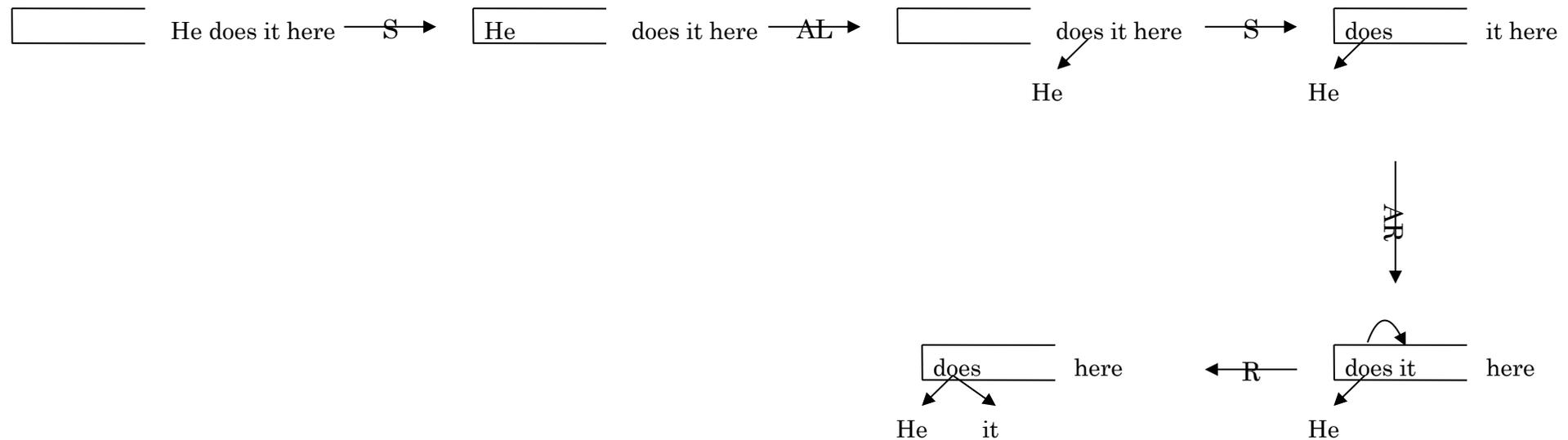
Greedy Local Method

- An Example



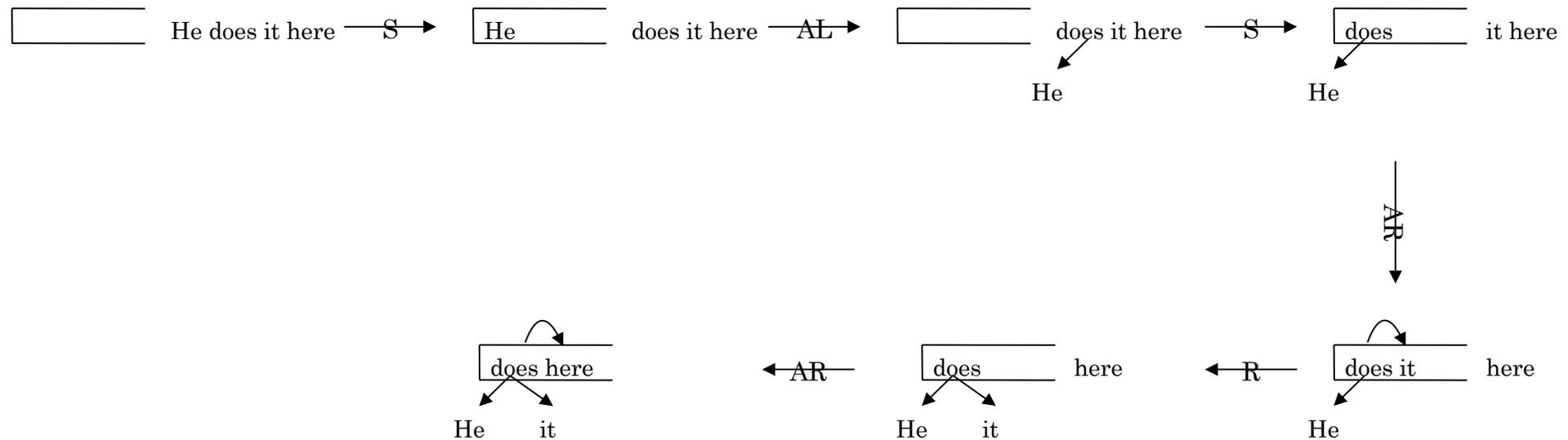
Greedy Local Method

- An Example



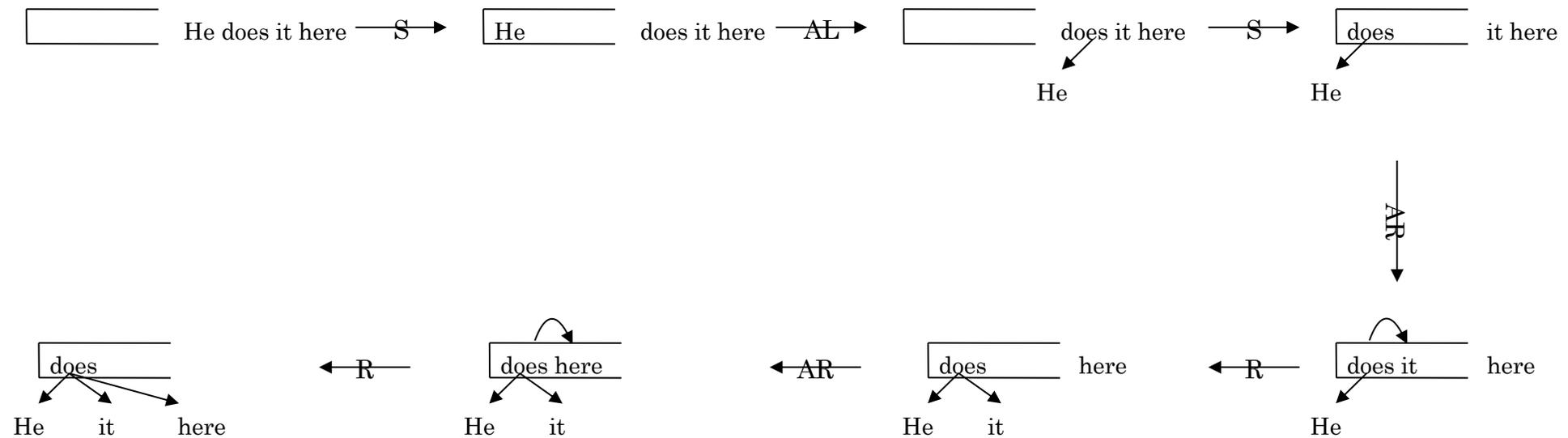
Greedy Local Method

- An Example



Greedy Local Method

- An Example



Problem of Greedy Local Modeling

- In a **globally** optimal action sequence, each action may not necessarily be the optimal choice **locally**.
- Result in **error propagation**.
- Model does not see incorrect states during training.

Contents

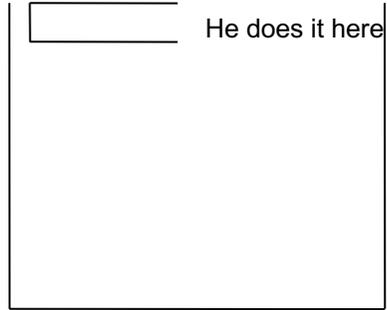
- **11.1 Transition-based Structured Prediction**
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

- In a **globally** optimal action sequence, each action may not necessarily be the optimal choice **locally**.
- Result in **error propagation**.

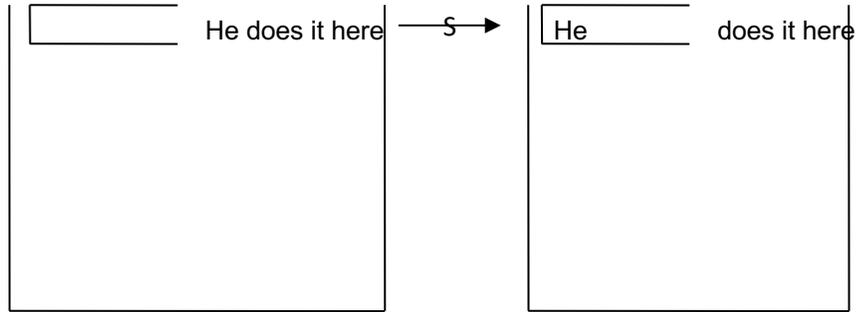
Solution:

- Beam search use a globally trained model:
 - Given an input X , a global transition-based model calculates $score(A|X)$ directly, where $A_{1:|A|} = a_1 a_2 \dots a_{|A|}$ is a sequence of transition actions a_i for building an output structure for X .

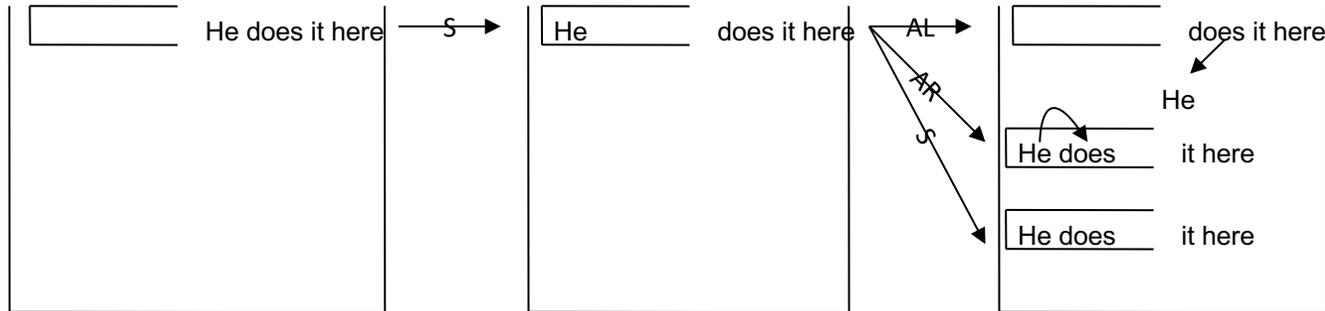
Beam search decoding example



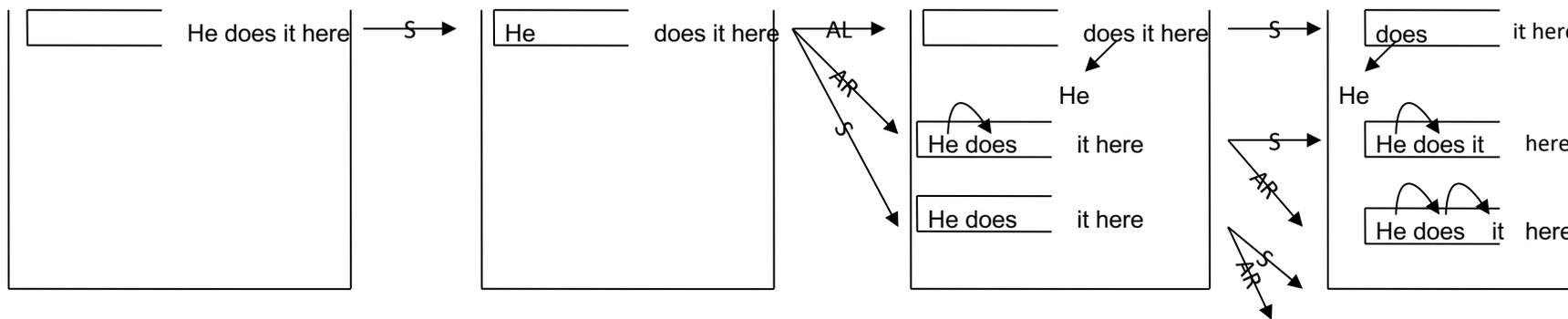
Beam search decoding example



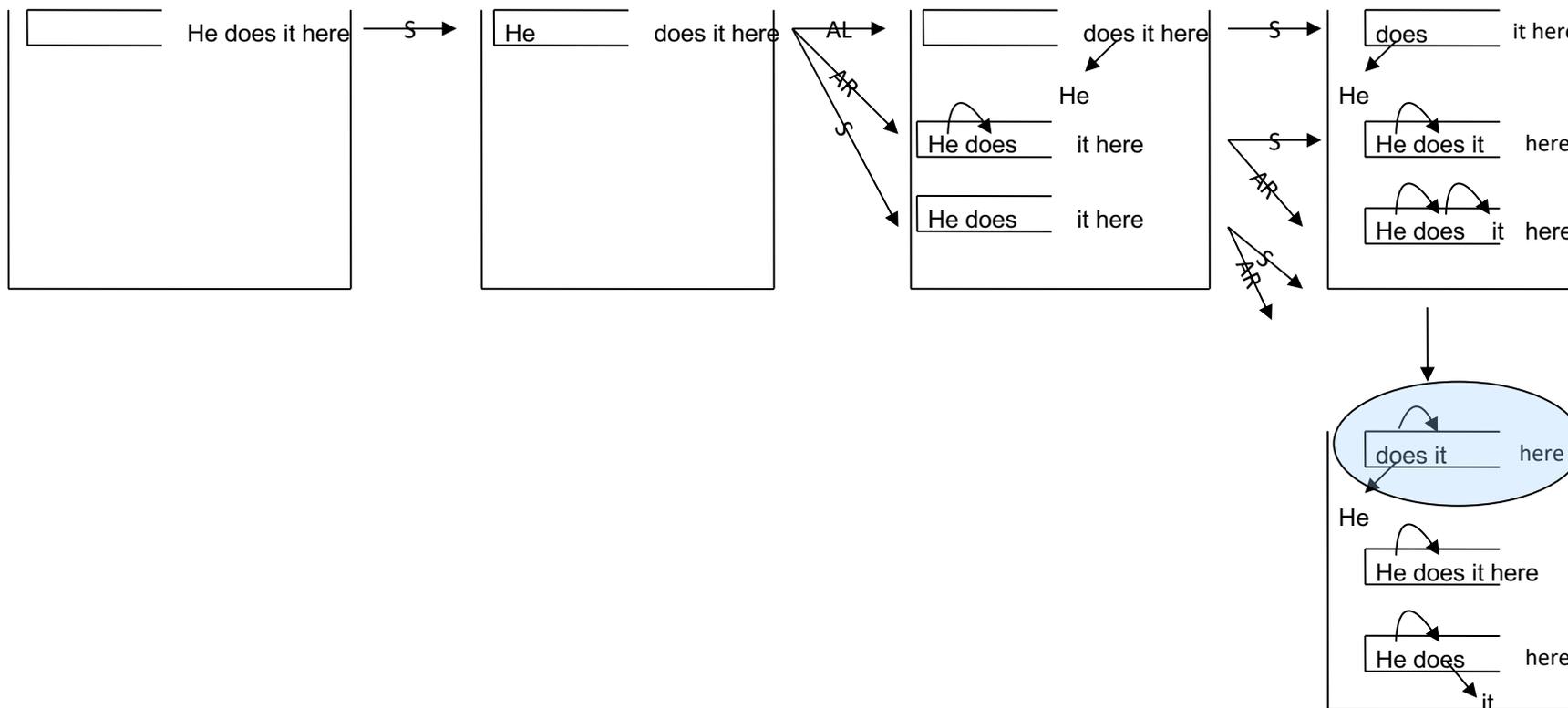
Beam search decoding example



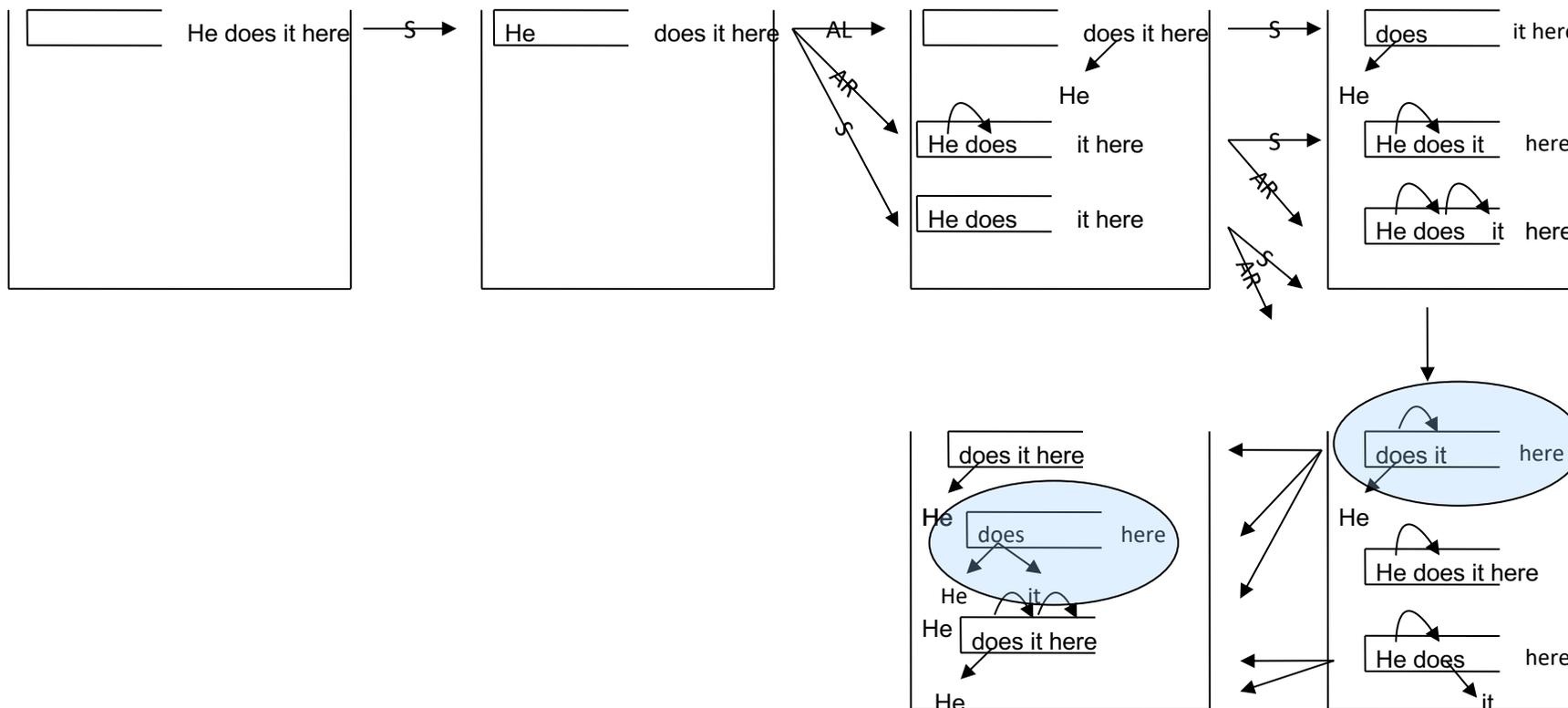
Beam search decoding example



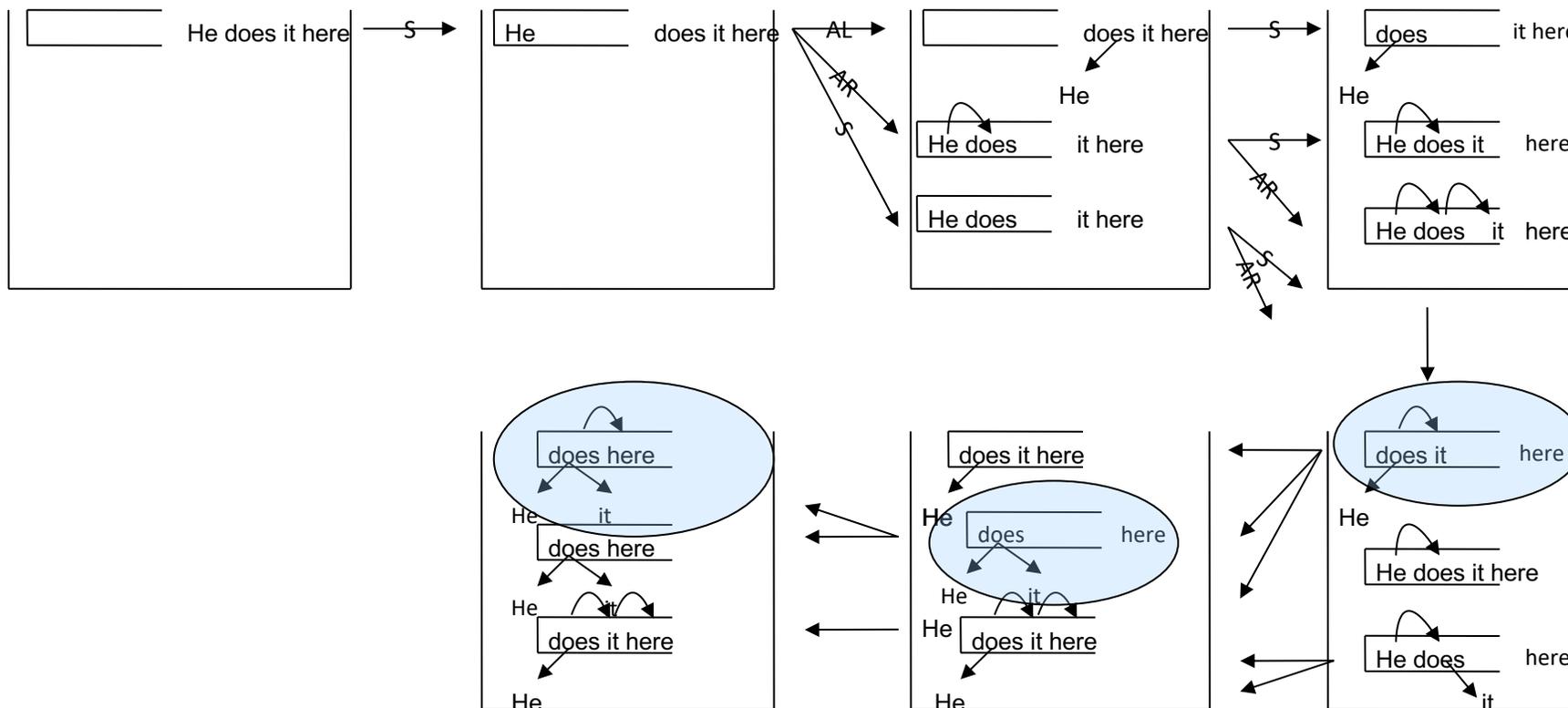
Beam search decoding example



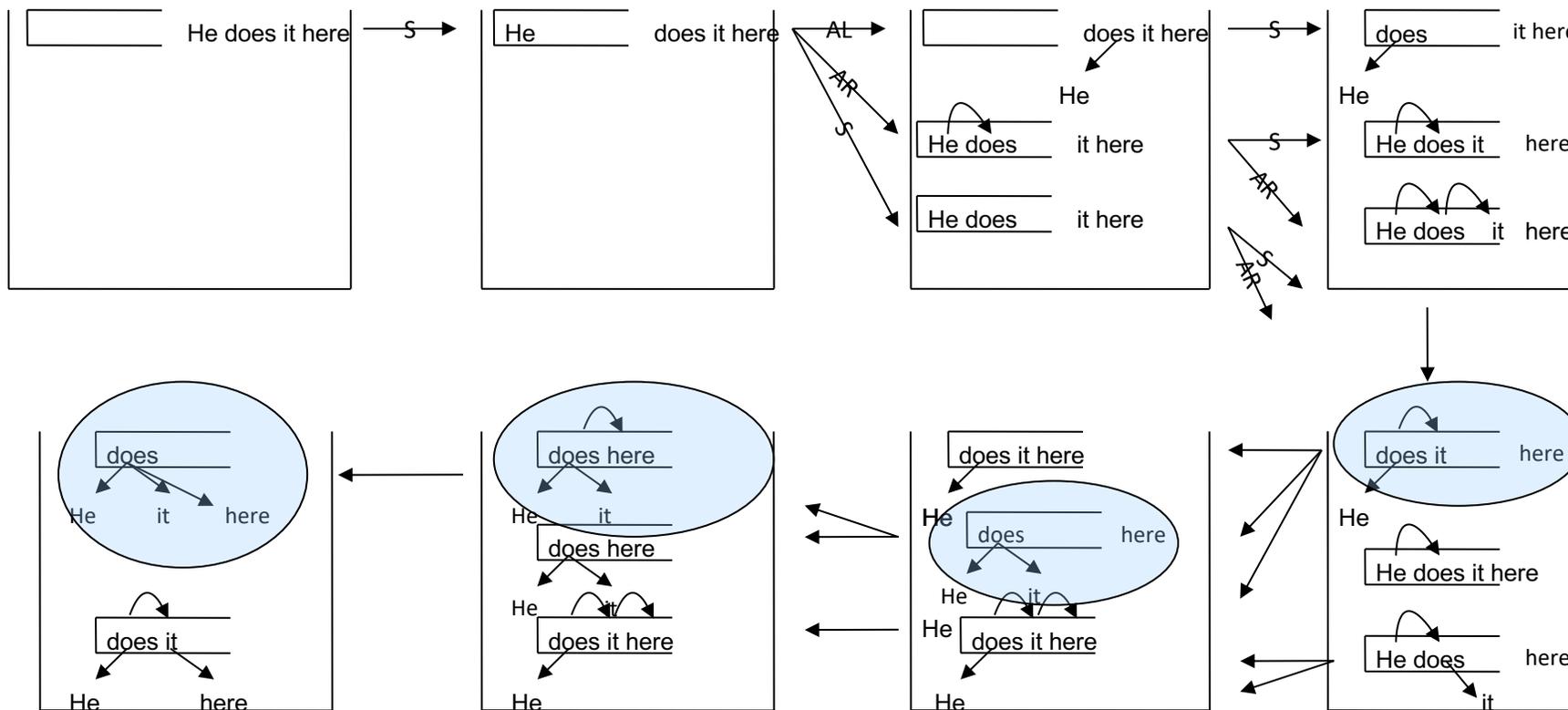
Beam search decoding example



Beam search decoding example



Beam search decoding example



Solution

- Use global linear model to calculate $score(A|X)$

$$score(A|X) = \vec{\theta} \cdot \vec{\phi}(A, X)$$

↓ decompose $\vec{\phi}(A|X)$ for incremental decoding

$$\vec{\phi}(A|X) = \sum_{i=1}^{|A|} \vec{\phi}(a_i, s_{i-1})$$



$$\begin{aligned} score(A|X) &= \vec{\theta} \cdot \vec{\phi}(A, X) = \vec{\theta} \cdot \left(\sum_{i=1}^{|A|} \vec{\phi}(a_i, s_{i-1}) \right) \\ &= \sum_{i=1}^{|A|} (\vec{\theta}) \cdot \left\{ \vec{\phi}(a_i, s_{i-1}) \right\} \end{aligned}$$

Beam search decoding algorithm

Inputs: $\vec{\theta}$ —discriminative linear model parameters;
 X — task input;
 K — beam size;

Initialization: $agenda \leftarrow [(\text{STARTSTATE}(X), 0)]$;

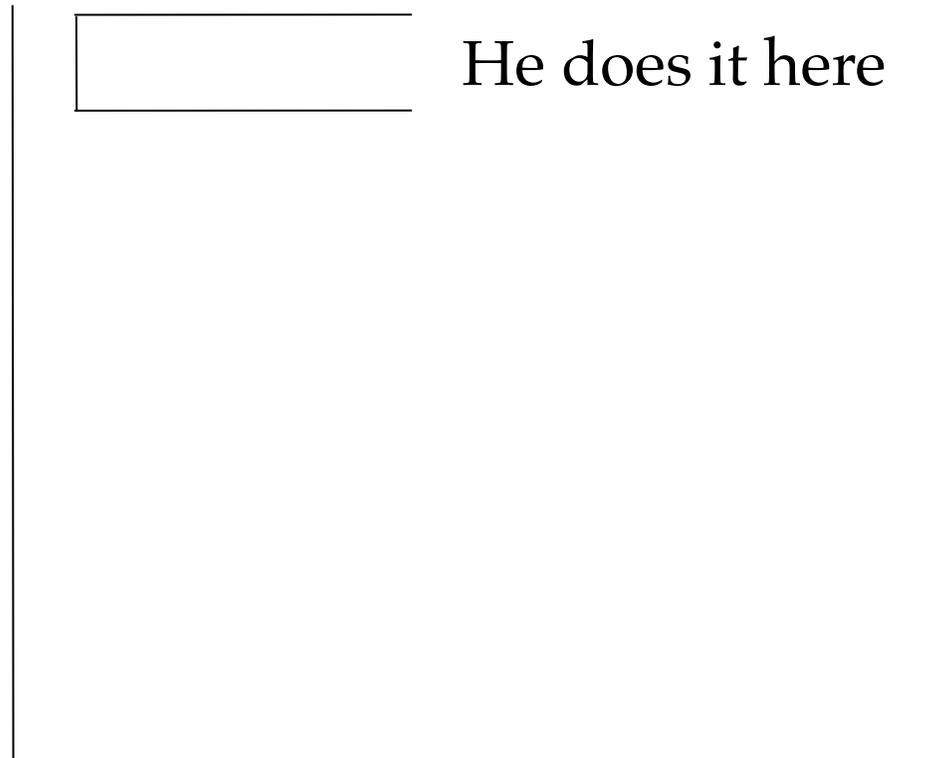
Algorithm:

```
while not ALLTERMINAL( $agenda$ ) do
   $to\_expand \leftarrow agenda$ ;
   $agenda \leftarrow []$ ;
  for  $(state, score) \in to\_expand$  do
    for  $a \in \text{POSSIBLEACTIONS}(state)$  do
       $new\_state \leftarrow \text{EXPAND}(state, a)$ ;
       $new\_score \leftarrow score + \vec{\theta} \cdot \vec{\phi}(state, a)$ ;
      APPEND( $agenda$ ,  $(new\_state, new\_score)$ );
   $agenda \leftarrow \text{TOP-K}(agenda, K)$ ;
```

Output: $\text{TOP-K}(agenda, 1)[0]$;

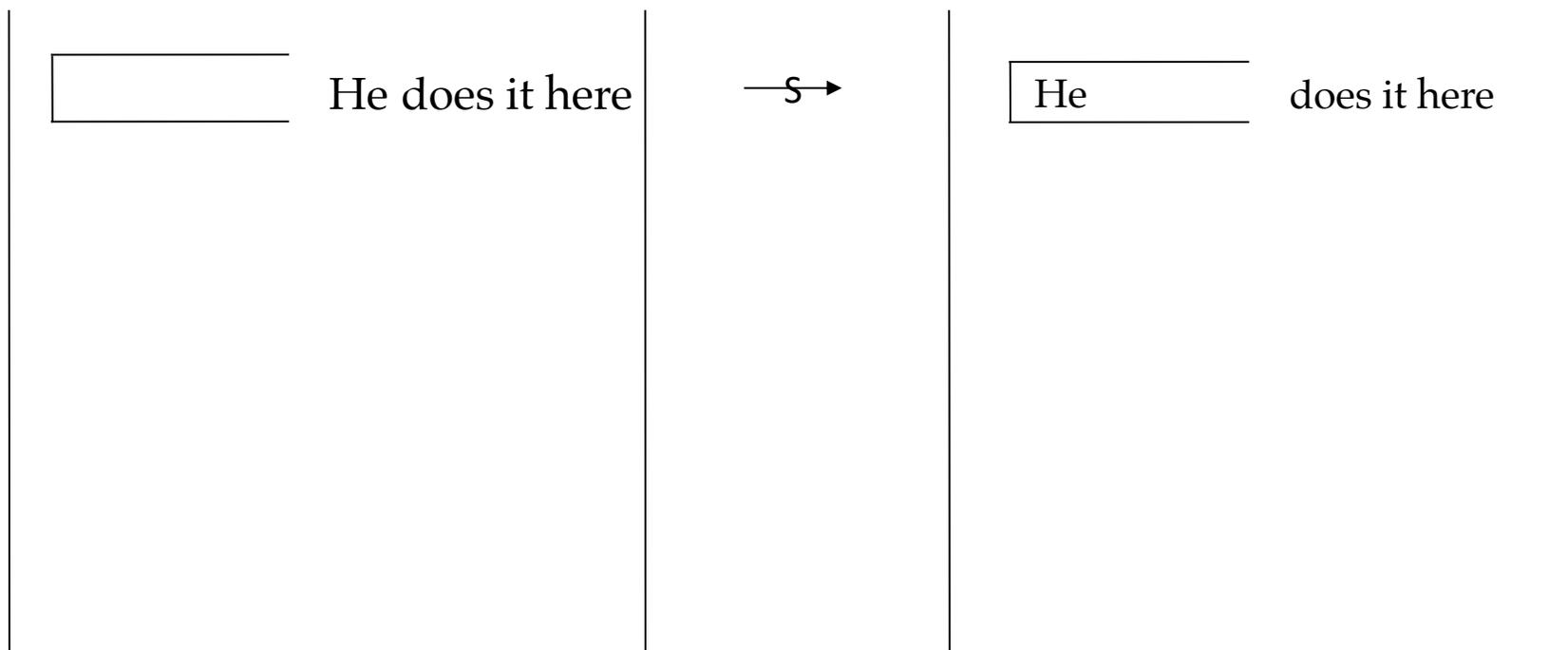
Beam search decoding example

- Dependency parsing



Beam search decoding example

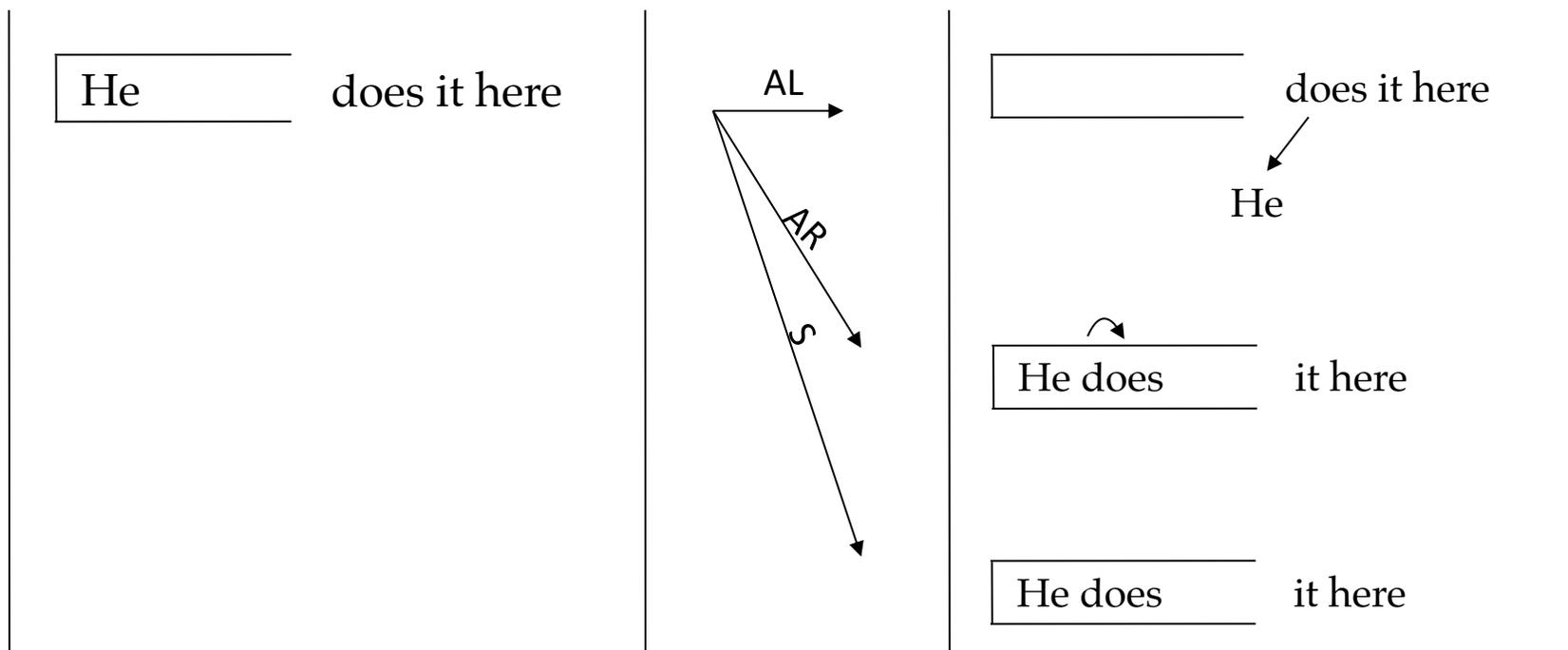
- Dependency parsing



S-SHIFT

Beam search decoding example

- Dependency parsing

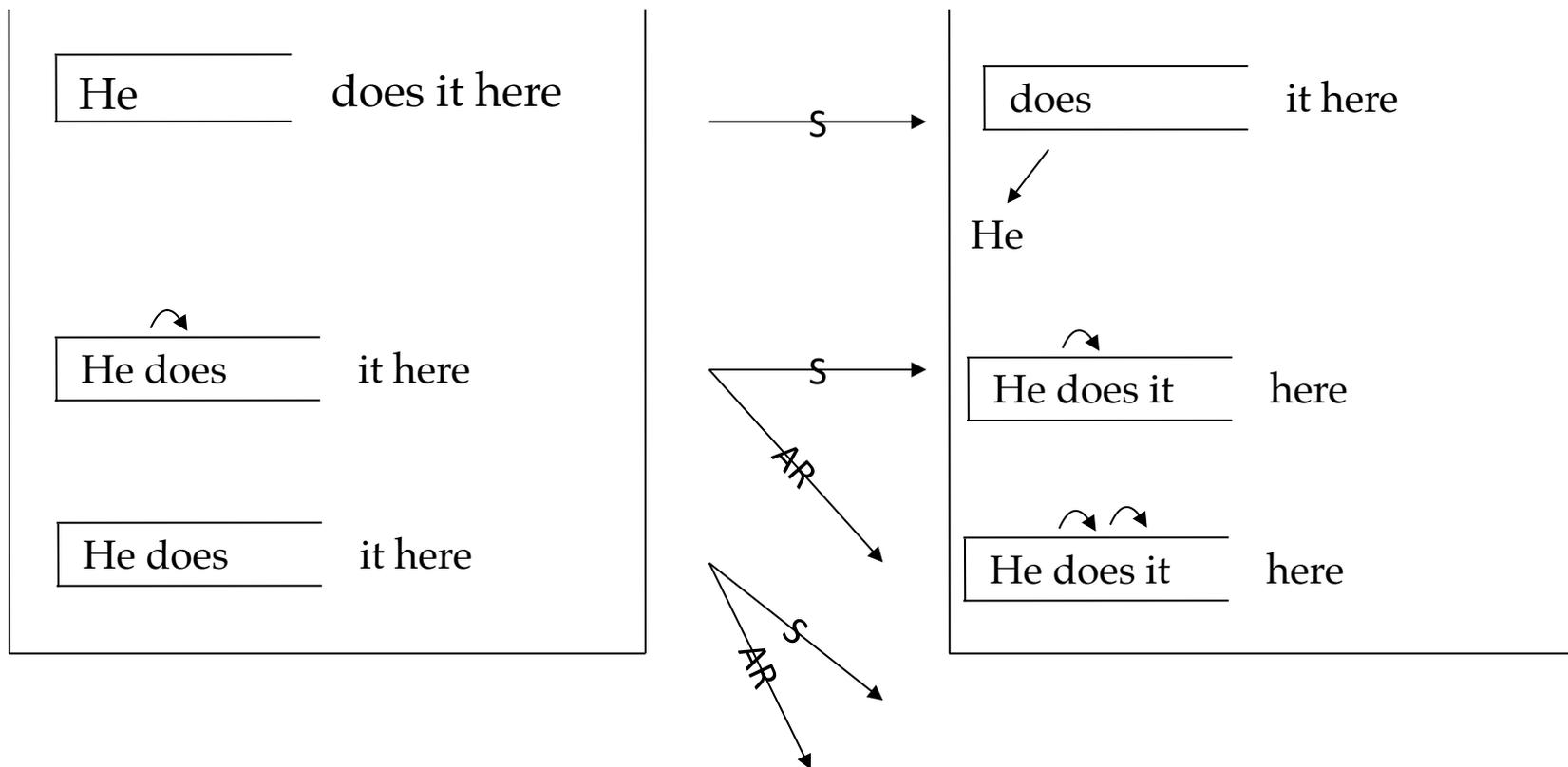


AL-LEFT-ARC

AR-RIGHT-ARC

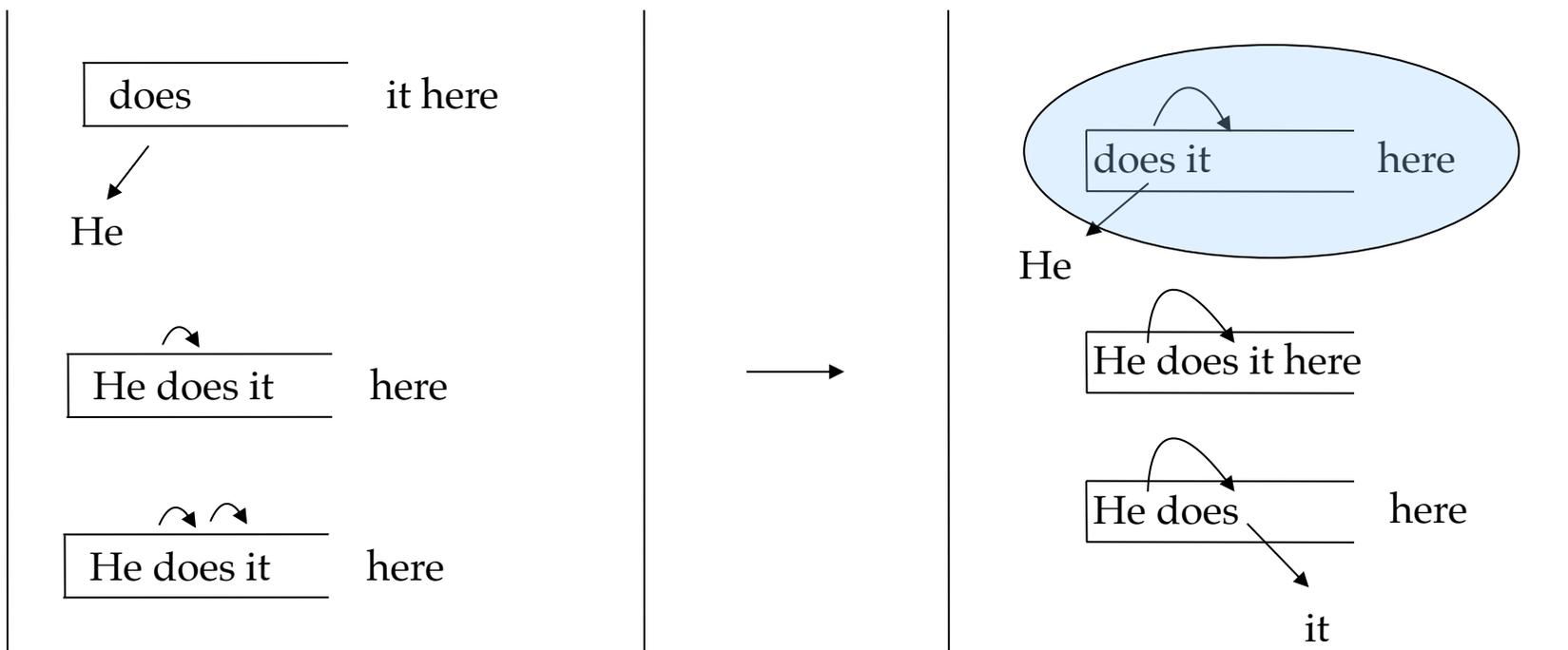
Beam search decoding example

- Dependency parsing



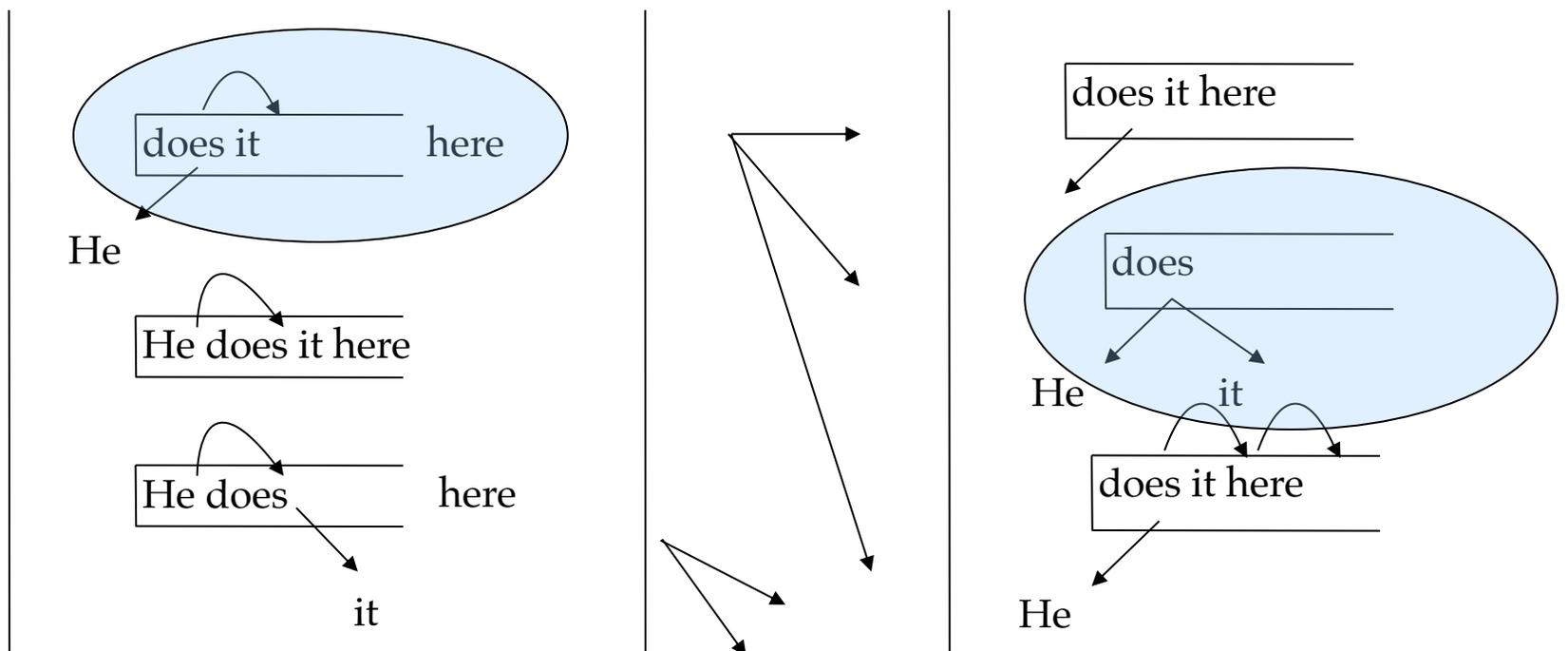
Beam search decoding example

- Dependency parsing



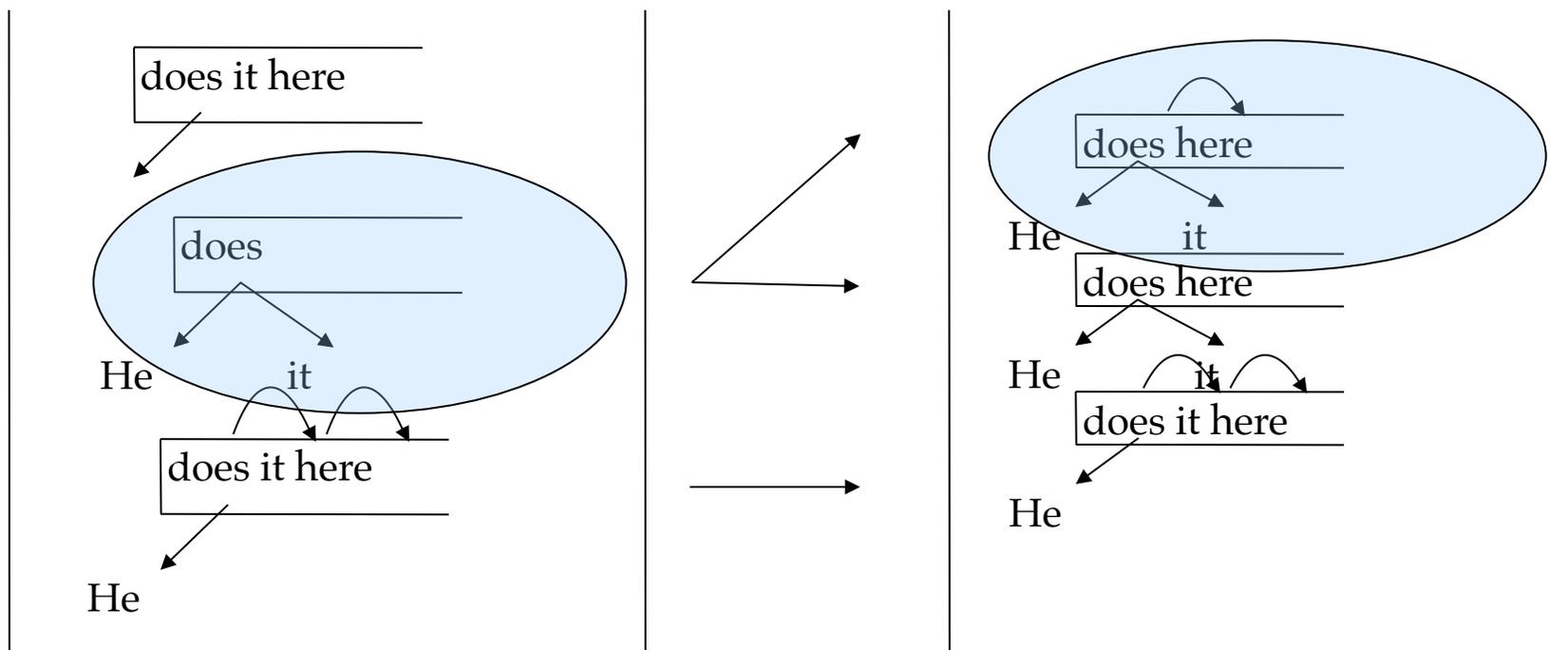
Beam search decoding example

- Dependency parsing



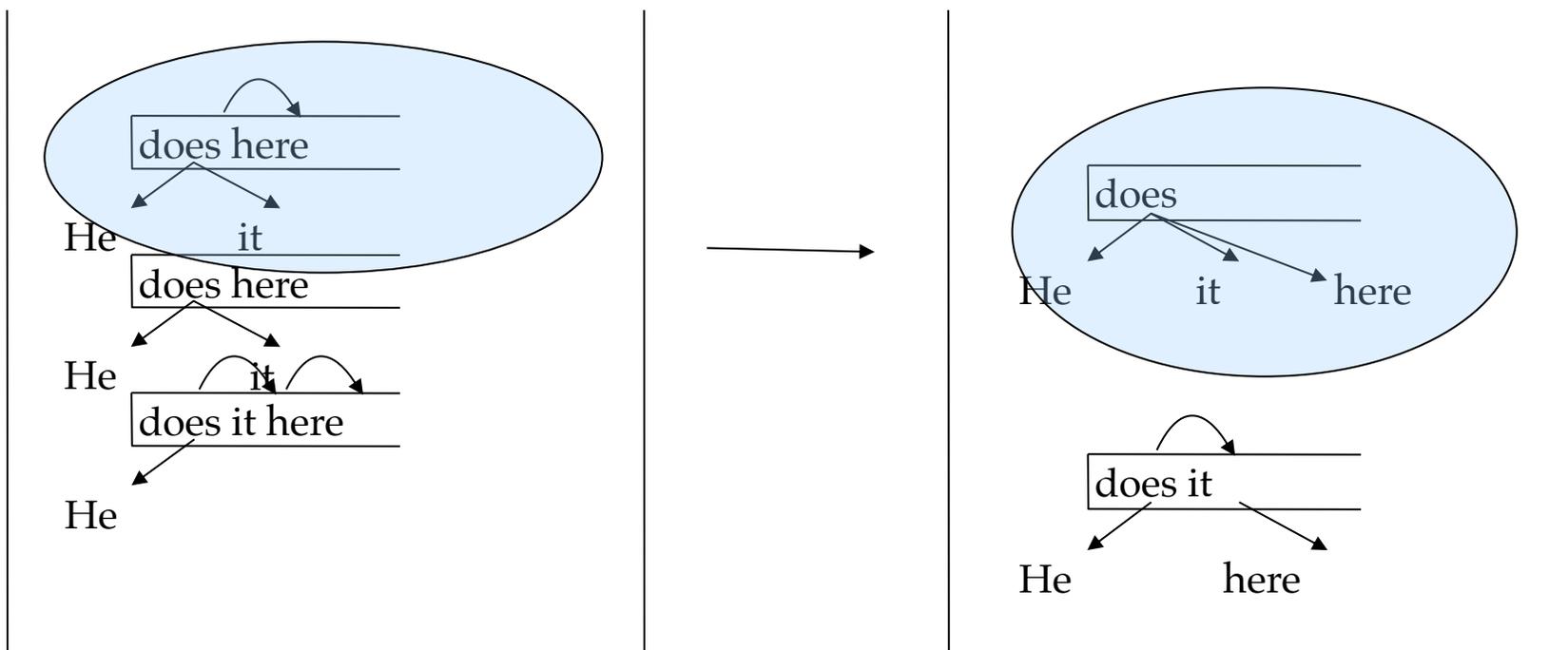
Beam search decoding example

- Dependency parsing

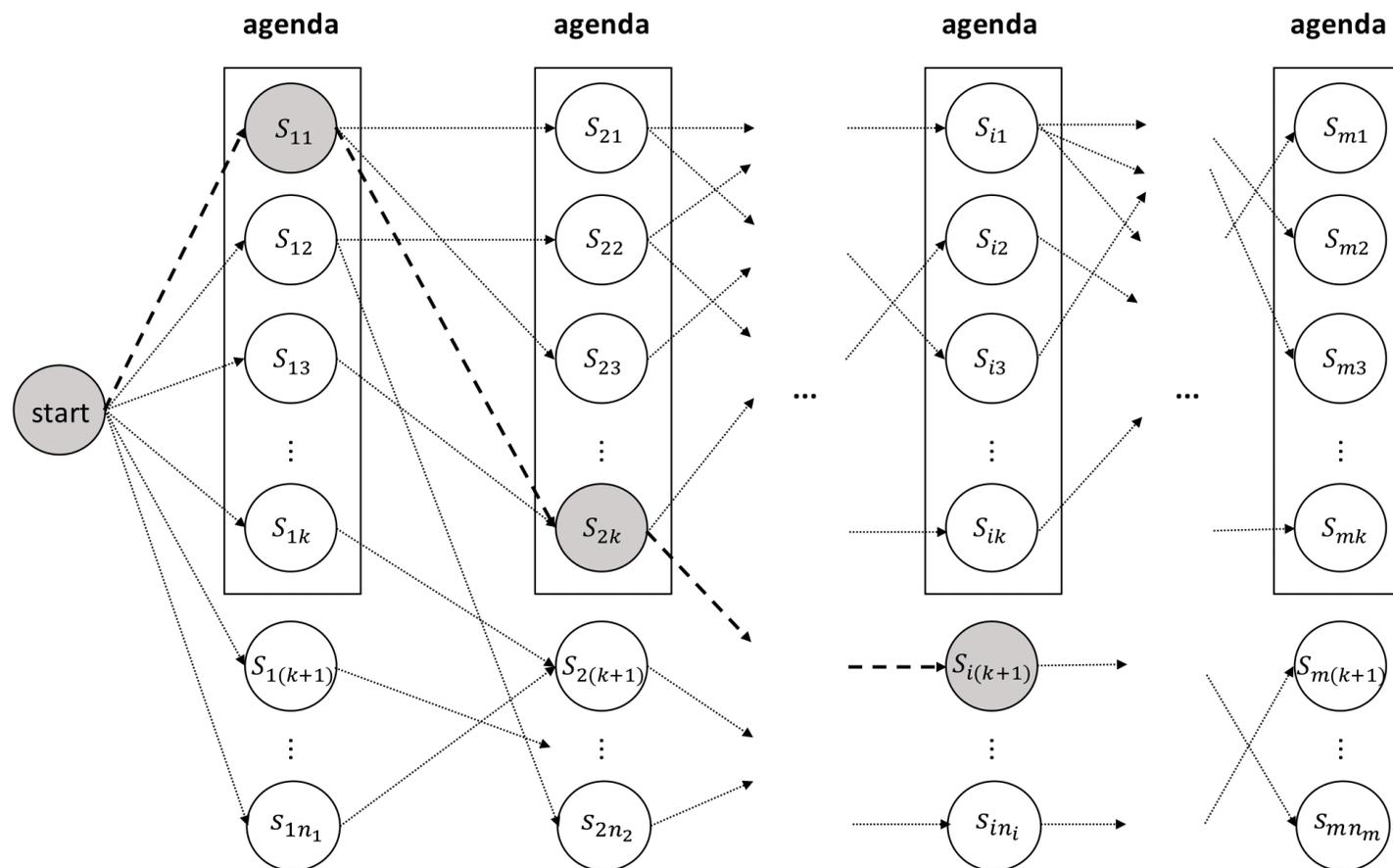


Beam search decoding example

- Dependency parsing



Beam search decoding example



Gold Sequence of Action: $\text{start} \rightarrow S_{11} \rightarrow S_{2k} \rightarrow \dots \rightarrow S_{i(k+1)}$

Beam search training algorithm

Inputs: D — gold standard training set;
 K — beam size;
 T — number of training iterations;

Initialisation: $\vec{\theta} \leftarrow 0$;

Algorithm:

```
for  $t \in [1, \dots, T]$  do
  for  $(X, Y) \in D$  do
     $G \leftarrow \text{GOLDACTIONSEQ}(X, Y)$ ;
     $\text{agenda} \leftarrow [(\text{STARTSTATE}(X), 0)]$ ;
     $\text{gold\_state} \leftarrow \text{STARTSTATE}(X)$ ;
     $i \leftarrow 0$ ;
    while not ALLTERMINAL( $\text{agenda}$ ) do
       $i \leftarrow i + 1$ ;
       $\text{to\_expand} \leftarrow \text{agenda}$ ;
       $\text{agenda} \leftarrow []$ ;
      for  $(\text{state}, \text{score}) \in \text{to\_expand}$  do
        for  $a \in \text{POSSIBLEACTIONS}(\text{state})$  do
           $\text{new\_state} \leftarrow \text{EXPAND}(\text{state}, a)$ ;
           $\text{new\_score} \leftarrow \text{score} + \vec{\theta} \cdot \vec{\phi}(\text{state}, a)$ ;
          APPEND( $\text{agenda}, (\text{new\_state}, \text{new\_score})$ );
       $\text{agenda} \leftarrow \text{TOP-K}(\text{agenda}, K)$ ;
       $\text{gold\_state} \leftarrow \text{EXPAND}(\text{gold\_state}, G[i])$ ;
      if not CONTAIN( $\text{agenda}, \text{gold\_state}$ ) then
         $\text{pos} \leftarrow \text{gold\_state}$ ;
         $\text{neg} \leftarrow \text{TOP-K}(\text{agenda}, 1)[0]$ ;
         $\vec{\theta} \leftarrow \vec{\theta} + \vec{\phi}(\text{pos}) - \vec{\phi}(\text{neg})$ ;
        continue( $(W, G) \in D$ )
    if  $\text{gold\_state} \neq \text{TOP-K}(\text{agenda}, 1)[0]$  then
       $\vec{\theta} \leftarrow \vec{\theta} + \vec{\phi}(\text{gold\_state}) - \vec{\phi}(\text{TOP-K}(\text{agenda}, 1)[0])$ ;
```

Output: $\vec{\theta}$;

Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- **11.2 Transition-based Constituent Parsing**
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Transition-based constituent parsing

- Three steps in transition-based modeling
 - Find the state transition process
 - Define global feature vector
 - Apply the standard learning and search framework

Transition-based constituent parsing

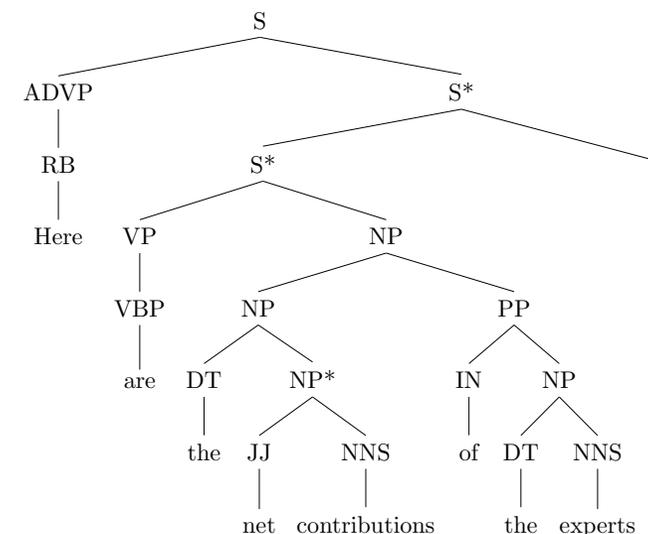
- **State**

σ : stack of partially constituent outputs.

β : buffer of the next incoming words.

- **actions**

SHIFT, REDUCE-L/R-X, UNARY-X, IDLE



Transition-based constituent parsing

Shift-reduce constituent parsing

Axiom:	$([], W_{1:n})$	SHIFT:	$\frac{(\sigma, w_0 \beta)}{(\sigma w_0, \beta)}$
Goal:	$(\sigma, [])$		
	$(\sigma s_1 s_0, \beta)$		$(\sigma s_1 s_0, \beta)$
REDUCE-L-X:	$(\sigma \begin{array}{c} X \\ \swarrow \quad \searrow \\ s_1 \quad s_0 \end{array}, \beta)$	REDUCE-R-X:	$(\sigma \begin{array}{c} X \\ / \quad \searrow \\ s_1 \quad s_0 \end{array}, \beta)$
	$(\sigma s_0, \beta)$		$(\sigma, [])$
UNARY-X:	$(\sigma \begin{array}{c} X \\ \downarrow \\ s_0 \end{array}, \beta)$	IDLE:	$(\sigma, [])$

Transition-based constituent parsing

Shift-reduce constituent parsing

- Example
 - Shift

stack



buffer

DT	ADJ	NN	VV	ADJ	NNS	.
The	little	boy	likes	red	tomatoes	.

Transition-based constituent parsing

Shift-reduce constituent parsing

- Example
 - Shift

stack

	DT
	The

buffer

ADJ	NN	VV	ADJ	NNS	.
little	boy	likes	red	tomatoes	.

Shift-reduce constituent parsing

- Example
 - Shift

stack

DT	ADJ
The	little

buffer

NN	VV	ADJ	NNS	.
boy	likes	red	tomatoes	.

Transition-based constituent parsing

Shift-reduce constituent parsing

- Example
 - Reduce-R-NP

stack

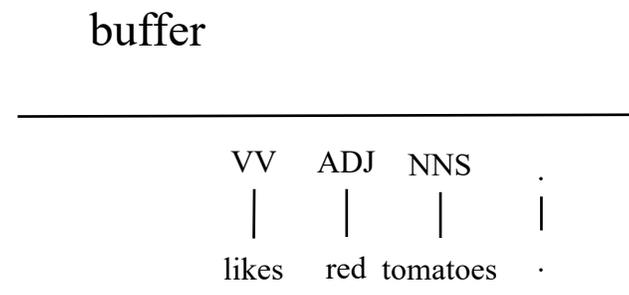
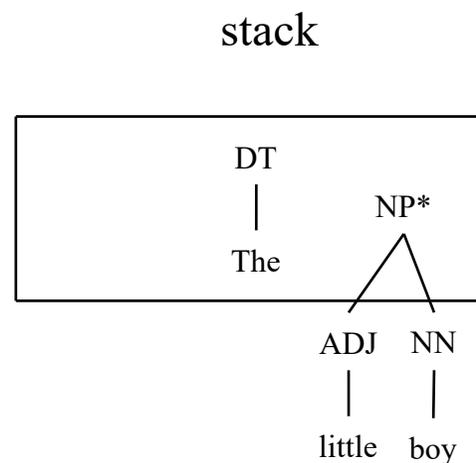
DT	ADJ	NN
The	little	boy

buffer

VV	ADJ	NNS	.
likes	red	tomatoes	.

Shift-reduce constituent parsing

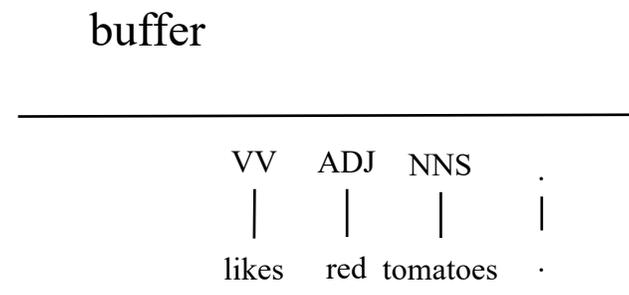
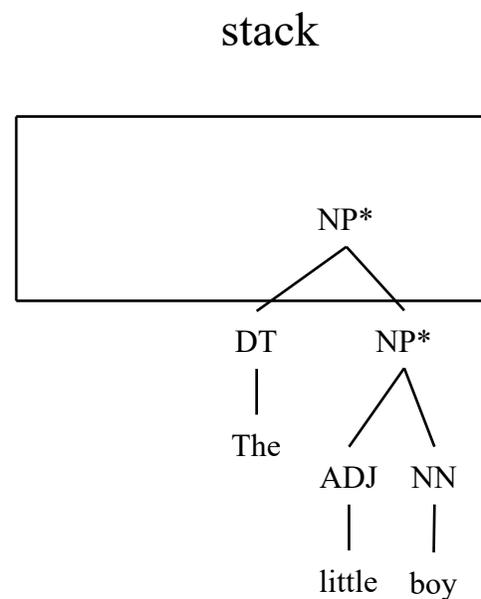
- Example
 - Reduce-R-NP



Transition-based constituent parsing

Shift-reduce constituent parsing

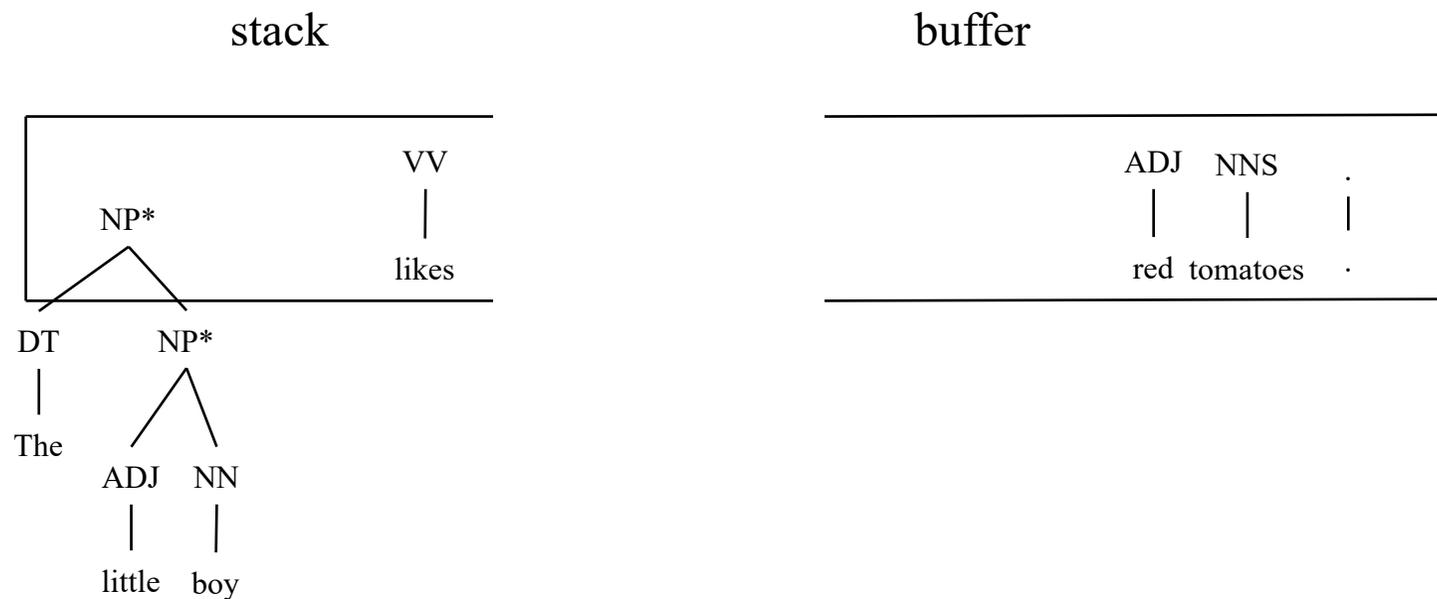
- Example
 - Shift



Transition-based constituent parsing

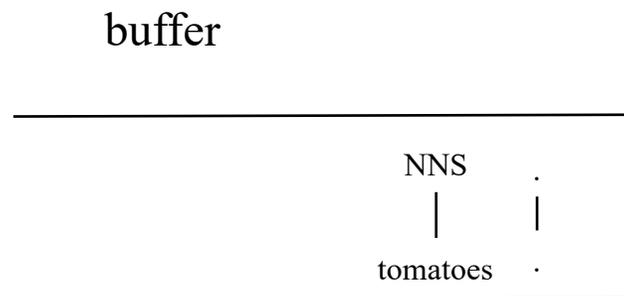
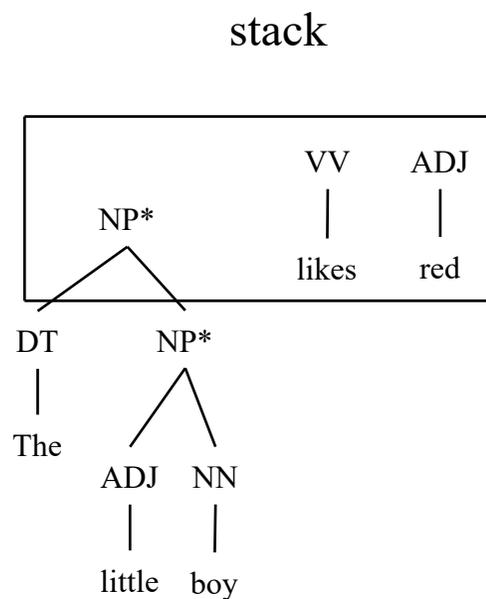
Shift-reduce constituent parsing

- Example
 - Shift



Shift-reduce constituent parsing

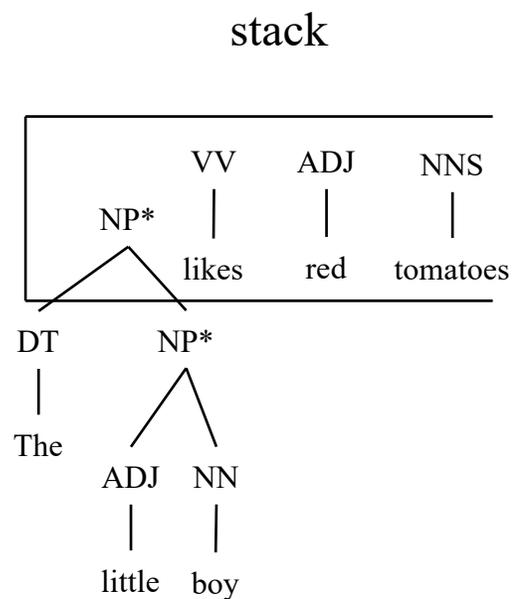
- Example
 - Shift



Transition-based constituent parsing

Shift-reduce constituent parsing

- Example
 - Reduce-R-NP

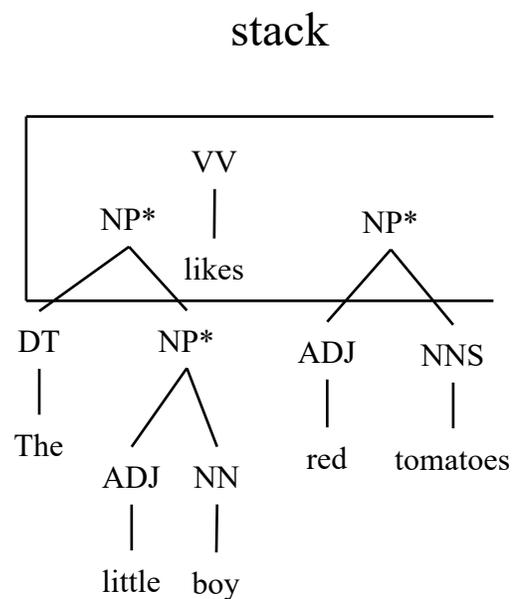


buffer

.

Shift-reduce constituent parsing

- Example
 - Reduce-L-VP



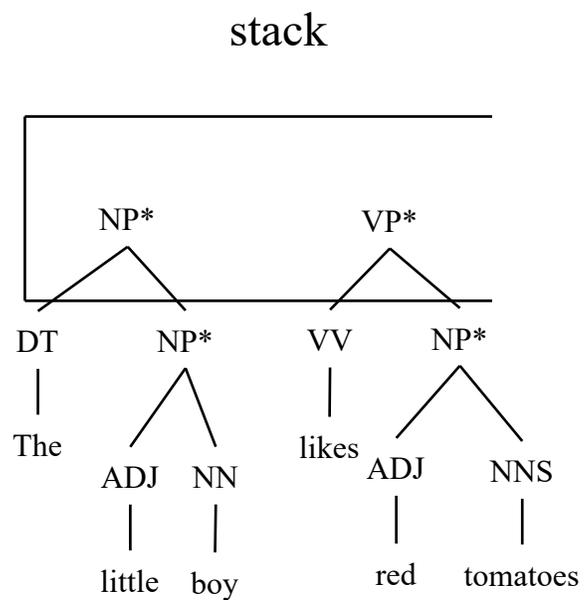
buffer



Transition-based constituent parsing

Shift-reduce constituent parsing

- Example
 - Shift

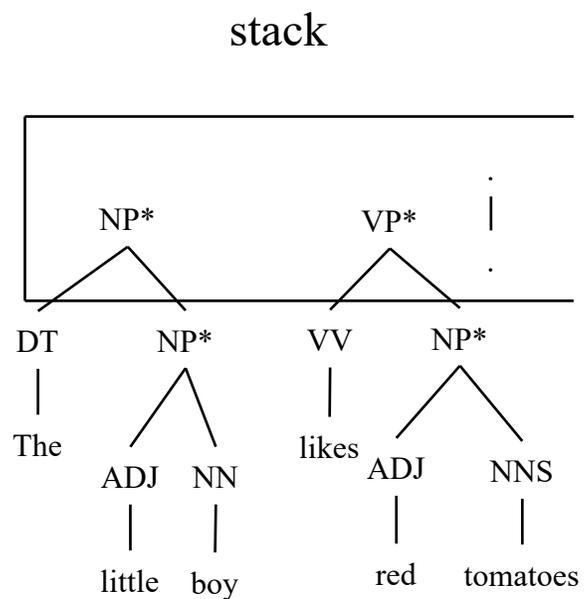


buffer

.

Shift-reduce constituent parsing

- Example
 - Reduce-L-S

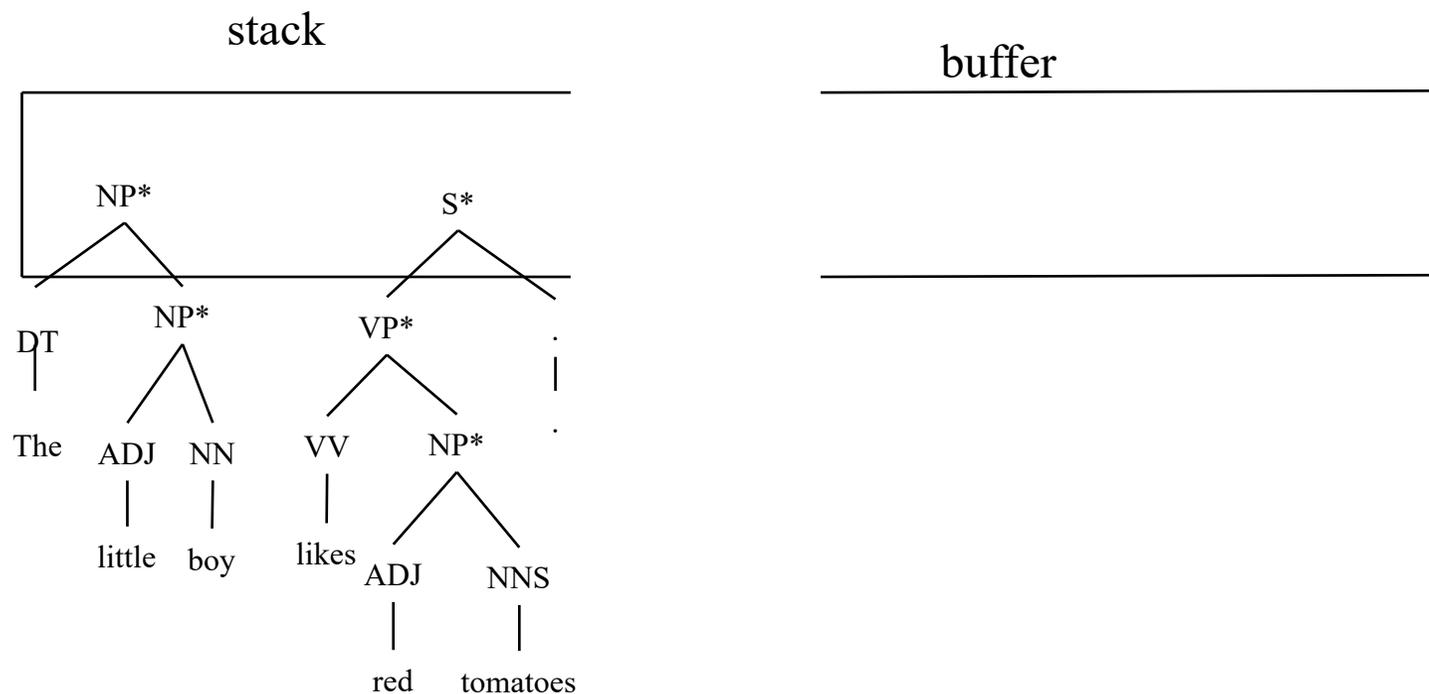


buffer

Transition-based constituent parsing

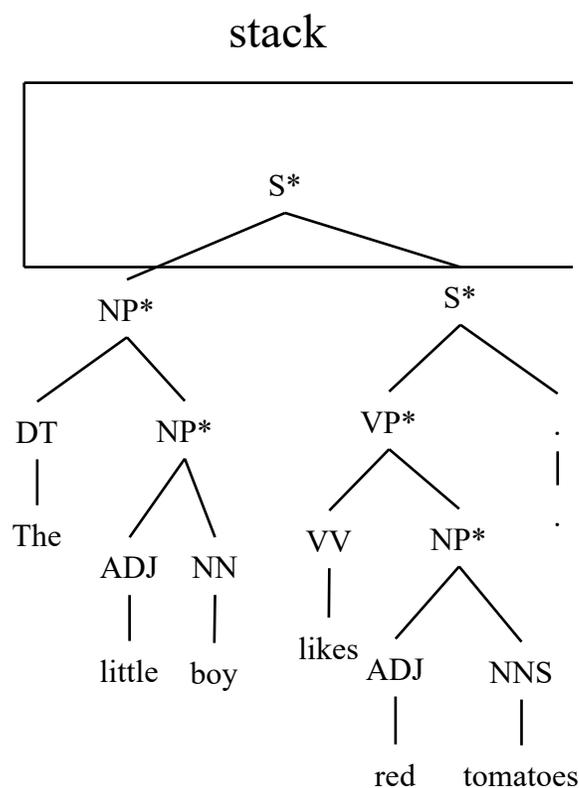
Shift-reduce constituent parsing

- Example
 - Reduce-R-S



Shift-reduce constituent parsing

- Example
 - Terminate

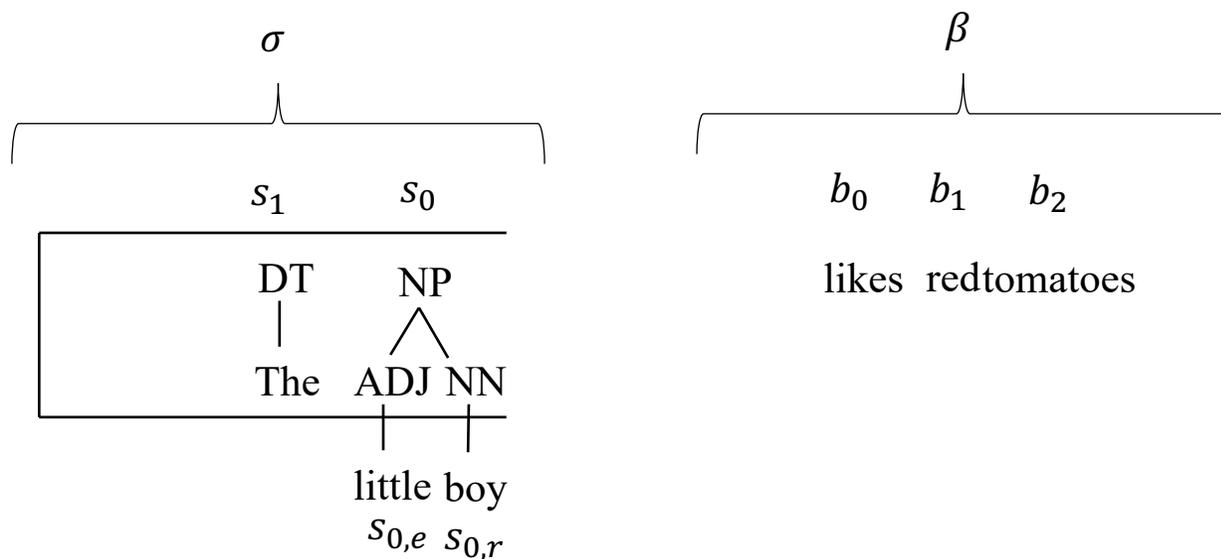


buffer



Transition-based constituent parsing

- Feature Templates for Shift-reduce Constituent Parser
- Example state



- Feature Templates for Shift-reduce Constituent Parser

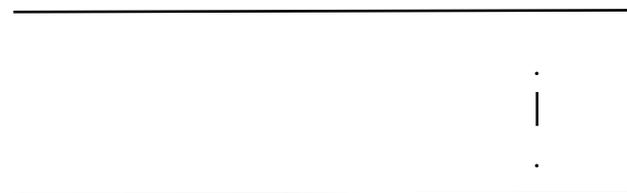
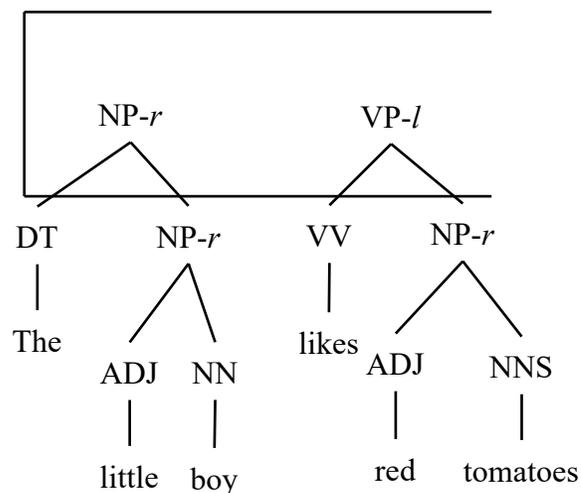
Feature type	Feature Template
unigrams	$s_0ps_0c, s_0ws_0c, s_1ps_1c, s_1ws_1c, s_2ps_2c, s_2ws_2c, s_3ps_3c, s_3ws_3c, b_0wb_0p, b_1wb_1p, b_2wb_2p, b_3wb_3p, s_{0.l}ws_{0.l}c, s_{0.r}ws_{0.r}c, s_{0.u}ws_{0.u}c, s_{1.l}ws_{1.l}c, s_{1.r}ws_{1.r}c, s_{1.u}ws_{1.u}c$
bigrams	$s_0ws_1w, s_0ws_1c, s_0cs_1w, s_0cs_1c, s_0wb_0w, s_0wb_0p, s_0cb_0w, s_0cb_0p, b_0wb_1w, b_0wb_1p, b_0pb_1w, b_0pb_1p, s_1wb_0w, s_1wb_0p, s_1cb_0w, s_1cb_0p$
trigrams	$s_0cs_1cs_2c, s_0ws_1cs_2c, s_0cs_1wb_0p, s_0cs_1cs_2w, s_0cs_1cb_0p, s_0ws_1cb_0p, s_0cs_1wb_0p, s_0cs_1cb_0w$

s_i : top node of the stack; b_i : front word on the buffer; xw : the word form of x ;
 xp : the POS tag; xc : the constituent label of a non-terminal node x ;
 xl, xr, xu : the left child, the right child and the unary child of x , respectively.

Transition-based constituent parsing

- Feature Templates for Shift-reduce Constituent Parser
- Example

In the 10-th step of the example sentence "*The little boy likes red tomatoes*", the $s_0 ps_0 c$ feature is $VBZ | VP$ and the $s_1 ps_1 c$ feature is $NN | NP$.

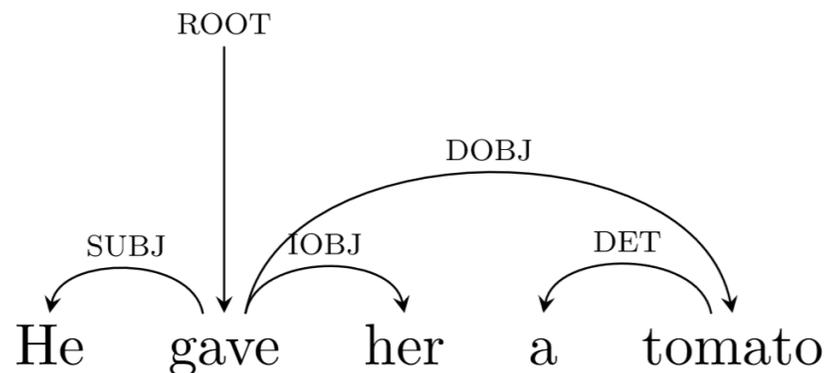


Contents

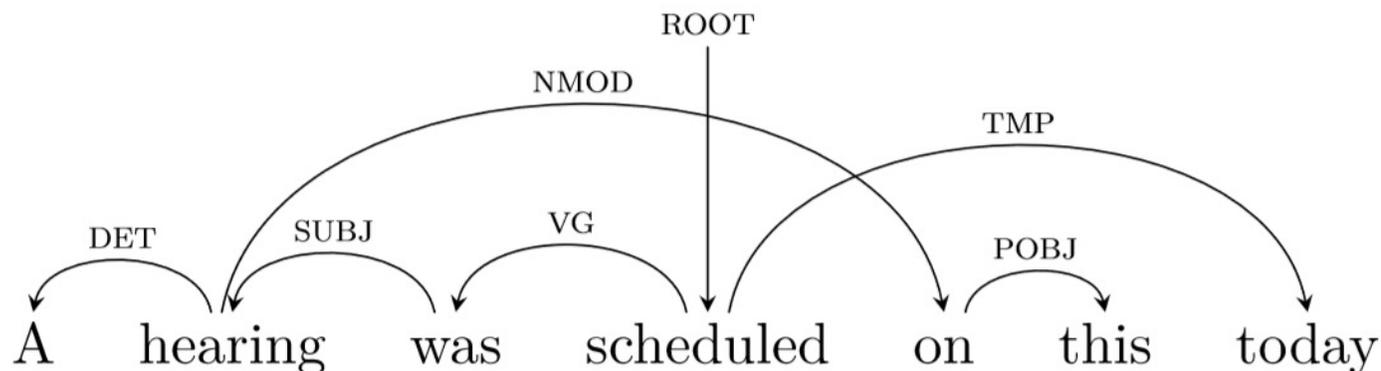
- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- **11.3 Shift-reduce Dependency Parsing**
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Shift-reduce dependency parsing

- Projective dependency tree



- Non-projective dependency tree



Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- **11.3 Shift-reduce Dependency Parsing**
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Shift-reduce dependency parsing

- Arc-standard Parsing

- **State**

σ : stack of partially dependency outputs.

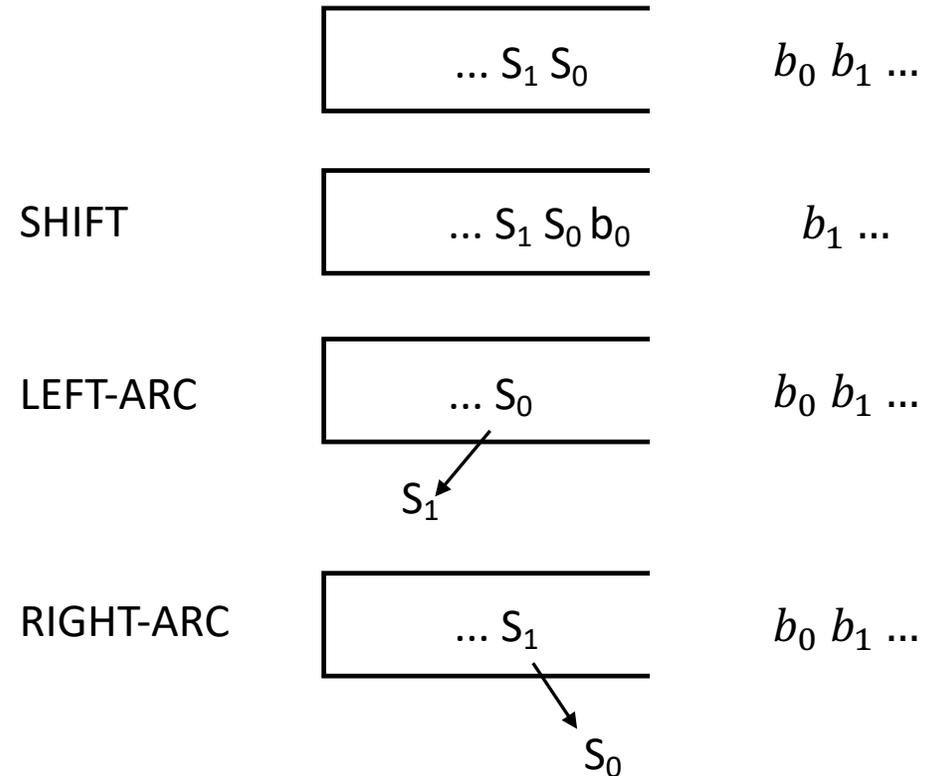
β : buffer of the next incoming words.

- **actions**

SHIFT, LEFT-ARC-X, RIGHT-ARC-X

Shift-reduce dependency parsing

Actions



Shift-reduce dependency parsing

Arc-standard Dependency Parsing

Axiom:	$([], W_{1:n}, \phi)$	LEFT-ARC-X:	$\frac{(\sigma s_1 s_0, \beta, A)}{(\sigma s_0, \beta, A \cup \{s_1 \overset{x}{\curvearrowright} s_0\})}$
Goal:	$([s_0], [], A)$		
SHIFT:	$\frac{(\sigma, b_0 \beta, A)}{(\sigma b_0, \beta, A)}$	RIGHT-ARC-X:	$\frac{(\sigma s_1 s_0, \beta, A)}{(\sigma s_1, \beta, A \cup \{s_1 \overset{x}{\curvearrowright} s_0\})}$

State: (σ, β, A) ;

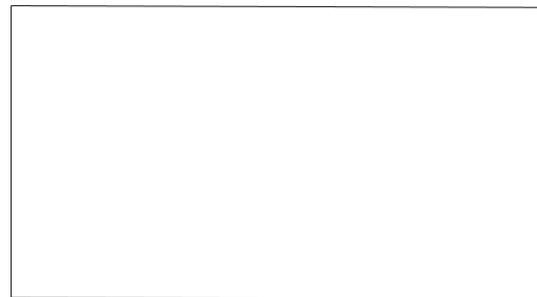
σ : stack; β : buffer; A : the set of dependency arcs that have been constructed

Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: SHIFT



He gave her a tomato

Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: SHIFT

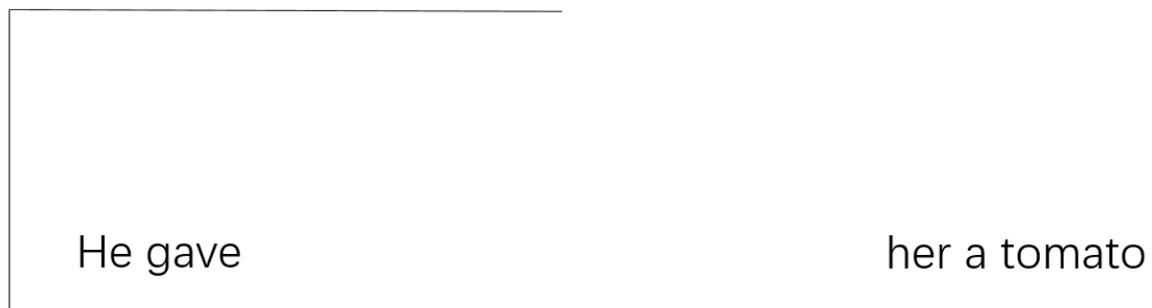


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: LEFT-ARC-SUBJ



Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: SHIFT

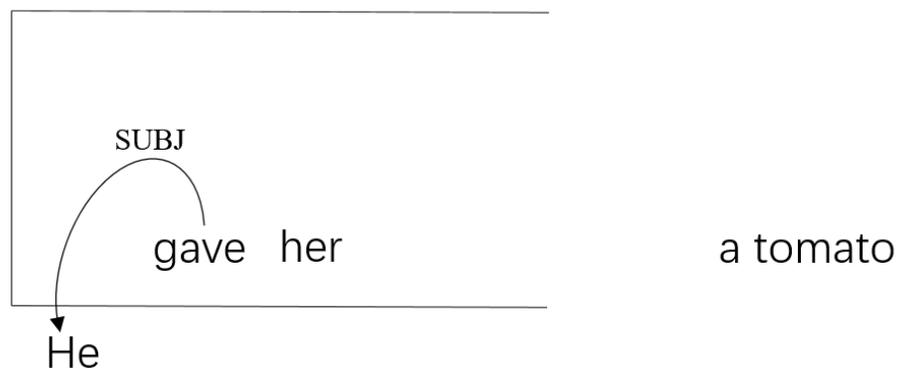


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: RIGHT-ARC-IOBJ

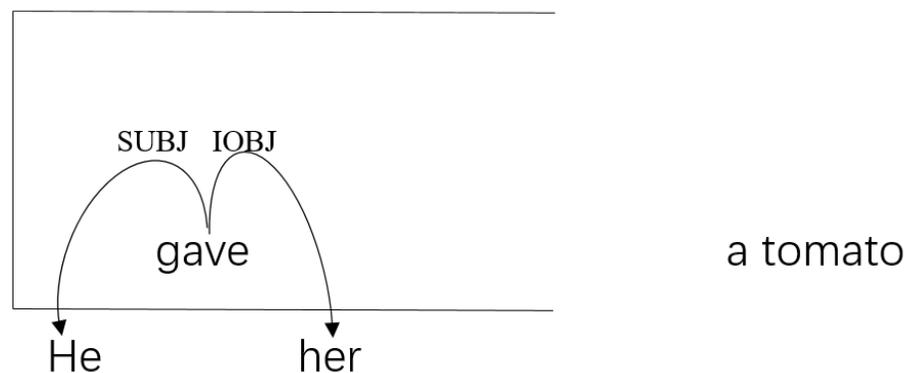


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: SHIFT

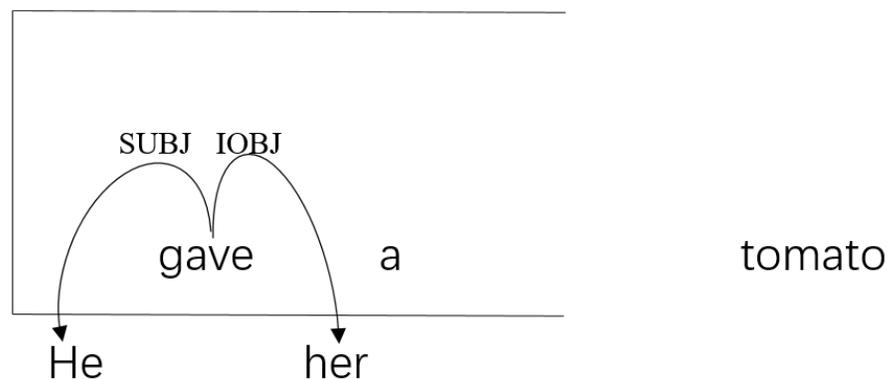


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: SHIFT

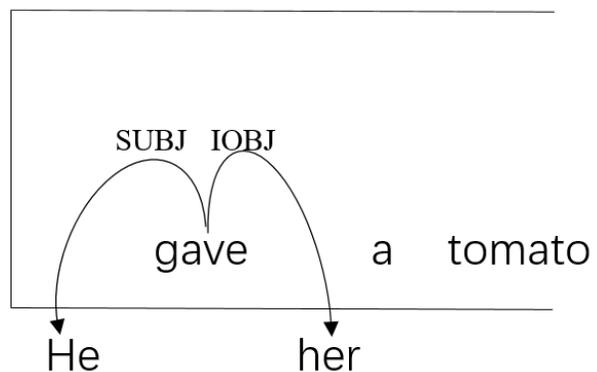


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: LEFT-ARC-DET

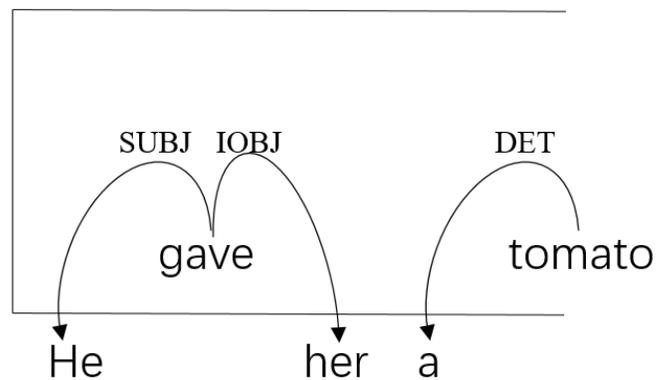


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

Next action: RIGHT-ARC-DOBJ

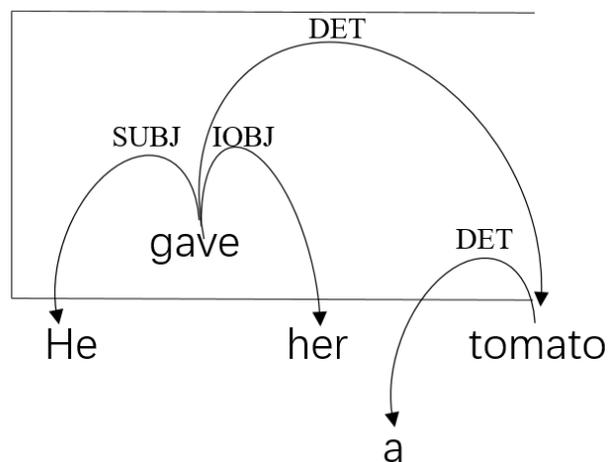


Arc-standard dependency parsing

Example

Sentence “He gave her a tomato”

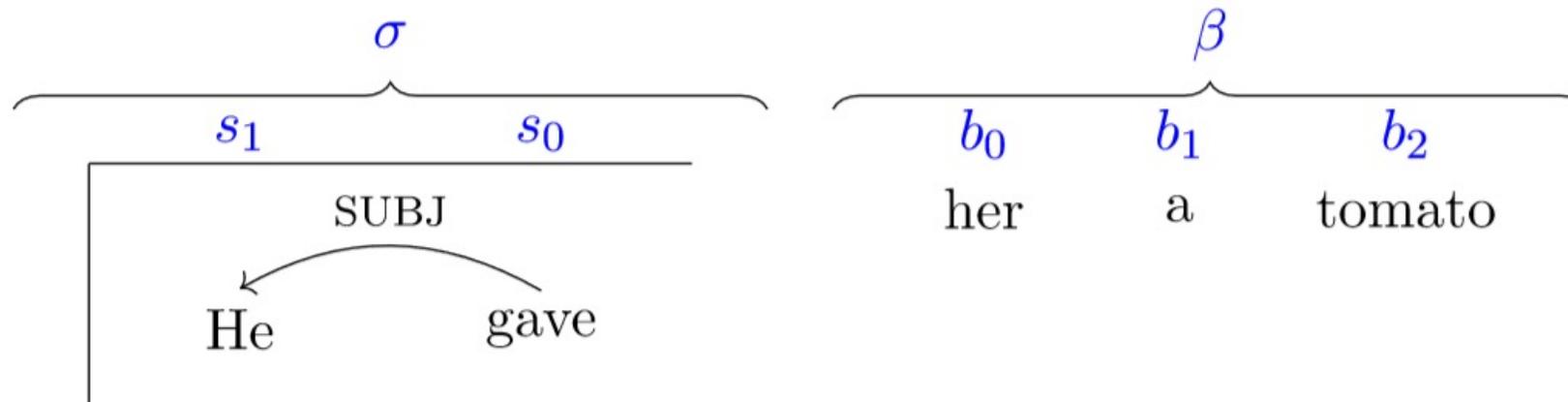
Next action: END



Shift-reduce dependency parsing

Arc-standard Dependency Parsing

- Example state



Feature Templates for arc-standard dependency parsing

Feature Type	Feature Template	Feature Type	Feature Template
from single words	$s_0wp; s_0w; s_0p; b_0wp;$ $b_0w; b_0p; b_1wp; b_1w;$ $b_1p; b_2wp; b_2w; b_2p;$	valency	$s_0wv_r; s_0pv_r; s_0wv_l;$ $s_0pv_l; b_0wv_l; b_0pv_l;$
from word pairs	$s_0wpb_0wp; s_0wpb_0w;$ $s_0wpb_0p; s_0wb_0wp;$ $s_0pb_0wp; s_0wb_0w;$ $s_0pb_0p; b_0pb_1p;$	unigrams	$s_{0.h}w; s_{0.h}p; s_{0.l};$ $s_{0.l}w; s_{0.l}p; s_{0.l}l;$ $s_{0.r}w; s_{0.r}p; s_{0.r}l;$ $b_{0.l}w; b_{0.l}p; b_{0.l}l;$
from three words	$b_0pb_1pb_2p; s_0pb_0pb_1p;$ $s_{0.hp}s_0pb_0p;$ $s_0ps_0.lpb_0p;$ $s_0ps_0.rpb_0p;$ $s_0pb_0pb_0.lp;$	third-order	$s_{0.h_2}w; s_{0.h_2}p; s_{0.hl};$ $s_{0.l_2}p; s_{0.l_2}l; s_{0.r_2}w;$ $s_{0.r_2}p; s_{0.r_2}l; b_{0.l_2}w;$ $b_{0.l_2}p; b_{0.l_2}l;$ $s_0ps_0.lps_0.l_2p;$ $s_0ps_0.rps_0.r_2p;$ $s_0ps_0.hps_0.h_2p;$ $b_0pb_0.lpb_0.l_2p;$
distance	$s_0wd; s_0pd; b_0wd;$ $b_0pd; s_0wb_0wd;$ $s_0pb_0pd;$	label set	$s_0ws_r; s_0ps_r; s_0ws_l;$ $s_0ps_l; n_0ws_l; n_0ps_l$

Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- **11.3 Shift-reduce Dependency Parsing**
 - 11.3.1 Arc-standard Dependency Parsing
 - **11.3.2 Arc-eager Projective Parsing**
 - 11.3.3 The Swap Action and Non-projective Trees
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Shift-reduce dependency parsing

- Arc-eager Parsing

- **State**

σ : stack of partially constituent outputs.

β : buffer of the next incoming words.

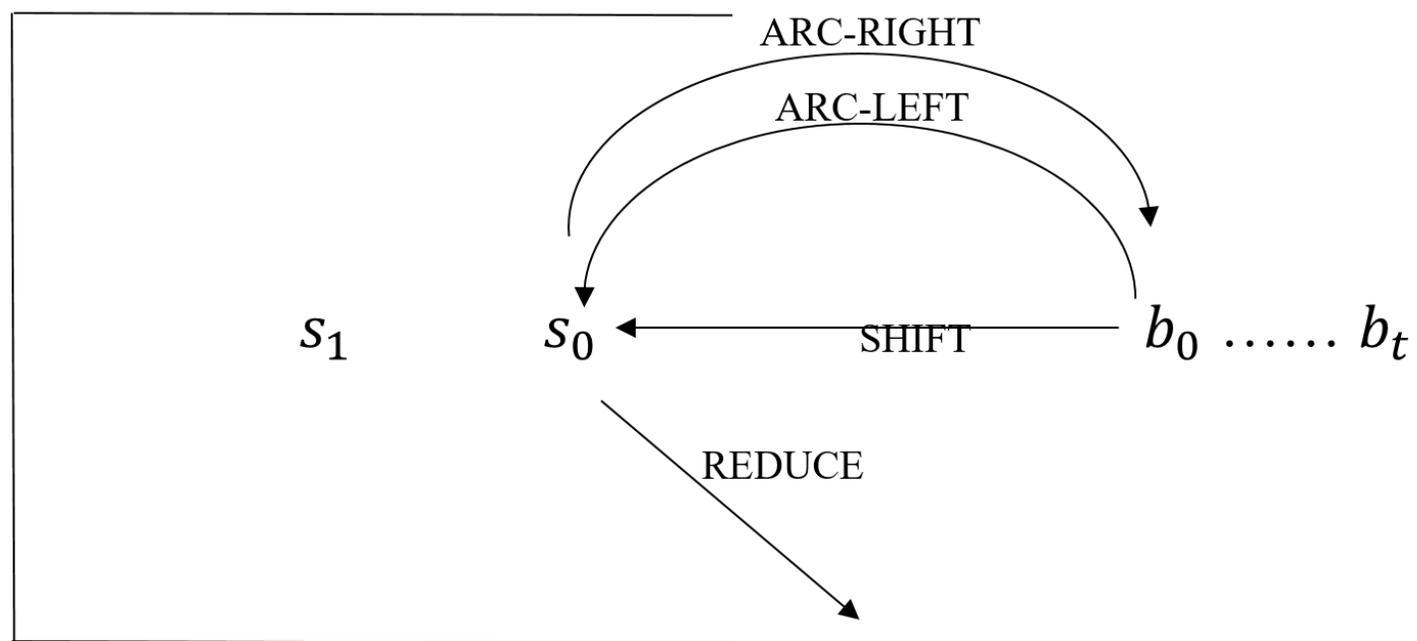
- **actions**

SHIFT, LEFT-ARC-X, RIGHT-ARC-X, REDUCE

Shift-reduce dependency parsing

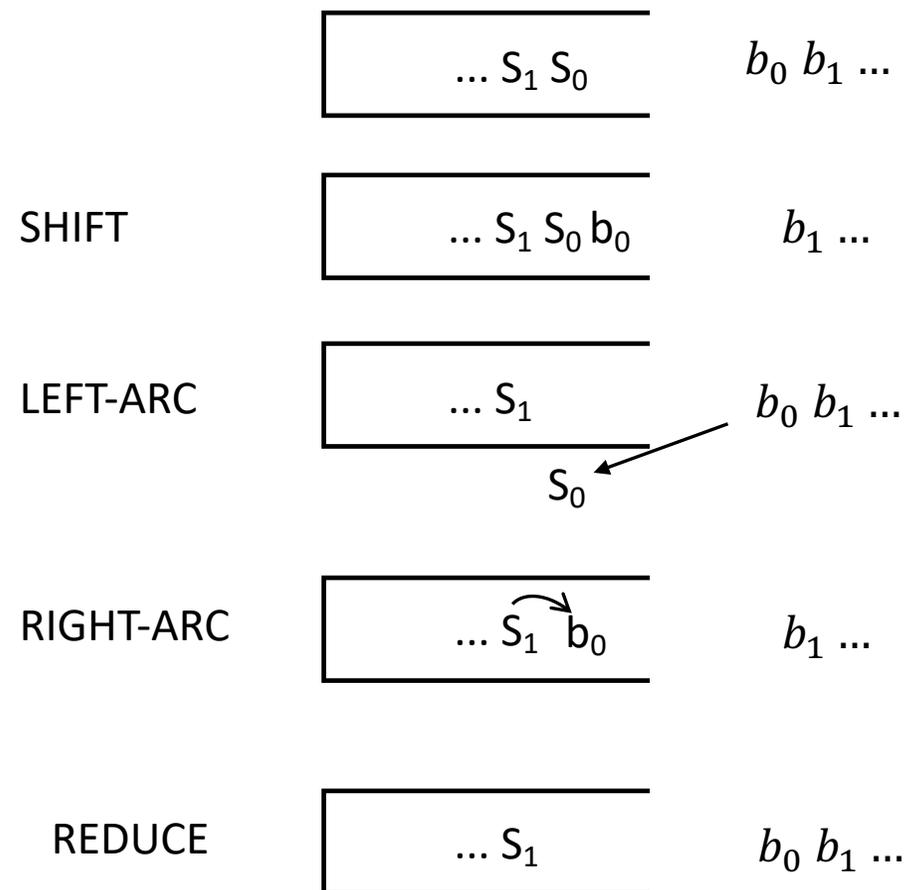
Arc-eager dependency parsing

- Example



Shift-reduce dependency parsing

Actions



Arc-eager dependency parsing

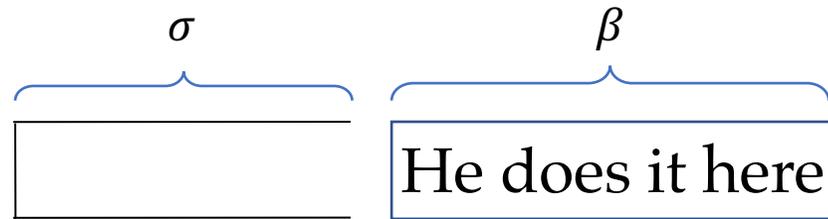
Axiom: $([], W_{1:n}, \phi)$	LEFT-ARC-X: $\frac{(\sigma s_0, b_0 \beta, A), \text{ such that } \neg(\exists(k, L) w_k \overset{L}{\curvearrowright} s_0 \in A)}{(\sigma, b_0 \beta, A \cup \{s_0 \overset{X}{\curvearrowright} b_0\})}$
Goal: $([s_0], [], A)$	
SHIFT: $\frac{(\sigma, b_0 \beta, \phi)}{(\sigma b_0, \beta, \phi)}$	RIGHT-ARC-X: $\frac{(\sigma s_0, b_0 \beta, A)}{(\sigma s_0 b_0, \beta, A \cup \{s_0 \overset{X}{\curvearrowright} b_0\})}$
	REDUCE: $\frac{(\sigma s_0, \beta, A), \text{ such that } (\exists(k, L) w_k \overset{L}{\curvearrowright} s_0 \in A)}{(\sigma, \beta, A)}$

Main differences comparing with arc-standard:

- *Left – Arc – X/Right – Arc – X*: construct a dependency arc from the front word on the buffer to the top word on the stack, but not as that (from the top two words on the stack) in arc-standard.
- *Reduce*: pop the top word off the stack.

Example (Dependency Parsing)

- NEXT ACTION: *Shift*

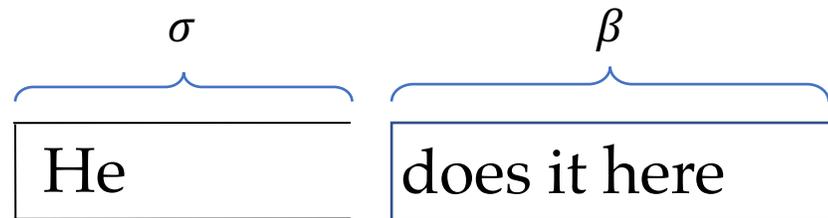


σ : stack

β : buffer

Example (Dependency Parsing)

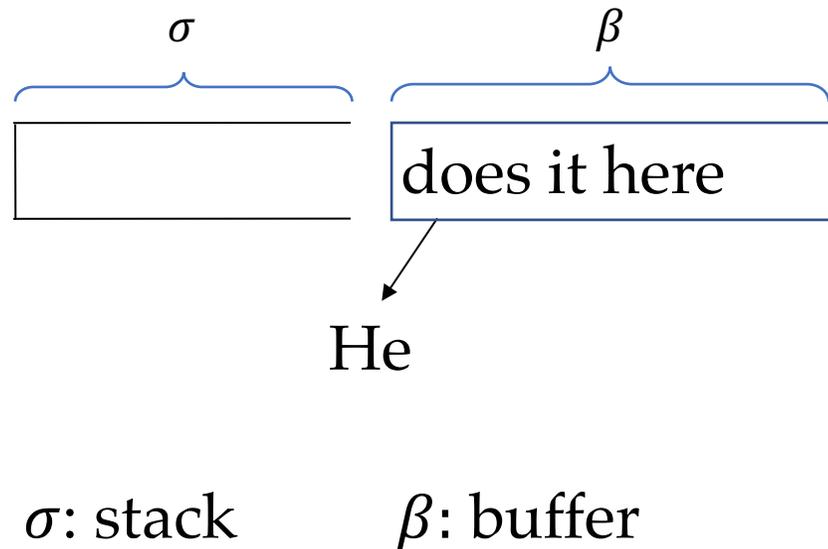
- NEXT ACTION: *Left – Arc – SUBJ*



σ : stack β : buffer

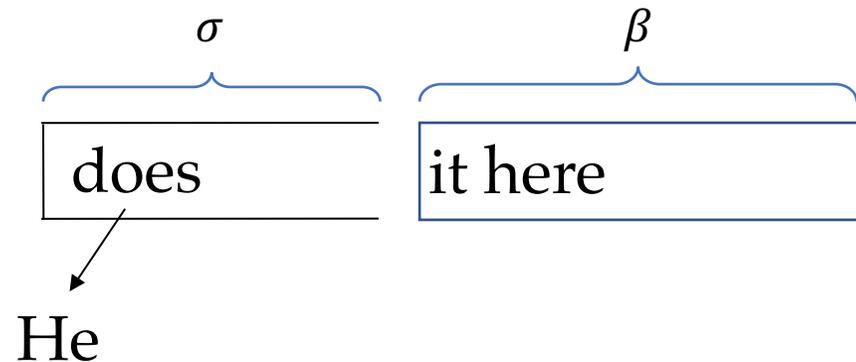
Example (Dependency Parsing)

- NEXT ACTION: *Shift*



Example (Dependency Parsing)

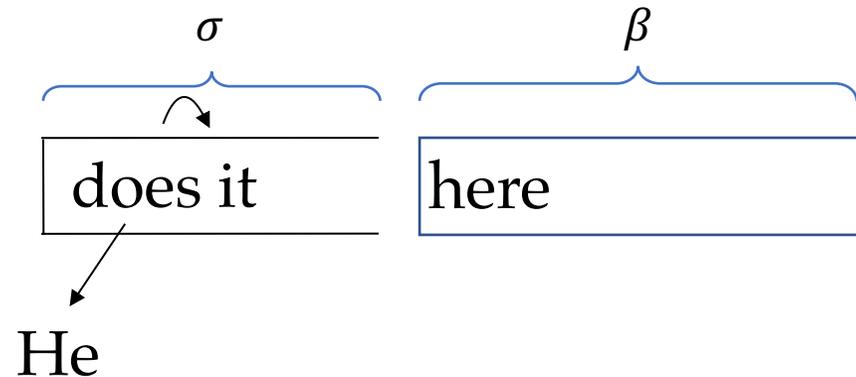
- NEXT ACTION: *Right – Arc – OBJ*



σ : stack β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Reduce*

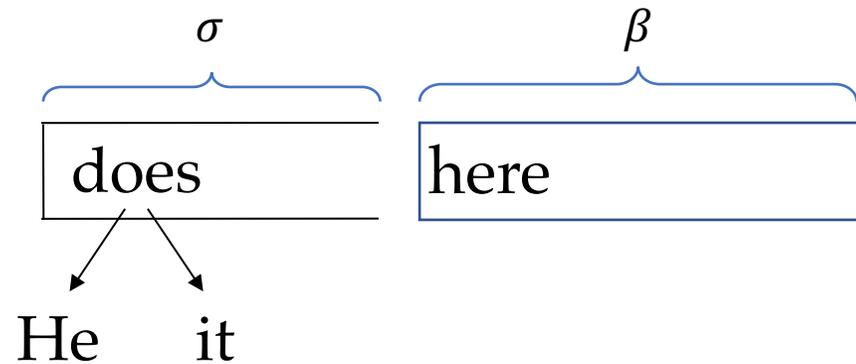


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Right – Arc – MOD*

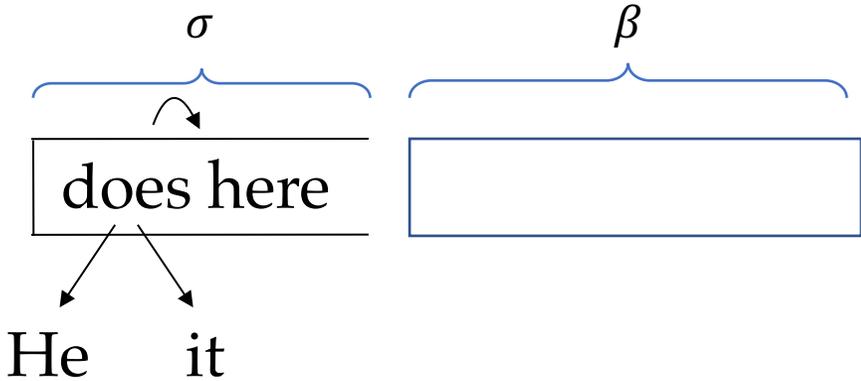


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Reduce*

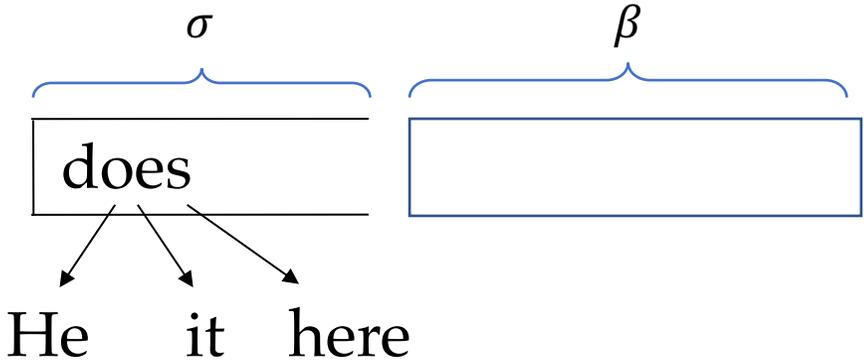


σ : stack

β : buffer

Example (Dependency Parsing)

- NEXT ACTION: *Finish*



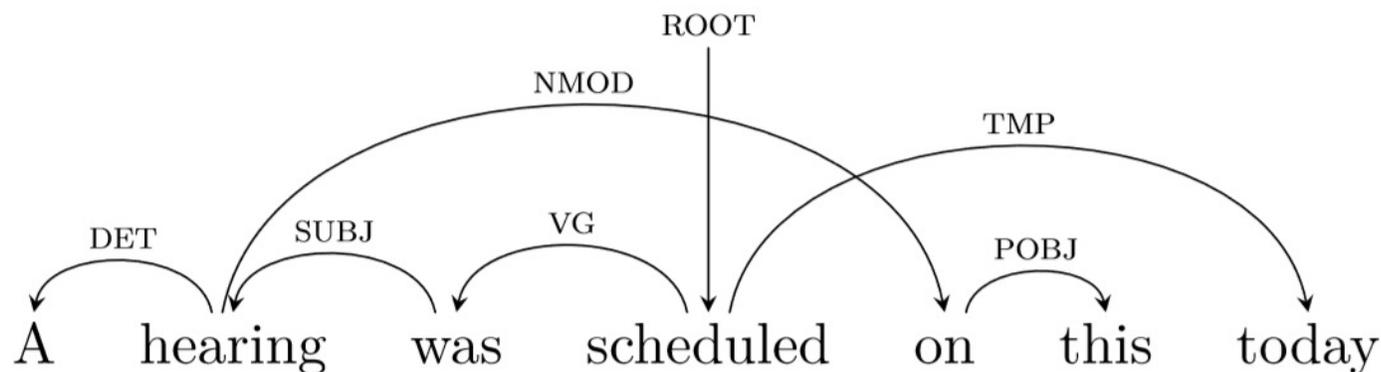
σ : stack β : buffer

Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- **11.3 Shift-reduce Dependency Parsing**
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - **11.3.3 The Swap Action and Non-projective Trees**
- 11.4 Joint Models
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

Shift-reduce dependency parsing

- **The Swap Action and Non-projective Trees**
- To allow constructing non-projective trees, the arc-standard system can be extended by adding a new action:
 - *Swap*: remove the second top word from the stack, pushing it onto the buffer front.



Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

Axiom: $([], W_{1:n}, \phi)$	LEFT-ARC-X: $\frac{([\sigma s_1 s_0], \beta, A)}{([\sigma s_0], \beta, A \cup \{s_1 \overset{x}{\curvearrowright} s_0\})}$
Goal: $([s_0], [], A)$	
SHIFT: $(\sigma, [b_0 \beta], \phi)$	RIGHT-ARC-X: $\frac{([\sigma s_1 s_0], \beta, A)}{([\sigma s_1], \beta, A \cup \{s_1 \overset{x}{\curvearrowright} s_0\})}$
$([\sigma b_0], \beta, \phi)$	
	SWAP: $\frac{([\sigma s_1 s_0], \beta, A), \text{ such that } \text{IDX}(s_1) < \text{IDX}(s_0)}{([\sigma s_0], [s_1 \beta], A)}$

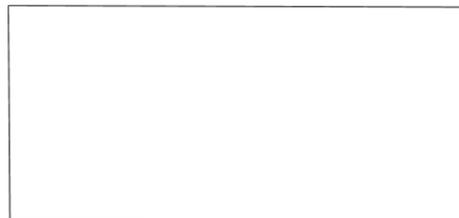
State: (σ, β, A) ; σ : stack; β : buffer; A : the set of dependency arcs that have been constructed; $\text{IDX}(w)$: return the index of w in the sentence $W_{1:n}$.

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



A hearing was scheduled on this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT

A

hearing was scheduled on this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: LEFT-ARC-DET

A hearing

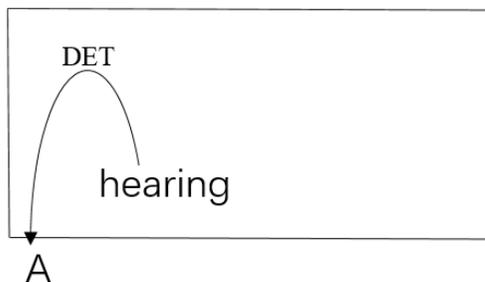
was scheduled on this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



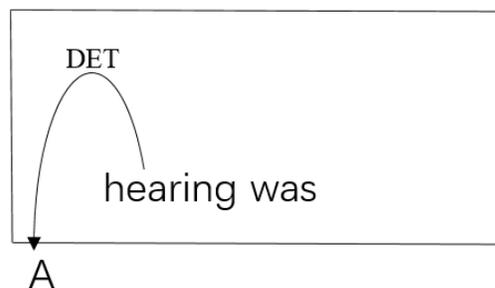
was scheduled on this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



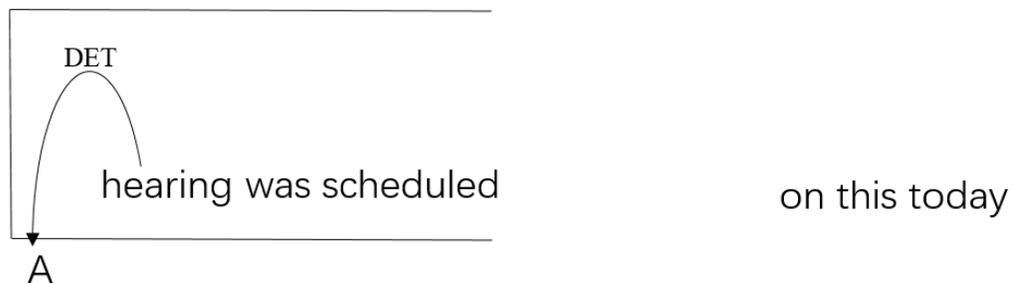
scheduled on this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT

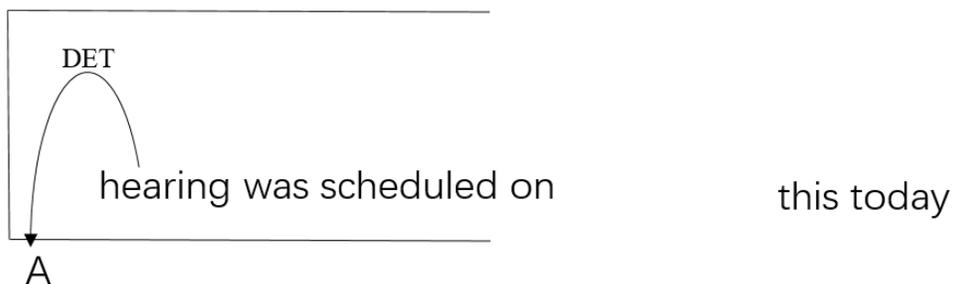


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SWAP

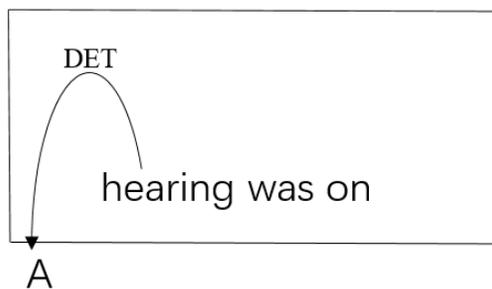


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SWAP



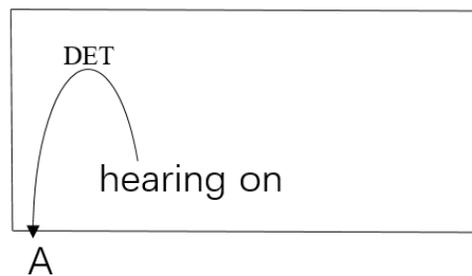
scheduled this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



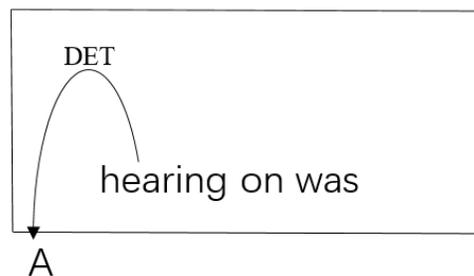
was scheduled this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



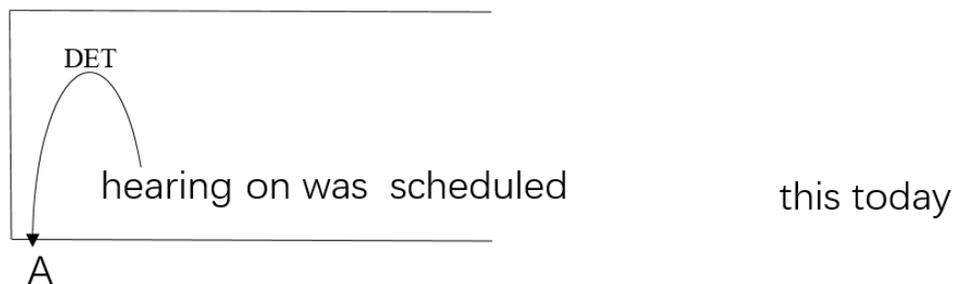
scheduled this today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT

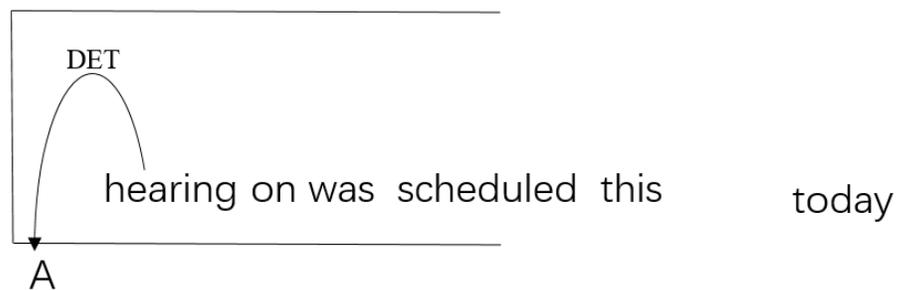


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SWAP

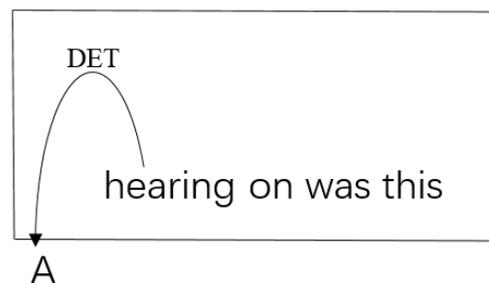


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SWAP



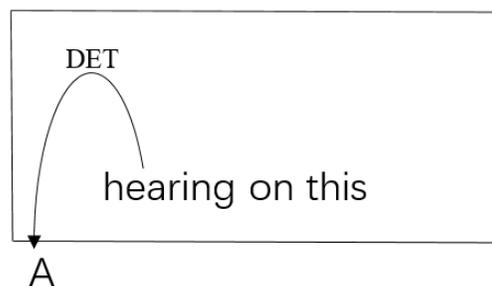
scheduled today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: RIGHT-ARC-POBJ



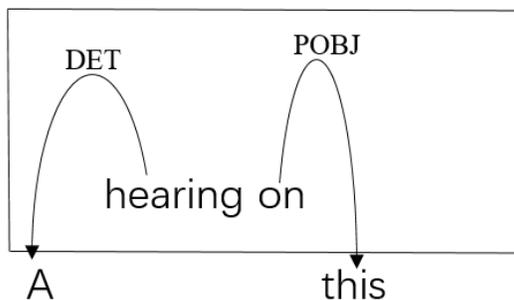
was scheduled today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: RIGHT-ARC-NMOD



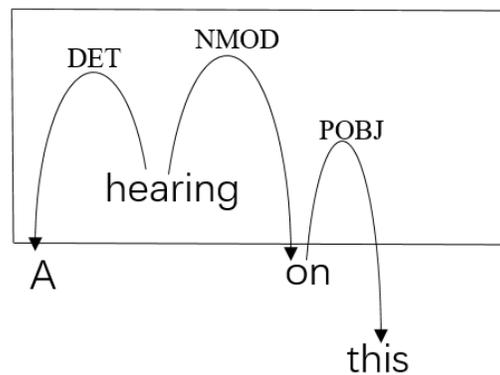
was scheduled today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT



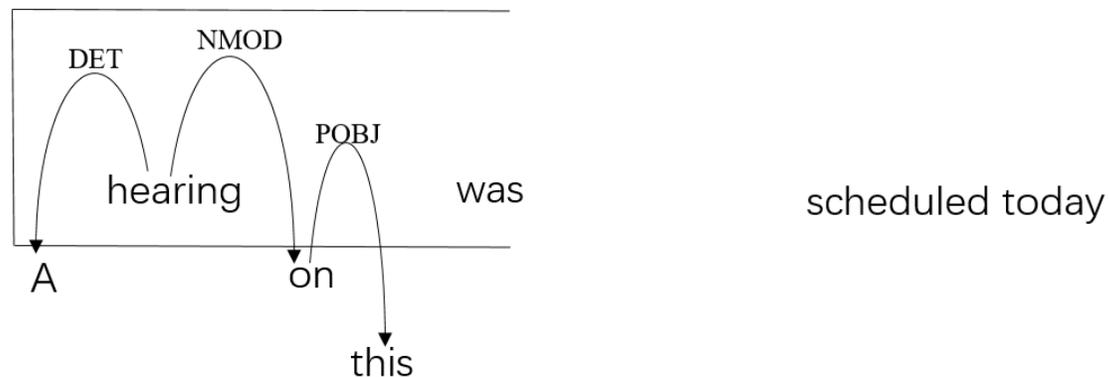
was scheduled today

Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: LEFT-ARC-SUBJ

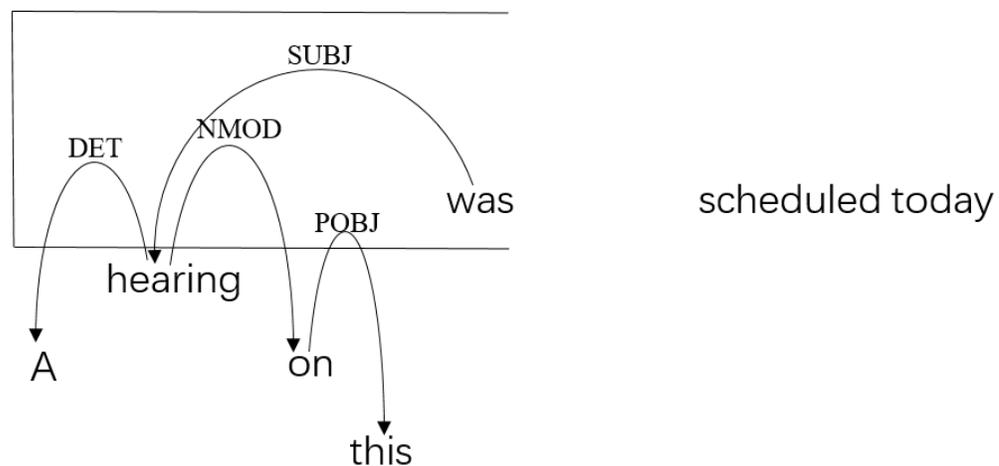


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT

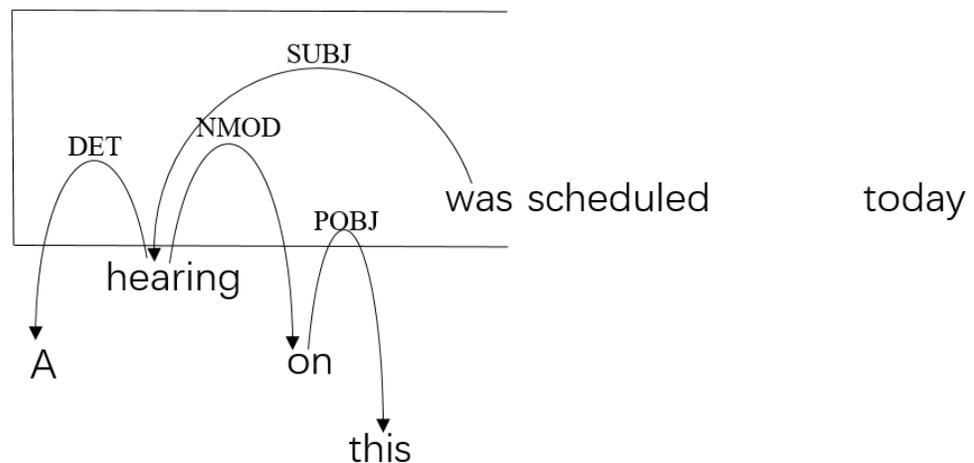


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: SHIFT

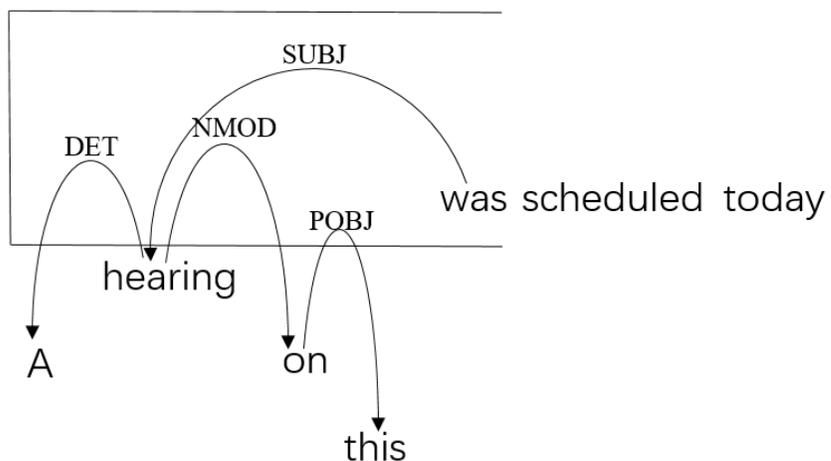


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: RIGHT-ARC-TMP

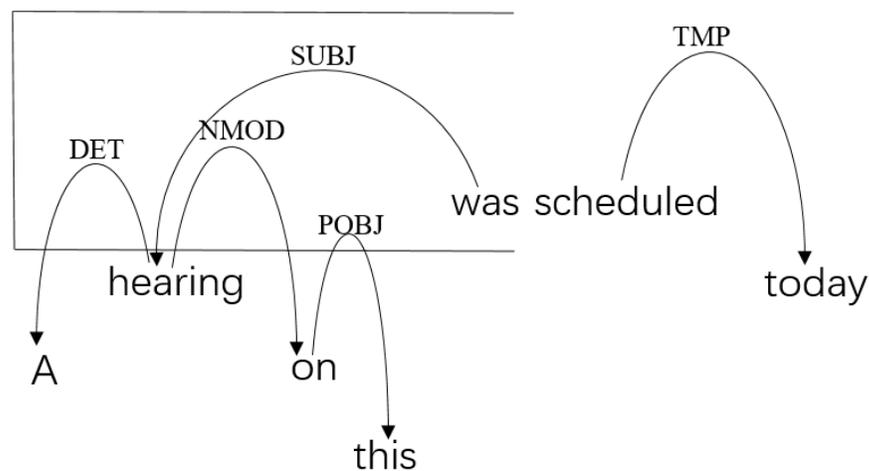


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: LEFT-ARC-VG

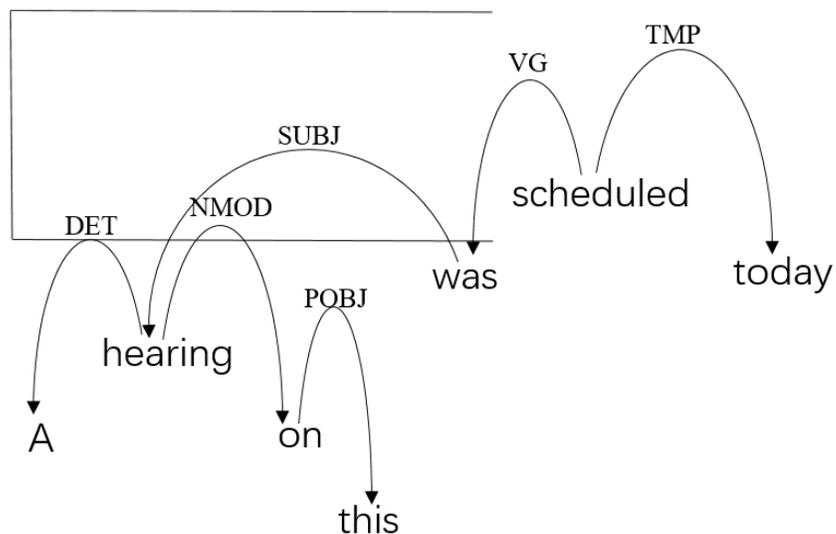


Shift-reduce dependency parsing

Arc-standard Dependency Parsing with Swap Action

- Example

Next action: END



Contents

- 11.1 Transition-based Structured Prediction
 - 11.1.1 Greedy Local Modelling
 - 11.1.2 Structured Modelling
- 11.2 Transition-based Constituent Parsing
 - 11.2.1 Shift-reduce Constituent Parsing
 - 11.2.2 Feature Templates
- 11.3 Shift-reduce Dependency Parsing
 - 11.3.1 Arc-standard Dependency Parsing
 - 11.3.2 Arc-eager Projective Parsing
 - 11.3.3 The Swap Action and Non-projective Trees
- **11.4 Joint Models**
 - 11.4.1 Joint POS-tagging and Dependency Parsing
 - 11.4.2 Joint Word Segmentation, POS-tagging and Dependency Parsing

- Motivations
 - Cross-task information sharing
 - Reduction of error propagation
- Example
 - For joint POS-tagging and syntactic parsing, the *SHIFT* action of arc-standard algorithm can be replaced with *SHIFT – X* action, where *X* refers to the POS label.

Joint POS-tagging and Dependency Parsing

- Input: $W_{1:n} = w_1, \dots, w_n$
- Output: $(T_{1:n}, A)$, where t_i is POS for w_i .
- The arc-standard algorithm can be extended by replacing *SHIFT* action with *SHIFT – X* action
 - *SHIFT – X*, which removes the front word from the buffer, assigning the POS label X to the word, and pushing it onto the stack.

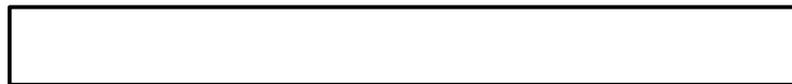
Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: SHIFT—PRP

Stack [S]



Buffer [B]

He₁ won₂ the₃ game₄

Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: SHIFT——VBD

Stack [S]



Buffer [B]

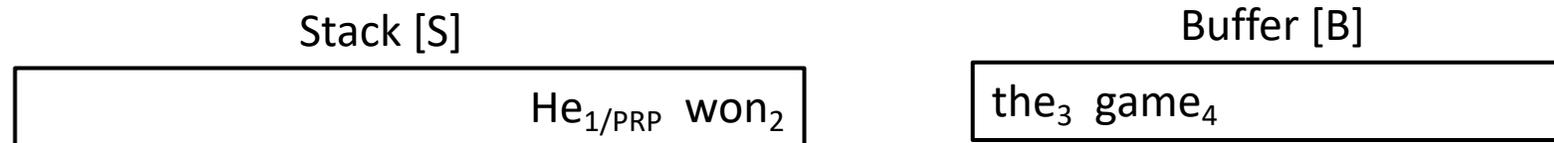


Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: LEFT—ARC—NSUBJ

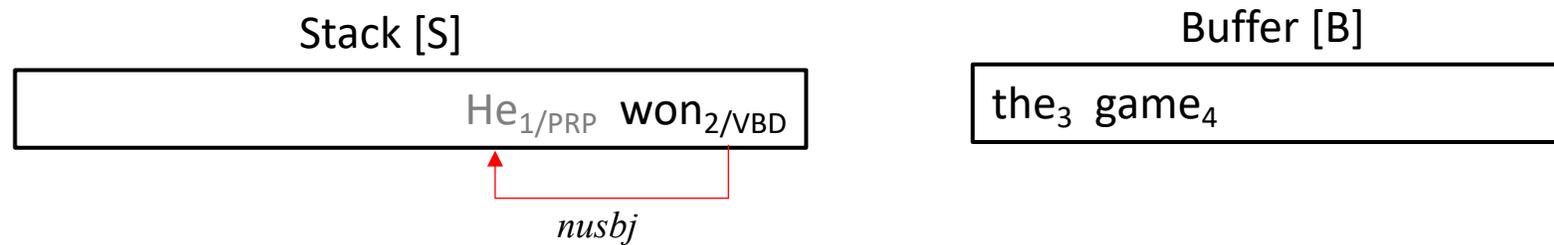


Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: SHIFT—DT

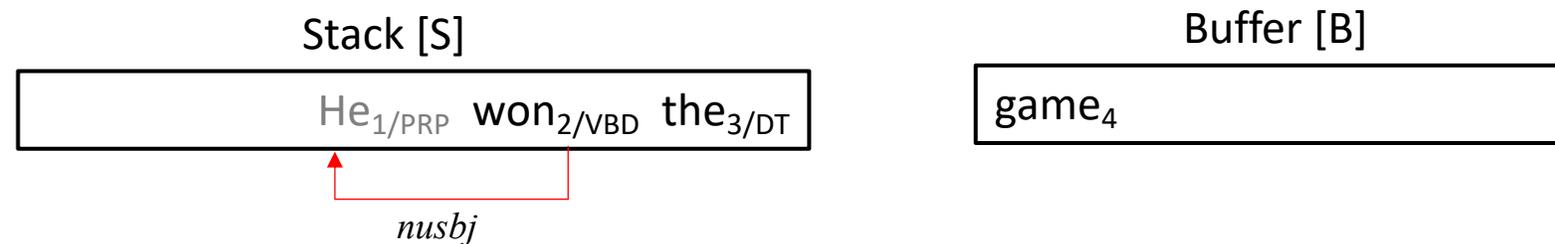


Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: SHIFT—NN

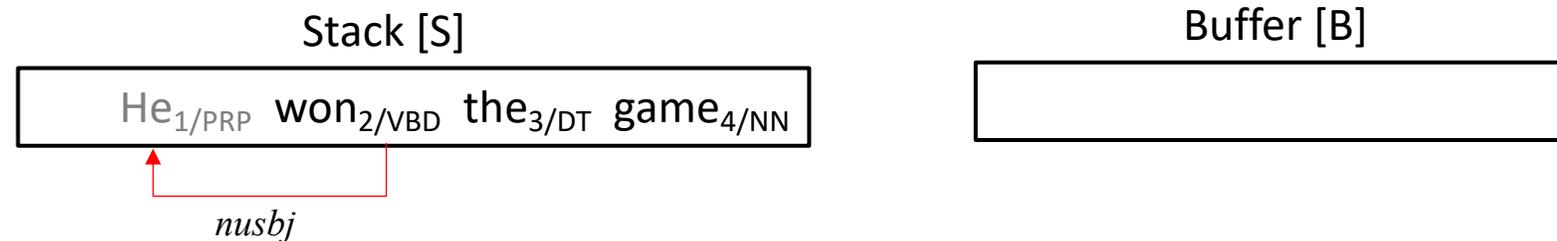


Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: LEFT—ARC—DET



Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: RIGHT—ARC—DOBJ

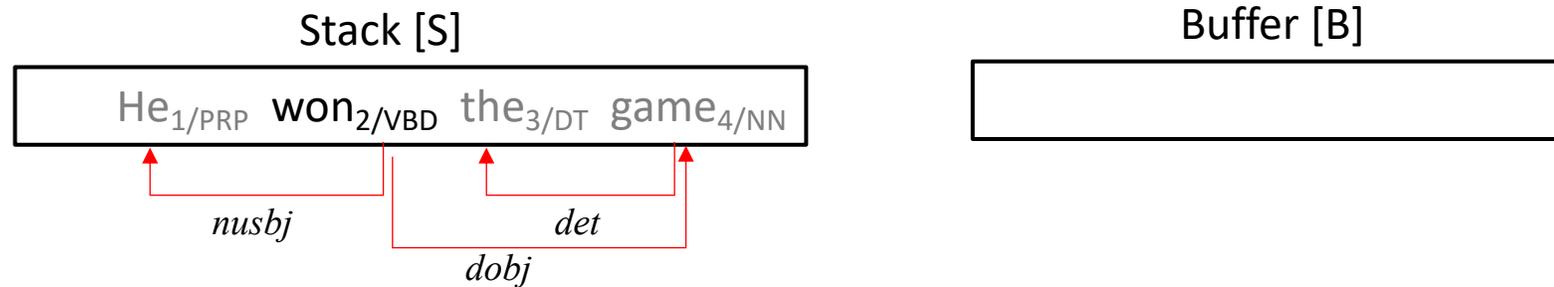


Joint Models

Joint POS-tagging and Dependency Parsing

- Example

Next action: END



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Input: $C_{1:n} = C_1 \dots C_n$
- Output: $(W_{1:m}, T_{1:m}, A)$ where $w_1 \dots w_m$ are words, $t_1 \dots t_n$ are POS, and A is the set of dependency arcs.
- State:
 - σ -- partially built outputs
 - δ -- words
 - β -- incoming characters
 - A_c -- character dependencies
 - A_w -- word dependencies

Joint Word Segmentation, POS-tagging and Dependency Parsing

Axiom:	$([], [], C_{1:n}, \phi, \phi)$
Goal:	$([S_0], [], [], A_c, A_w)$
LEFT-ARC-C:	$\frac{(\sigma, \delta d_0, b_0 \beta, A_c, A_w) \text{ such that } \neg(\exists d \in \delta, d \curvearrowright d_0 \in A_c)}{(\sigma, \delta, b_0 \beta, A_c \cup \{d_0 \curvearrowright b_0\}, A_w)}$
LEFT-ARC-X:	$\frac{(\sigma s_0, [d_0], \beta, A_c, A_w) \text{ such that } \neg(\exists s \in \sigma, s \curvearrowright^l s_0 \in A_w)}{(\sigma, [d_0], \beta, A_c, A_w \cup \{s_0 \curvearrowright d_0\})}$
SHIFT:	$\frac{(\sigma, [d_0], \beta, A_c, A_w)}{(\sigma d_0, [], \beta, A_c, A_w)}$
SHIFT-C:	$\frac{(\sigma, \delta, b_0 \beta, A_c, A_w)}{(\sigma, \delta b_0, \beta, A_c, A_w)}$
RIGHT-ARC-C:	$\frac{(\sigma, \delta d_0, b_0 \beta, A_c, A_w)}{(\sigma, \delta d_0 b_0, \beta, A_c \cup \{d_0 \curvearrowright b_0\}, A_w)}$
RIGHT-ARC-X:	$\frac{(\sigma s_0, [d_0], \beta, A_c, A_w \cup \{s_0 \curvearrowright d_0\})}{(\sigma s_0 d_0, [], \beta, A_c, A_w)}$
POP-X:	$\frac{(\sigma, [d_0], \beta, A_c, A_w)}{(\sigma, [\text{SUBTREE}(d_0, A_c)/X], \beta, A_c, A_w)}$
REDUCE-C:	$\frac{(\sigma, \delta d_0, \beta, A_c, A_w) \text{ such that } \exists d \in \delta, d \curvearrowright d_0 \in A_c}{(\sigma, \delta, \beta, A_c, A_w)}$
REDUCE:	$\frac{(\sigma s_0, \delta, \beta, A_c, A_w) \text{ such that } \exists s \in \sigma, s \curvearrowright s_0 \in A_w}{(\sigma, \delta, \beta, A_c, A_w)}$

State: $(\sigma, \delta, \beta, A_c, A_w)$ σ : stack; β :buffer; δ : partial-word buffer;

A_c : the set of character dependencies; A_w : the set of word dependencies ¹⁶⁴

Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
0			我, 来, 到, 会, 客, 室			<i>SHIFT-C</i>

--	--

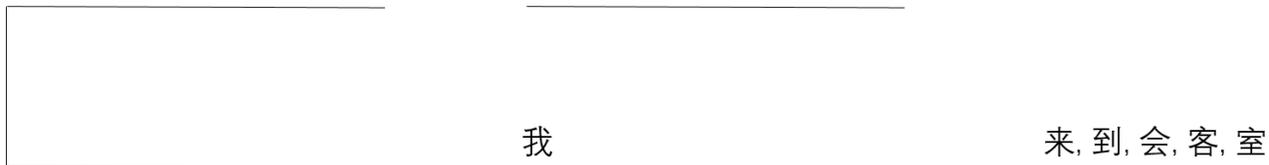
我, 来, 到, 会, 客, 室

Joint Models

Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
1		我	来, 到, 会, 客, 室			<i>POP-PN</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
2		我/PN	来, 到, 会, 客, 室			<i>SHIFT-W</i>



我/PN

来, 到, 会, 客, 室

Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
3	我/PN		来, 到, 会, 客, 室			<i>SHIFT-C</i>

我/PN

来, 到, 会, 客, 室

来, 到, 会, 客, 室

Joint Models

Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
4	我/PN	来	到, 会, 客, 室			<i>LEFTARC-C</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

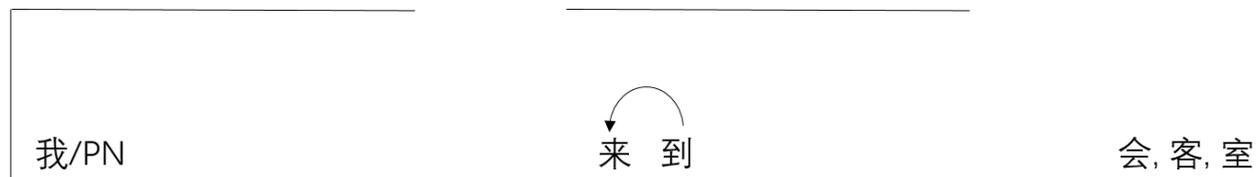
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
5	我/PN		到, 会, 客, 室	来←到		<i>SHIFT</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

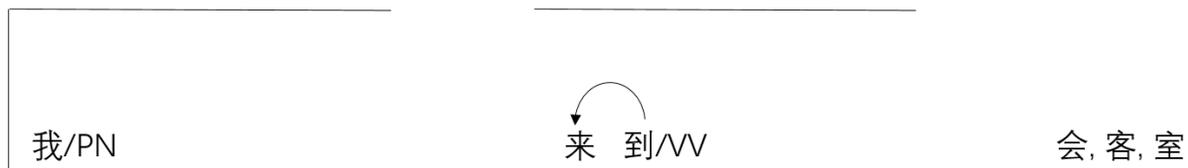
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
6	我/PN	到	到, 会, 客, 室	来←到		<i>POP-VV</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

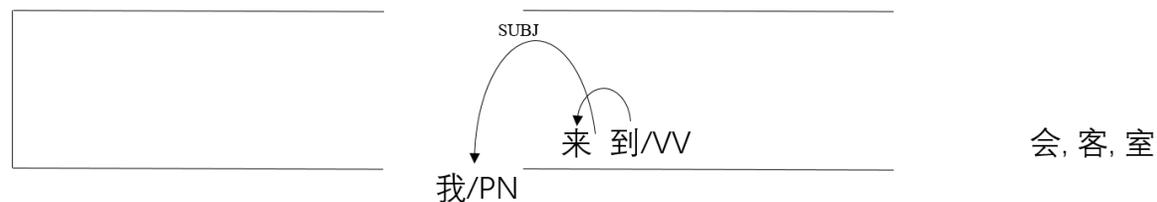
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
7	我/PN	来到/VV	会, 客, 室	来←到		<i>LEFTARC-SUBJ</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

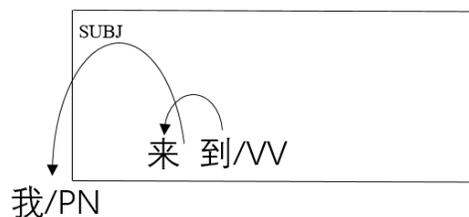
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
8		来到/VV	会, 客, 室	来←到	我/PN← <i>SUBJ</i> 来到/VV	<i>SHIFT</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
9	来到/VV		会, 客, 室	来←到	我/PN← <i>SUBJ</i> 来到/VV	<i>SHIFT-C</i>



会, 客, 室

Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

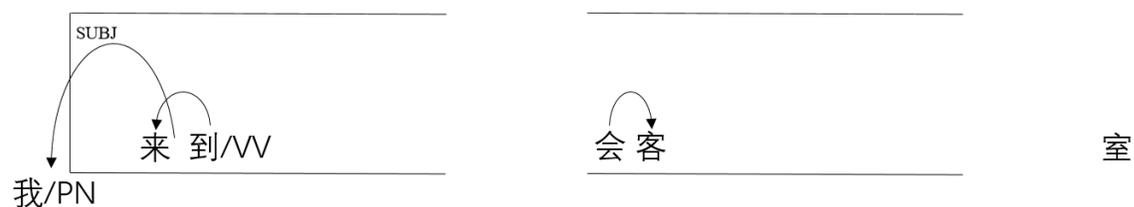
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
10	来到/VV	会	客, 室	来←到	我/PN← <i>SUBJ</i> 来到/VV	<i>RIGHTARC-C</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

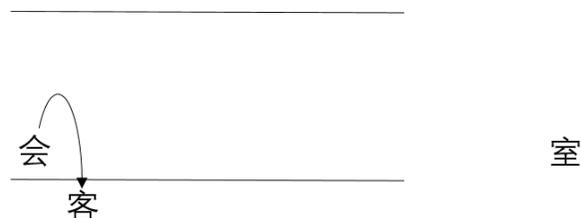
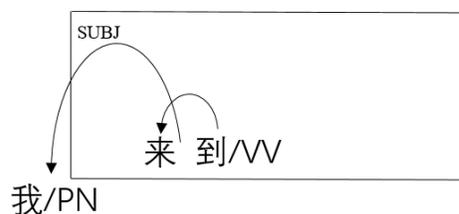
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
11	来到/VV	会, 客	室	来←到, 会→客	我/PN← <i>SUBJ</i> 来到/VV	<i>REDUCE-C</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

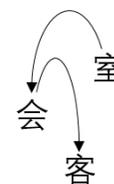
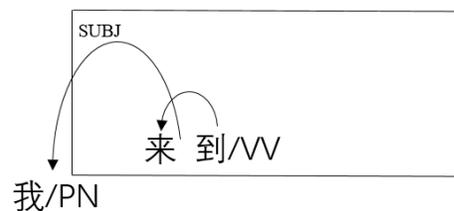
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
12	来到/VV	会	室	来←到, 会→客	我/PN← <i>SUBJ</i> 来到/VV	<i>LEFTARC-C</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

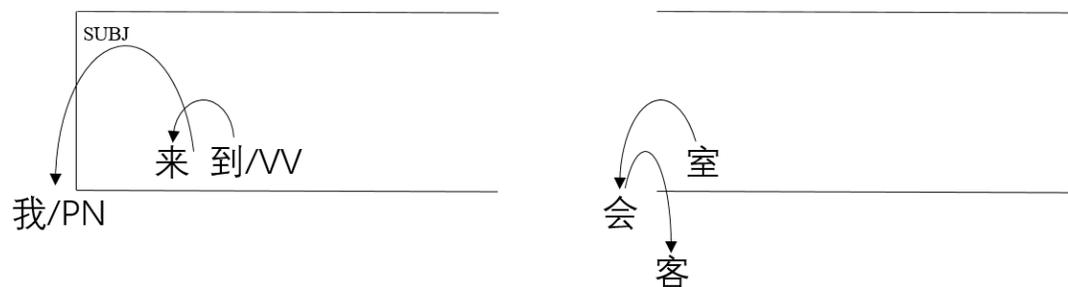
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
13	来到/VV		室	来←到, 会→客, 会←室	我/PN← <i>SUBJ</i> 来到/VV	<i>SHIFT-C</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

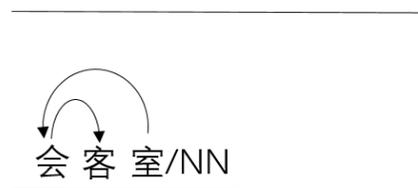
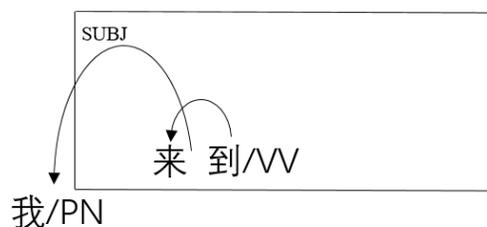
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
14	来到/VV	室		来←到, 会→客, 会←室	我/PN← <i>SUBJ</i> 来到/VV	<i>POP-NN</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

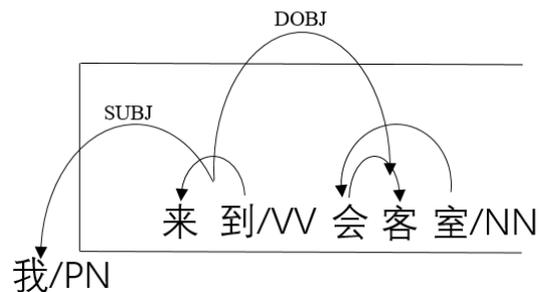
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
15	来 到/VV	会客 室/NN		来←到, 会→客, 会 ←室	我/PN← <i>SUBJ</i> 来 到/VV	<i>RIGHTARC-DOBJ</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

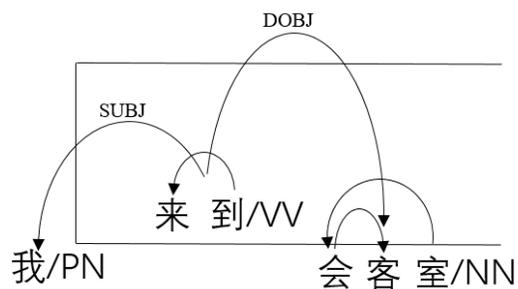
<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
16	来到/VV, 会客室/NN			来← 到, 会→ 客, 会←室	我/PN← <i>SUBJ</i> 来到/VV, 来到/VV <i>DOBJ</i> →会客室/NN	<i>REDUCE</i>



Joint Word Segmentation, POS-tagging and Dependency Parsing

- Example

<i>Step</i>	σ	δ	β	A_c	A_w	<i>Action</i>
17	来 到/VV			来←到, 会→客, 会 ←室	我/PN← <i>SUBJ</i> 来到/VV, 来到/VV <i>DOBJ</i> →会客 室/NN	<i>END</i>



Summary

- What is transition-based method?
- Apply transition-based methods on different tasks.
- Joint modeling with transition-based methods.