

# Statistical Machine Translation

-tree-based models (cont.)-

**Artem Sokolov**

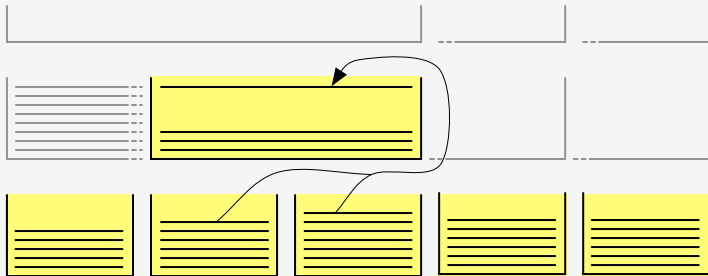
Computerlinguistik

Universität Heidelberg

Sommersemester 2015

material from P. Koehn, S. Riezler, D. Altshuler

- For each span, a stack of (partial) translations is maintained
- Bottom-up: a higher stack is filled, once underlying stacks are complete

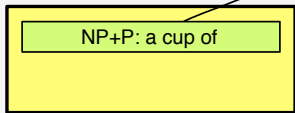
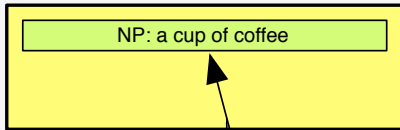


**Input:** Foreign sentence  $\mathbf{f} = f_1, \dots, f_{l_f}$ , with syntax tree

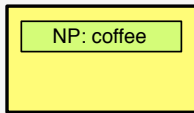
**Output:** English translation  $\mathbf{e}$

```
1: for all spans [start,end] (bottom up) do
2:   for all sequences  $s$  of hypotheses and words in span [start,end] do
3:     for all rules  $r$  do
4:       if rule  $r$  applies to chart sequence  $s$  then
5:         create new hypothesis  $c$ 
6:         add hypothesis  $c$  to chart
7:       end if
8:     end for
9:   end for
10: end for
11: return English translation  $\mathbf{e}$  from best hypothesis in span  $[0, l_f]$ 
```

Applying rule creates new hypothesis



apply rule:  
 $NP \rightarrow NP \text{ Kaffee}$  ;  $NP \rightarrow NP+P \text{ coffee}$

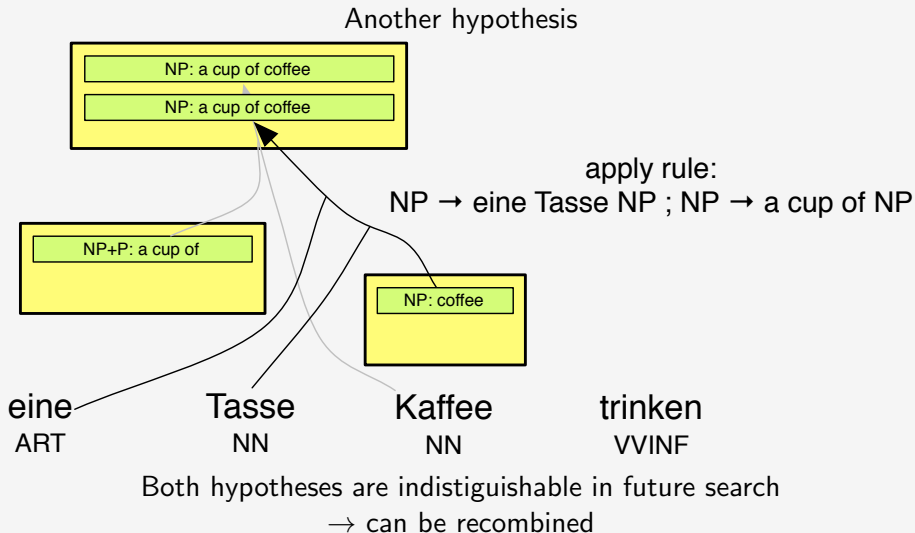


eine  
ART

Tasse  
NN

Kaffee  
NN

trinken  
VVINF

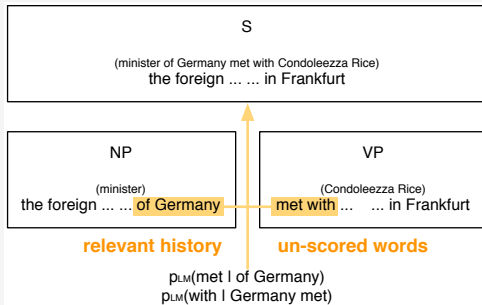


Hypotheses have to match in

- span of input words covered
- output constituent label
- first  $n-1$  output words (not properly scored, since they lack context)
- last  $n-1$  output words (still affect scoring of subsequently added words, just like in phrase-based decoding)

( $n$  is the order of the  $n$ -gram language model)

When merging hypotheses, internal language model contexts are absorbed



## Search space pruning

- recombination
- stack pruning

## Algorithmic techniques

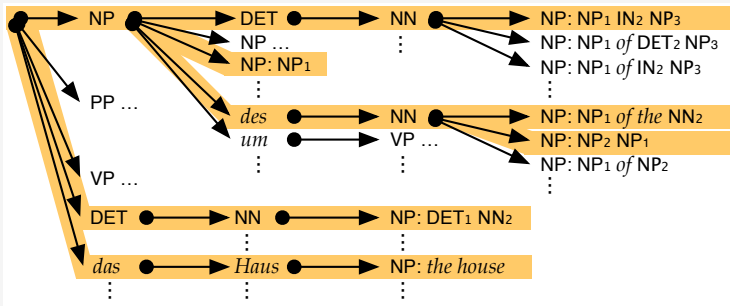
- prefix tree
- Earley's parsing algorithm
- cube pruning

- Number of hypotheses in each chart cell explodes
- ⇒ need to discard bad hypotheses  
e.g., keep  $n = 100$  best only
- Different stacks for different output constituent labels
- ⇒ keep at least  $m$  different ( $m=2,3,..$ )



- Many subspan sequences
    - for all** sequences  $s$  of hypotheses and words in span [start,end]
  - Many rules
    - for all** rules  $r$
  - Checking if a rule applies not trivial
    - rule  $r$  applies to chart sequence  $s$
- ⇒ Unworkable

- Easy:
  - ➔ given a rule
  - ➔ check if and how it can be applied
- But there are too many rules (millions) to check them all
- Instead:
  - ➔ given the underlying chart cells and input words
  - ➔ find which rules apply



## Highlighted Rules

$NP \rightarrow NP_1 \text{ DET}_2 \text{ NN}_3 \mid NP_1 \text{ IN}_2 \text{ NN}_3$   
 $NP \rightarrow NP_1 \mid NP_1$   
 $NP \rightarrow NP_1 \text{ des } \text{NN}_2 \mid NP_1 \text{ of the } \text{NN}_2$   
 $NP \rightarrow NP_1 \text{ des } \text{NN}_2 \mid NP_2 \text{ NP}_1$   
 $NP \rightarrow \text{DET}_1 \text{ NN}_2 \mid \text{DET}_1 \text{ NN}_2$   
 $NP \rightarrow \text{das Haus} \mid \text{the house}$

- CFGs are ubiquitous for describing (syntactic) structure in NLP
- parsing algorithms are core of NL analysis systems
- recognition vs. parsing:
  - ➔ recognition - deciding the membership in the language
  - ➔ parsing - recognition + producing a parse tree for it
- parsing has more time complexity than recognition
- an input may have exponentially many parses

- one of the earliest recognition and parsing algorithms
- standard CKY can only recognize languages defined by CFGs in Chomsky Normal Form (CNF).
- based on a dynamic programming

- considers every possible consecutive subsequence of letters and sets  $K \in T[i, j]$  if the sequence of letters starting from  $i$  to  $j$  can be generated from the non-terminal  $K$
- once it has considered sequences of length 1, it goes on to sequences of length 2, and so on
- for subsequences of length 2 and greater, it considers every possible partition of the subsequence into two halves, and checks to see if there is some production  $A \rightarrow BC$  such that  $B$  matches the 1st half and  $C$  matches the 2nd half. If so, it records  $A$  as matching the whole subsequence
- once completed, the sentence is recognized by the grammar if the entire string is matched by the start symbol

- any portion of the input string spanning  $i$  to  $j$  can be split at  $k$ , and structure can then be built using sub-solutions spanning  $i$  to  $k$  and sub-solutions spanning  $k$  to  $j$
- solution to problem  $[i, j]$  can be constructed from solution to sub problem  $[i, k]$  and solution to sub problem  $[k, j]$

# CKY Algorithm for Deciding CFL

Consider the grammar  $G$  given by:

$$S \rightarrow \varepsilon \mid AB \mid XB$$

$$T \rightarrow AB \mid XB$$

$$X \rightarrow AT$$

$$A \rightarrow a$$

$$B \rightarrow b$$



# CKY Algorithm for Deciding CFL

Now look at  $aaabbb$  :

$$S \rightarrow \varepsilon \mid AB \mid XB$$

$$T \rightarrow AB \mid XB$$

$$X \rightarrow AT$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$a$   $a$   $a$   $b$   $b$   $b$

# CKY Algorithm for Deciding CFL

1) Write variables for all length 1 substrings.

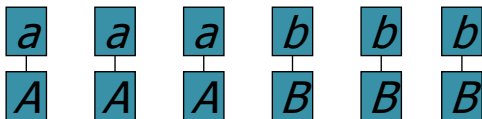
$S \rightarrow \varepsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



# CKY Algorithm for Deciding CFL

2) Write variables for all length 2 substrings.

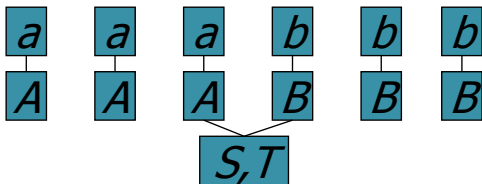
$S \rightarrow \epsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



# CKY Algorithm for Deciding CFL

3) Write variables for all length 3 substrings.

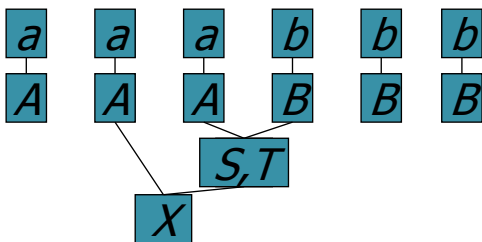
$S \rightarrow \varepsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



# CKY Algorithm for Deciding CFL

4) Write variables for all length 4 substrings.

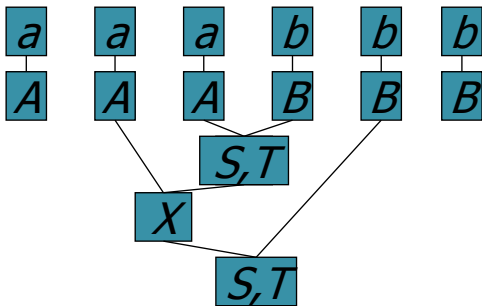
$S \rightarrow \epsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



# CKY Algorithm for Deciding CFL

5) Write variables for all length 5 substrings.

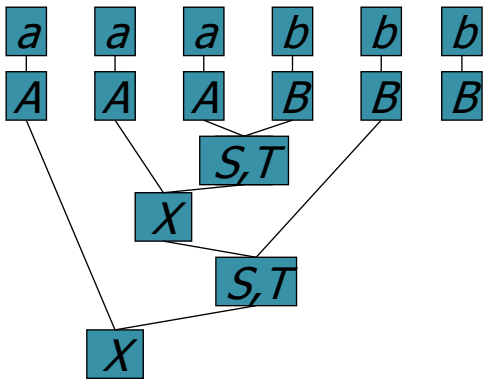
$S \rightarrow \varepsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



# CKY Algorithm for Deciding CFL

6) Write variables for all length 6 substrings.

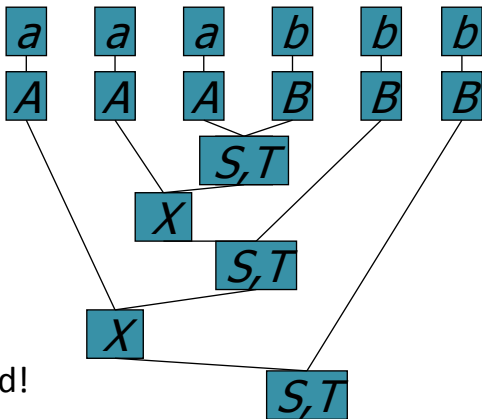
$S \rightarrow \varepsilon \mid AB \mid XB$

$T \rightarrow AB \mid XB$

$X \rightarrow AT$

$A \rightarrow a$

$B \rightarrow b$



S is included so  
*aabbbb* accepted!

# The CKY Algorithm

**function CKY (word  $w$ , grammar  $P$ ) returns table**

**for**  $i \leftarrow$  from 1 to  $\text{LENGTH}(w)$  **do**

$\text{table}[i-1, i] \leftarrow \{A \mid A \rightarrow w_i \in P\}$

**for**  $j \leftarrow$  from 2 to  $\text{LENGTH}(w)$  **do**

**for**  $i \leftarrow$  from  $j-2$  down to 0 **do**

**for**  $k \leftarrow i + 1$  to  $j - 1$  **do**

$\text{table}[i,j] \leftarrow \text{table}[i,j] \cup \{A \mid A \rightarrow BC \in P,$   
                     $B \in \text{table}[i,k], C \in \text{table}[k,j]\}$

If the start symbol  $S \in \text{table}[0,n]$  then  $w \in L(G)$



# CKY Algorithm for Deciding CFL

The table chart used by the algorithm:

$j$ $i$	1	2	3	4	5	6
0						
1						
2						
3						
4						
5						

# CKY Algorithm for Deciding CFL

1. Variables for length 1 substrings.

j \ i	1	2	3	4	5	6
0	<b>a</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>b</b>	<b>b</b>
1	<i>A</i>					
2		<i>A</i>				
3			<i>A</i>			
4				<i>B</i>		
5					<i>B</i>	
6						<i>B</i>

# CKY Algorithm for Deciding CFL

## 2. Variables for length 2 substrings.

$j \backslash i$	1	2	3	4	5	6
0	<b>a</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>b</b>	<b>b</b>
1	A	-				
2		A	-			
3			A	S, T		
4				B	-	
5					B	-
6						B

# CKY Algorithm for Deciding CFL

## 3. Variables for length 3 substrings.

$j \backslash i$	1	2	3	4	5	6
0	<b>a</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>b</b>	<b>b</b>
1	A	-	-			
2		A	-	X		
3			A	S, T	-	
4				B	-	-
5					B	-
6						B

# CKY Algorithm for Deciding CFL

## 4. Variables for length 4 substrings.

$j \backslash i$	1	2	3	4	5	6
	<b>a</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>b</b>	<b>b</b>
0	A	-	-	-		
1		A	-	X	S,T	
2			A	S,T	-	-
3				B	-	-
4					B	-
5						B

# CKY Algorithm for Deciding CFL

## 5. Variables for length 5 substrings.

$j \backslash i$	1	2	3	4	5	6
	a	a	a	b	b	b
0	A	-	-	-	X	
1		A	-	X	S,T	-
2			A	S,T	-	-
3				B	-	-
4					B	-
5						B

# CKY Algorithm for Deciding CFL

6. Variables for *aaabbb*. ACCEPTED!

j \ i	1	2	3	4	5	6
	a	a	a	b	b	b
0	A	-	-	-	X	S,T
1		A	-	X	S,T	-
2			A	S,T	-	-
3				B	-	-
4					B	-
5						B

# CKY Space and Time Complexity

## Time complexity:

- Three nested “for” loop each one of  $O(n)$  size.
- Lookup for  $r = |N|$  pair rules at each step.

**Time complexity –  $O(r^2n^3) = O(n^3)$**

## Space complexity:

- A three dimensions table at size  $n*n*r$      **or**
- A  $n*n$  table with lists up to size of  $r$

**Space complexity –  $O(rn^2) = O(n^2)$**



- doesn't require the grammar to be in CNF.
  - ➔ grammar intended to reflect actual structure of language
  - ➔ conversion to CNF completely destroys the parse structure
- efficiency:
  - ➔ usually moves left-to-right (prefix trees!)
  - ➔ faster than  $O(n^3)$  for many grammars
  - ➔ uses a parse table as CKY, so can backtrack


- **dotted rule**
  - ➔ a partially constructed constituent, w/ the dot indicating what has been found and what is still predicted
  - ➔ generated from ordinary grammar rules (no CNF!)
- maintains a **set of states**, for each position in the input

# The Dotted Rules

With dotted rules, an entry in the chart records:

- Which rule has been used in the analysis
- Which part of the rule has already been found (left of the dot).
- Which part is still predicted to be found and will combine into a complete parse (right of the dot).
- the start and end position of the material left of the dot.

Example:      $A \rightarrow X_1 X_2 \dots \bullet C \dots X_m$



# Parsing Operations

The Earley algorithm has three main operations:

**Predictor:** an **incomplete** entry looks for a symbol to the right of its dot. if there is no matching symbol in the chart, one is predicted by adding all matching rules with an initial dot.

**Scanner:** an **incomplete** entry looks for a symbol to the right of the dot. this prediction is compared to the input, and a complete entry is added to the chart if it matches.

**Completer:** a **complete** edge is combined with an incomplete entry that is looking for it to form another complete entry.

# Parsing Operations

- **Predictor:** If state  $[A \rightarrow X_1 \dots \bullet C \dots X_m, j] \in S_i$  then for every rule of the form  $C \rightarrow Y_1 \dots Y_k$ , add to  $S_i$  the state  $[C \rightarrow \bullet Y_1 \dots Y_k, i]$
- **Scanner:** If state  $[A \rightarrow X_1 \dots \bullet a \dots X_m, j] \in S_i$  and the next input word is  $x_{i+1} = a$ , then add to  $S_{i+1}$  the state  $[A \rightarrow X_1 \dots a \bullet \dots X_m, j]$
- **Completer:** If state  $[A \rightarrow X_1 \dots X_m \bullet, j] \in S_i$  then for every state in  $S_j$  of form  $[B \rightarrow X_1 \dots \bullet A \dots X_k, l]$ , add to  $S_i$  the state  $[B \rightarrow X_1 \dots A \bullet \dots X_k, l]$

# The Earley Recognition Algorithm

The Main Algorithm: parsing input  $w=w_1w_2\dots w_n$

1.  $S_0 = \{[S \rightarrow \bullet P (0)]\}$

2. For  $0 \leq i \leq n$  do:

Process each item  $s \in S_i$  in order by applying to it a *single* applicable operation among:

(a) Predictor (adds new items to  $S_i$ )

(b) Completer (adds new items to  $S_i$ )

(c) Scanner (adds new items to  $S_{i+1}$ )

3. If  $S_{i+1} = \emptyset$  Reject the input.

4. If  $i = n$  and  $[S \rightarrow P \bullet (0)] \in S_n$  then Accept the input.

# Earley Algorithm Example

Consider the following grammar for arithmetic expressions:

$S \rightarrow P$  (the start rule)

$P \rightarrow P + M$

$P \rightarrow M$

$M \rightarrow M * T$

$M \rightarrow T$

$T \rightarrow \text{number}$

With the input:  $2 + 3 * 4$

# Earley Algorithm Example

Sequence(0) • 2 + 3 \* 4

(1)  $S \rightarrow \bullet P(0)$  # start rule



# Earley Algorithm Example

Sequence(0) • 2 + 3 \* 4

- (1)  $S \rightarrow \bullet P$  (0)                   # start rule
- (2)  $P \rightarrow \bullet P + M$  (0)               # predict from (1)
- (3)  $P \rightarrow \bullet M$  (0)                   # predict from (1)

# Earley Algorithm Example

Sequence(0) • 2 + 3 \* 4

- (1)  $S \rightarrow \bullet P$  (0) # start rule
- (2)  $P \rightarrow \bullet P + M$  (0) # predict from (1)
- (3)  $P \rightarrow \bullet M$  (0) # predict from (1)
- (4)  $M \rightarrow \bullet M * T$  (0) # predict from (3)
- (5)  $M \rightarrow \bullet T$  (0) # predict from (3)

# Earley Algorithm Example

Sequence(0) • 2 + 3 \* 4

- (1)  $S \rightarrow \bullet P$  (0) # start rule
- (2)  $P \rightarrow \bullet P + M$  (0) # predict from (1)
- (3)  $P \rightarrow \bullet M$  (0) # predict from (1)
- (4)  $M \rightarrow \bullet M * T$  (0) # predict from (3)
- (5)  $M \rightarrow \bullet T$  (0) # predict from (3)
- (6)  $T \rightarrow \bullet \text{number}$  (0) # predict from (5)

# Earley Algorithm Example

Sequence(1) 2 • + 3 \* 4

(1) T → number • (0)    # scan from S(0)(6)

# Earley Algorithm Example

Sequence(1) 2 • + 3 \* 4

(1)  $T \rightarrow \text{number} \bullet (0)$  # scan from  $S(0)(6)$

(2)  $M \rightarrow T \bullet (0)$  # complete from  $S(0)(5)$

# Earley Algorithm Example

Sequence(1) 2 • + 3 \* 4

- (1)  $T \rightarrow \text{number} \bullet (0)$  # scan from  $S(0)(6)$
- (2)  $M \rightarrow T \bullet (0)$  # complete from  $S(0)(5)$
- (3)  $M \rightarrow M \bullet * T (0)$  # complete from  $S(0)(4)$
- (4)  $P \rightarrow M \bullet (0)$  # complete from  $S(0)(3)$

# Earley Algorithm Example

Sequence(1) 2 • + 3 \* 4

- (1)  $T \rightarrow \text{number} \bullet (0)$  # scan from S(0)(6)
- (2)  $M \rightarrow T \bullet (0)$  # complete from S(0)(5)
- (3)  $M \rightarrow M \bullet * T (0)$  # complete from S(0)(4)
- (4)  $P \rightarrow M \bullet (0)$  # complete from S(0)(3)
- (5)  $P \rightarrow P \bullet + M (0)$  # complete from S(0)(2)
- (6)  $S \rightarrow P \bullet (0)$  # complete from S(0)(1)

# Earley Algorithm Example

Sequence(2) 2 + • 3 \* 4

(1)  $P \rightarrow P + \bullet M$  (0)      # scan from S(1)(5)



# Earley Algorithm Example

Sequence(2) 2 + • 3 \* 4

- (1)  $P \rightarrow P + \bullet M$  (0)      # scan from S(1)(5)
- (2)  $M \rightarrow \bullet M * T$  (2)      # predict from (1)
- (3)  $M \rightarrow \bullet T$  (2)      # predict from (1)

# Earley Algorithm Example

Sequence(2) 2 + • 3 \* 4

- (1)  $P \rightarrow P + \bullet M$  (0) # scan from S(1)(5)
- (2)  $M \rightarrow \bullet M * T$  (2) # predict from (1)
- (3)  $M \rightarrow \bullet T$  (2) # predict from (1)
- (4)  $T \rightarrow \bullet \text{number}$  (2) # predict from (3)

# Earley Algorithm Example

Sequence(3) 2 + 3 • \* 4

(1) T  $\rightarrow$  number • (2) # scan from S(2)(4)

# Earley Algorithm Example

Sequence(3) 2 + 3 • \* 4

- (1)  $T \rightarrow \text{number} \bullet (2)$  # scan from  $S(2)(4)$
- (2)  $M \rightarrow T \bullet (2)$  # complete from  $S(2)(3)$

# Earley Algorithm Example

Sequence(3) 2 + 3 • \* 4

- (1)  $T \rightarrow \text{number} \bullet$  (2) # scan from S(2)(4)
- (2)  $M \rightarrow T \bullet$  (2) # complete from S(2)(3)
- (3)  $M \rightarrow M \bullet * T$  (2) # complete from S(2)(2)
- (4)  $P \rightarrow P + M \bullet$  (0) # complete from S(2)(1)

# Earley Algorithm Example

Sequence(3) 2 + 3 • \* 4

- (1)  $T \rightarrow \text{number} \bullet$  (2) # scan from S(2)(4)
- (2)  $M \rightarrow T \bullet$  (2) # complete from S(2)(3)
- (3)  $M \rightarrow M \bullet * T$  (2) # complete from S(2)(2)
- (4)  $P \rightarrow P + M \bullet$  (0) # complete from S(2)(1)
- (5)  $P \rightarrow P \bullet + M$  (0) # complete from S(0)(2)
- (6)  $S \rightarrow P \bullet$  (0) # complete from S(0)(1)

# Earley Algorithm Example

Sequence(4) 2 + 3 \* • 4

(1)  $M \rightarrow M * \bullet T$  (2) # scan from S(3)(3)

# Earley Algorithm Example

Sequence(4) 2 + 3 \* • 4

- (1)  $M \rightarrow M * \bullet T$  (2) # scan from S(3)(3)  
(2)  $T \rightarrow \bullet \text{number}$  (4) # predict from (1)



# Earley Algorithm Example

Sequence(5) 2 + 3 \* 4 •

(1)  $T \rightarrow \text{number} \bullet (4)$  # scan from  $S(4)(2)$

# Earley Algorithm Example

Sequence(5) 2 + 3 \* 4 •

- (1)  $T \rightarrow \text{number} \bullet (4)$  # scan from  $S(4)(2)$
- (2)  $M \rightarrow M * T \bullet (2)$  # complete from  $S(4)(1)$

# Earley Algorithm Example

Sequence(5) 2 + 3 \* 4 •

- (1)  $T \rightarrow \text{number} \bullet (4)$  # scan from  $S(4)(2)$
- (2)  $M \rightarrow M * T \bullet (2)$  # complete from  $S(4)(1)$
- (3)  $M \rightarrow M \bullet * T (2)$  # complete from  $S(2)(2)$
- (4)  $P \rightarrow P + M \bullet (0)$  # complete from  $S(2)(1)$

# Earley Algorithm Example

Sequence(5) 2 + 3 \* 4 •

- (1)  $T \rightarrow \text{number} \bullet (4)$  # scan from  $S(4)(2)$
- (2)  $M \rightarrow M * T \bullet (2)$  # complete from  $S(4)(1)$
- (3)  $M \rightarrow M \bullet * T (2)$  # complete from  $S(2)(2)$
- (4)  $P \rightarrow P + M \bullet (0)$  # complete from  $S(2)(1)$
- (5)  $P \rightarrow P \bullet + M (0)$  # complete from  $S(0)(2)$
- (6)  $S \rightarrow P \bullet (0)$  # complete from  $S(0)(1)$

# Earley Algorithm Example

Sequence(5) 2 + 3 \* 4 •

- (1)  $T \rightarrow \text{number} \bullet (4)$  # scan from  $S(4)(2)$
- (2)  $M \rightarrow M * T \bullet (2)$  # complete from  $S(4)(1)$
- (3)  $M \rightarrow M \bullet * T (2)$  # complete from  $S(2)(2)$
- (4)  $P \rightarrow P + M \bullet (0)$  # complete from  $S(2)(1)$
- (5)  $P \rightarrow P \bullet + M (0)$  # complete from  $S(0)(2)$
- (6)  $S \rightarrow P \bullet (0)$  # complete from  $S(0)(1)$

The state  $S \rightarrow P \bullet (0)$  represents a completed parse.

# Finding the parse tree

<u>Seq 0</u>	<u>Seq 1</u>	<u>Seq 2</u>	<u>Seq 3</u>	<u>Seq 4</u>	<u>Seq 5</u>
$\bullet 2 + 3 * 4$	$2 \bullet + 3 * 4$	$2 + \bullet 3 * 4$	$2 + 3 \bullet * 4$	$2 + 3 * \bullet 4$	$2 + 3 * 4 \bullet$
$S \rightarrow \bullet P (0)$	$T \rightarrow '2' \bullet (0)$	$P \rightarrow P + \bullet M (0)$	$T \rightarrow '3' \bullet (2)$	$M \rightarrow M * \bullet T (2)$	$T \rightarrow '4' \bullet (4)$
$P \rightarrow \bullet P + M (0)$	$M \rightarrow T \bullet (0)$	$M \rightarrow \bullet M * T (2)$	$M \rightarrow T \bullet (2)$	$T \rightarrow \bullet \text{num} (4)$	$M \rightarrow M * T \bullet (2)$
$P \rightarrow \bullet M (0)$	$M \rightarrow M \bullet * T (0)$	$M \rightarrow \bullet T (2)$	$M \rightarrow M \bullet * T (2)$		$M \rightarrow M \bullet * T (2)$
$M \rightarrow \bullet M * T (0)$	$P \rightarrow M \bullet (0)$	$T \rightarrow \bullet \text{num} (2)$	$P \rightarrow P + M \bullet (0)$		$P \rightarrow P + M \bullet (0)$
$M \rightarrow \bullet T (0)$	$P \rightarrow P \bullet + M (0)$		$P \rightarrow P \bullet + M (0)$		$P \rightarrow P \bullet + M (0)$
$T \rightarrow \bullet \text{num} (0)$	$S \rightarrow P \bullet (0)$		$S \rightarrow P \bullet (0)$		$S \rightarrow P \bullet (0)$

# Finding the parse tree

<u>Seq 0</u>	<u>Seq 1</u>	<u>Seq 2</u>	<u>Seq 3</u>	<u>Seq 4</u>	<u>Seq 5</u>
$\bullet 2 + 3 * 4$	$2 \bullet + 3 * 4$	$2 + \bullet 3 * 4$	$2 + 3 \bullet * 4$	$2 + 3 * \bullet 4$	$2 + 3 * 4 \bullet$
$S \rightarrow \bullet P (0)$	$T \rightarrow '2' \bullet (0)$	$P \rightarrow P + \bullet M (0)$	$T \rightarrow '3' \bullet (2)$	$M \rightarrow M * \bullet T (2)$	$T \rightarrow '4' \bullet (4)$
$P \rightarrow \bullet P + M (0)$	$M \rightarrow T \bullet (0)$	$M \rightarrow \bullet M * T (2)$	$M \rightarrow T \bullet (2)$	$T \rightarrow \bullet \text{num} (4)$	$M \rightarrow M * T \bullet (2)$
$P \rightarrow \bullet M (0)$	$M \rightarrow M \bullet * T (0)$	$M \rightarrow \bullet T (2)$	$M \rightarrow M \bullet * T (2)$		$M \rightarrow M \bullet * T (2)$
$M \rightarrow \bullet M * T (0)$	$P \rightarrow M \bullet (0)$	$T \rightarrow \bullet \text{num} (2)$	$P \rightarrow P + M \bullet (0)$		$P \rightarrow P + M \bullet (0)$
$M \rightarrow \bullet T (0)$	$P \rightarrow P \bullet + M (0)$		$P \rightarrow P \bullet + M (0)$		$P \rightarrow P \bullet + M (0)$
$T \rightarrow \bullet \text{num} (0)$	$S \rightarrow P \bullet (0)$		$S \rightarrow P \bullet (0)$		$S \rightarrow P \bullet (0)$

# Finding the parse tree

<u>Seq 0</u>	<u>Seq 1</u>	<u>Seq 2</u>	<u>Seq 3</u>	<u>Seq 4</u>	<u>Seq 5</u>
$\bullet 2 + 3 * 4$	$2 \bullet + 3 * 4$	$2 + \bullet 3 * 4$	$2 + 3 \bullet * 4$	$2 + 3 * \bullet 4$	$2 + 3 * 4 \bullet$
$S \rightarrow \bullet P$ (0)	$T \rightarrow '2' \bullet$ (0)	$P \rightarrow P + \bullet M$ (0)	$T \rightarrow '3' \bullet$ (2)	$M \rightarrow M * \bullet T$ (2)	$T \rightarrow '4' \bullet$ (4)
$P \rightarrow \bullet P + M$ (0)	$M \rightarrow T \bullet$ (0)	$M \rightarrow \bullet M * T$ (2)	$M \rightarrow T \bullet$ (2)	$T \rightarrow \bullet \text{num}$ (4)	$M \rightarrow M * T \bullet$ (2)
$P \rightarrow \bullet M$ (0)	$M \rightarrow M \bullet * T$ (0)	$M \rightarrow \bullet T$ (2)	$M \rightarrow M \bullet * T$ (2)		$M \rightarrow M \bullet * T$ (2)
$M \rightarrow \bullet M * T$ (0)	$P \rightarrow M \bullet$ (0)	$T \rightarrow \bullet \text{num}$ (2)	$P \rightarrow P + M \bullet$ (0)		$P \rightarrow P + M \bullet$ (0)
$M \rightarrow \bullet T$ (0)	$P \rightarrow P \bullet + M$ (0)		$P \rightarrow P \bullet + M$ (0)		$P \rightarrow P \bullet + M$ (0)
$T \rightarrow \bullet \text{num}$ (0)	$S \rightarrow P \bullet$ (0)		$S \rightarrow P \bullet$ (0)		$S \rightarrow P \bullet$ (0)



# Finding the parse tree

<u>Seq 0</u>	<u>Seq 1</u>	<u>Seq 2</u>	<u>Seq 3</u>	<u>Seq 4</u>	<u>Seq 5</u>
$\bullet 2 + 3 * 4$	$2 \bullet + 3 * 4$	$2 + \bullet 3 * 4$	$2 + 3 \bullet * 4$	$2 + 3 * \bullet 4$	$2 + 3 * 4 \bullet$
$S \rightarrow \bullet P$ (0)	$T \rightarrow '2' \bullet$ (0)	$P \rightarrow P + \bullet M$ (0)	$T \rightarrow '3' \bullet$ (2)	$M \rightarrow M * \bullet T$ (2)	$T \rightarrow '4' \bullet$ (4)
$P \rightarrow \bullet P + M$ (0)	$M \rightarrow T \bullet$ (0)	$M \rightarrow \bullet M * T$ (2)	$M \rightarrow T \bullet$ (2)	$T \rightarrow \bullet \text{num}$ (4)	$M \rightarrow M * T \bullet$ (2)
$P \rightarrow \bullet M$ (0)	$M \rightarrow M \bullet * T$ (0)	$M \rightarrow \bullet T$ (2)	$M \rightarrow M \bullet * T$ (2)		$M \rightarrow M \bullet * T$ (2)
$M \rightarrow \bullet M * T$ (0)	$P \rightarrow M \bullet$ (0)	$T \rightarrow \bullet \text{num}$ (2)	$P \rightarrow P + M \bullet$ (0)		$P \rightarrow P + M \bullet$ (0)
$M \rightarrow \bullet T$ (0)	$P \rightarrow P \bullet + M$ (0)		$P \rightarrow P \bullet + M$ (0)		$P \rightarrow P \bullet + M$ (0)
$T \rightarrow \bullet \text{num}$ (0)	$S \rightarrow P \bullet$ (0)		$S \rightarrow P \bullet$ (0)		$S \rightarrow P \bullet$ (0)

**Input:** Foreign sentence  $\mathbf{f} = f_1, \dots, f_{l_f}$ , with syntax tree

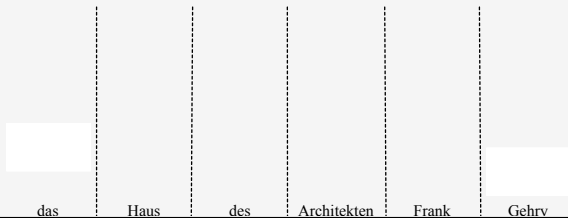
**Output:** English translation  $\mathbf{e}$

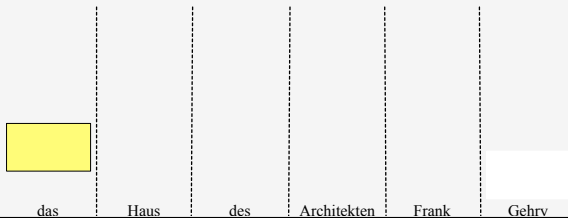
```

1: for  $i=0 \dots \text{length}(\mathbf{f})-1$  do // initialize chart
2:   store pointer to initial node in prefix tree in span  $[i,i]$ 
3: end for
4: for  $l=1 \dots l_f$  do // build chart from the bottom up
5:   for  $\text{start}=0 \dots l_f-1$  do // beginning of span
6:      $\text{end} = \text{start}+1$ 
7:     for  $\text{midpoint}=\text{start} \dots \text{end}-1$  do
8:       for all dotted rules  $d$  in span  $[\text{start},\text{midpoint}]$  do
9:         for all distinct head node nonterminals or input words  $h$  covering
           span  $[\text{midpoint}+1,\text{end}]$  do
10:          if extension  $d \rightarrow h$  exists in prefix tree then
11:             $d_{\text{new}} = d \rightarrow h$ 
12:            for all complete rules at  $d_{\text{new}}$  do
13:              apply rules
14:              store chart entries in span  $[\text{start},\text{end}]$ 
15:            end for
16:            if extension exist for  $d_{\text{new}}$  then
17:              store  $d_{\text{new}}$  in span  $[\text{start},\text{end}]$  // new dotted rule
18:            end if
19:          end if
20:        end for
21:      end for
22:    end for
23:  end for
24: end for
25: return English translation  $\mathbf{e}$  from best chart entry in span  $[0, l_f]$ 

```

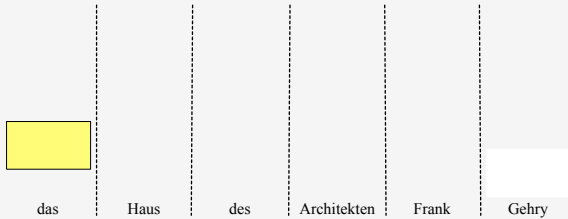
# Finding Applicable Rules in Prefix Tree





# Looking up Rules in the Prefix Tree

● — das 1



# Taking Note of the Dotted Rule

● — das ①

das ①

das

Haus

des

Architekten

Frank

Gehry

# Checking if Dotted Rule has Translations

● — das ① DET: the  
DET: that

das ①

das

Haus

des

Architekten

Frank

Gehry

# Applying the Translation Rules

das ① DET: the  
DET: that

DET: that  
DET: the  
das ①

das

Haus

des

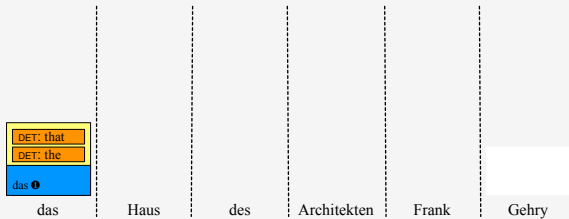
Architekten

Frank

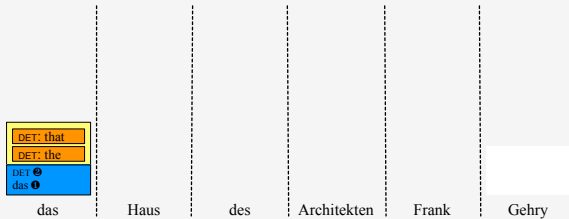
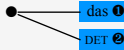
Gehry



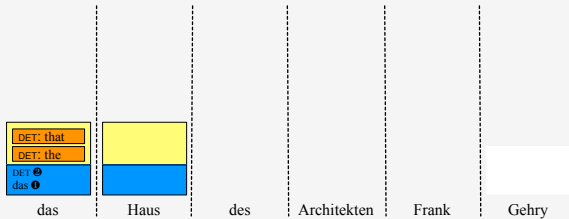
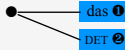
# Looking up Constituent Label in Prefix Tree



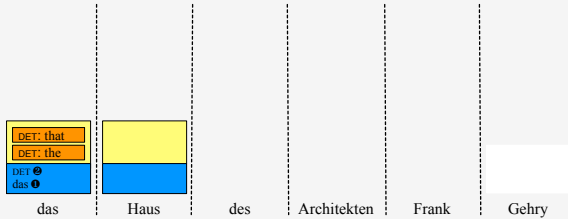
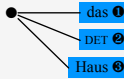
# Add to Span's List of Dotted Rules



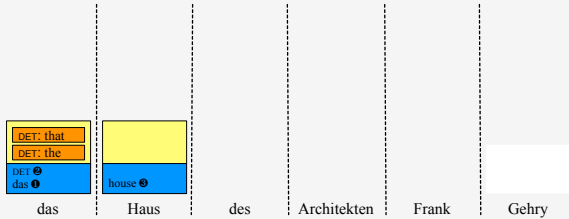
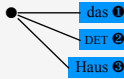
# Moving on to the Next Cell



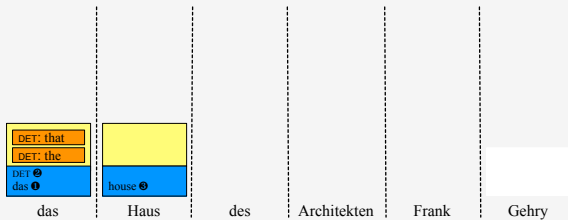
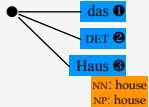
# Looking up Rules in the Prefix Tree



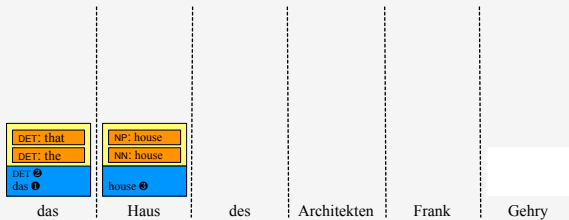
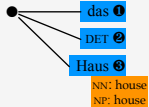
# Taking Note of the Dotted Rule



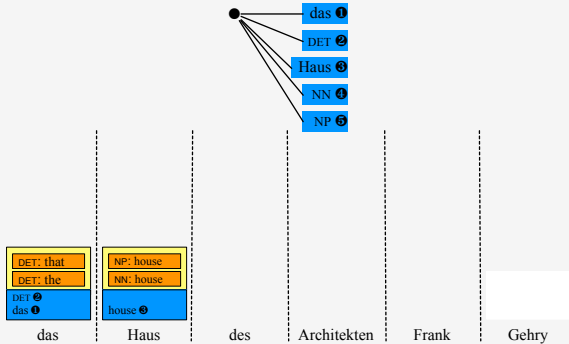
# Checking if Dotted Rule has Translations



# Applying the Translation Rules

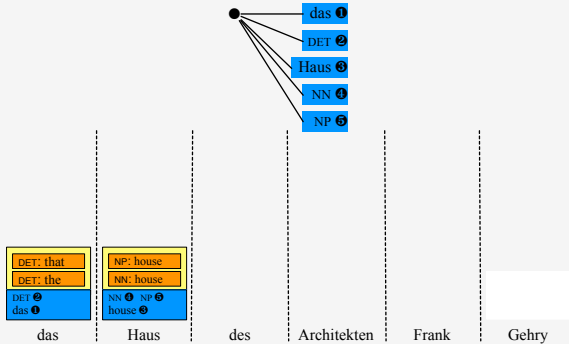


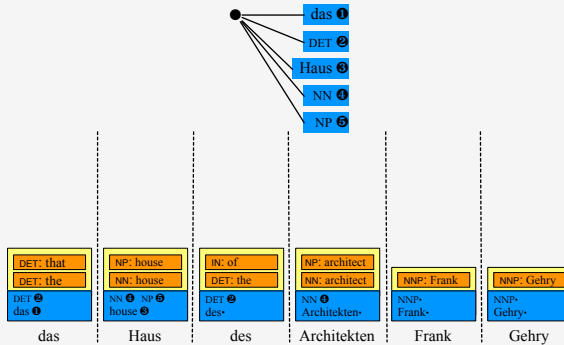
# Looking up Constituent Label in Prefix Tree



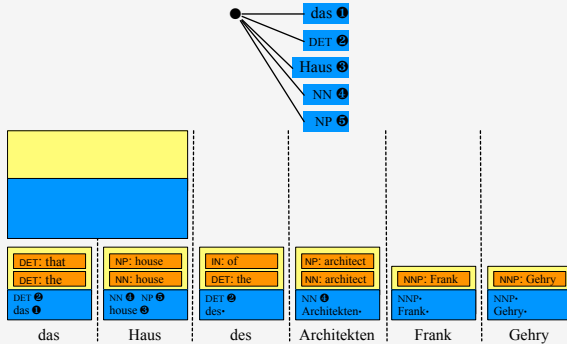


# Add to Span's List of Dotted Rules





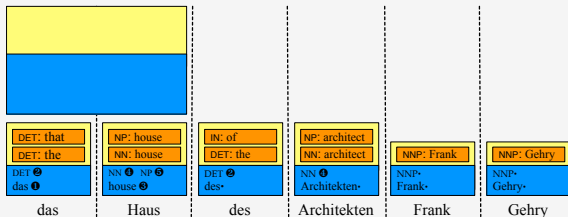
# Moving on to the Next Cell

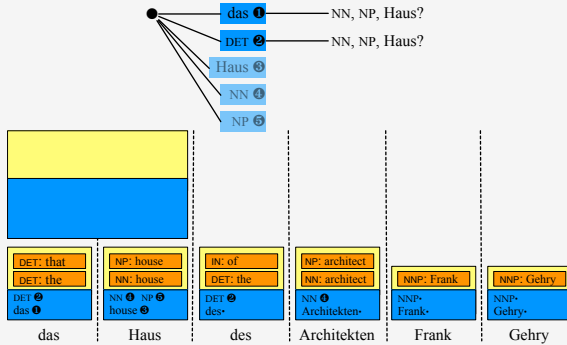


Cannot consume multiple words at once

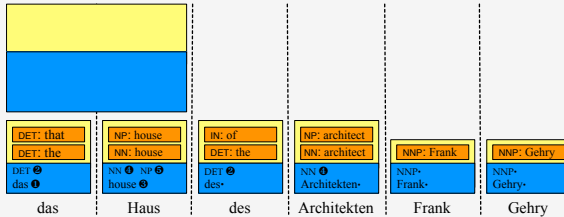
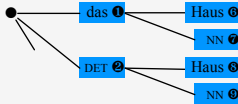
All rules are extensions of existing dotted rules

Here: only extensions of span over *das* possible

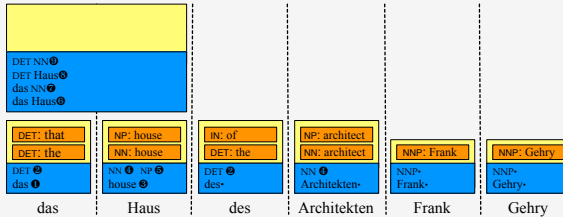
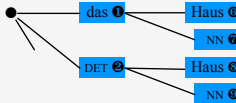




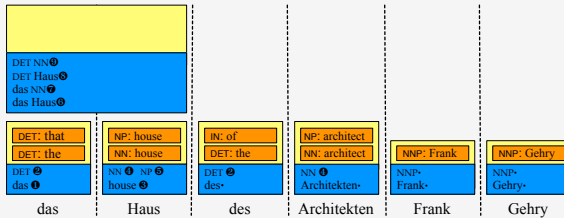
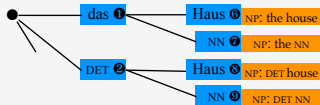
# Looking up Rules in the Prefix Tree



# Taking Note of the Dotted Rule

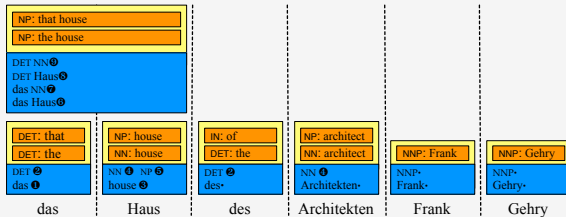
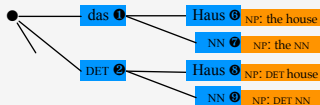


# Checking if Dotted Rules have Translations

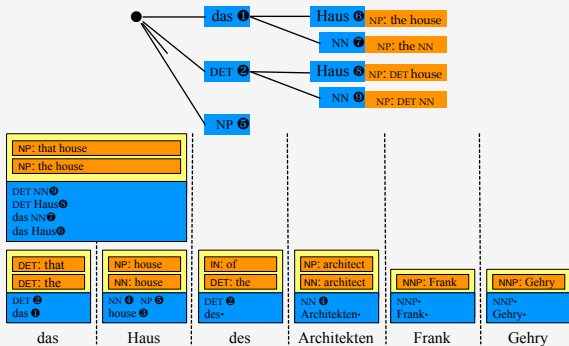




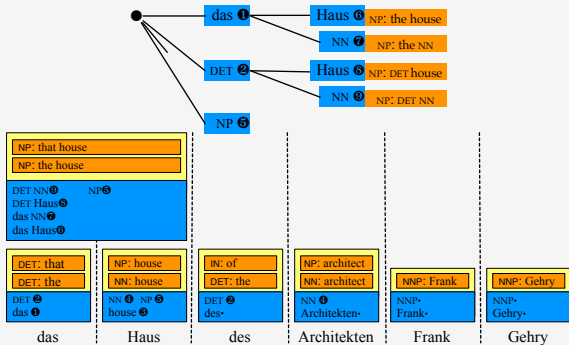
# Applying the Translation Rules



# Looking up Constituent Label in Prefix Tree



# Add to Span's List of Dotted Rules



- Complexity  $O(rn^3)$  with sentence length  $n$  and size of dotted rule list  $r$ 
  - ➔ may introduce maximum size for spans that do not start at beginning
  - ➔ may limit size of dotted rule list (very arbitrary)
  
- Does the list of dotted rules explode?
  
- Yes, if there are many rules with neighboring target-side non-terminals
  - ➔ such rules apply in many places
  - ➔ rules with words are much more restricted

- Some rules may apply in too many ways
- Neighboring input non-terminals

$$VP \rightarrow \textit{gibt} X_1 X_2 \mid \textit{gives} NP_2 \textit{ to} NP_1$$

- ➔ non-terminals may match many different pairs of spans
  - ➔ especially a problem for hierarchical models (no constituent label restrictions)
  - ➔ may be okay for syntax-models
- Three neighboring input non-terminals

$$VP \rightarrow \textit{trifft} X_1 X_2 X_3 \textit{ heute} \mid \textit{meets} NP_1 \textit{ today} PP_2 PP_3$$

- ➔ will get out of hand even for syntax models

- We know which rules apply
- We know where they apply (each non-terminal tied to a span)
- But there are still many choices
  - ➔ many possible translations
  - ➔ each non-terminal may match multiple hypotheses
  - ➔ number choices exponential with number of non-terminals

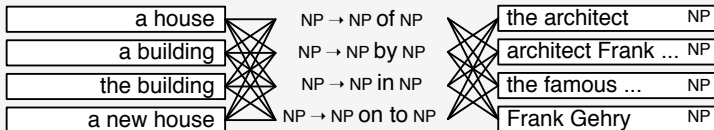
Found applicable rules  $PP \rightarrow des X \mid \dots NP \dots$



- Non-terminal will be filled any of  $h$  underlying matching hypotheses
  - Choice of  $t$  lexical translations
- $\Rightarrow$  Complexity  $O(ht)$

(note: we may not group rules by target constituent label, so a rule  $NP \rightarrow des X \mid the NP$  would also be considered here as well)

Found applicable rule  $NP \rightarrow X_1 \text{ des } X_2 \mid NP_1 \dots NP_2$



- Two non-terminal will be filled any of  $h$  underlying matching hypotheses each
  - Choice of  $t$  lexical translations
- ⇒ Complexity  $O(h^2t)$  — a three-dimensional “cube” of choices

(note: rules may also reorder differently)



		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0				
a building	1.3				
the building	2.2				
a new house	2.6				

Arrange all the choices in a “cube”

(here: a square, generally a orthotope, also called a hyperrectangle)

	1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house 1.0	2.1			
a building 1.3				
the building 2.2				
a new house 2.6				

- Hypotheses created in cube: (0,0)

# Add ("Pop") Hypothesis to Chart Cell

		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0	2.1			
a building	1.3				
the building	2.2				
a new house	2.6				

- Hypotheses created in cube:  $\epsilon$
- Hypotheses in chart cell stack: (0,0)

# Create Neighboring Hypotheses

		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0	2.1	2.5		
a building	1.3	2.7			
the building	2.2				
a new house	2.6				

- Hypotheses created in cube: (0,1), (1,0)
- Hypotheses in chart cell stack: (0,0)

		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0	2.1	2.5		
a building	1.3	2.7			
the building	2.2				
a new house	2.6				

- Hypotheses created in cube: (0,1)
- Hypotheses in chart cell stack: (0,0), (1,0)

# Create Neighboring Hypotheses

		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0	2.1	2.5	3.1	
a building	1.3	2.7	2.4		
the building	2.2				
a new house	2.6				

- Hypotheses created in cube: (0,1), (1,1), (2,0)
- Hypotheses in chart cell stack: (0,0), (1,0)

		1.5 in the ...	1.7 by architect ...	2.6 by the ...	3.2 of the ...
a house	1.0	2.1	2.5	3.1	
a building	1.3	2.7	2.4	3.0	
the building	2.2		3.8		
a new house	2.6				

- Hypotheses created in cube: (0,1), (1,2), (2,1), (2,0)
- Hypotheses in chart cell stack: (0,0), (1,0), (1,1)

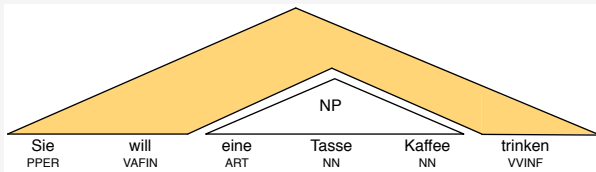
- Several groups of rules will apply to a given span
  - Each of them will have a cube
  - We can create a queue of cubes
- ⇒ Always pop off the most promising hypothesis, regardless of cube
- 
- May have separate queues for different target constituent labels



- 1: **for** all spans (bottom up) **do**
- 2:   extend dotted rules
- 3:   **for all** dotted rules **do**
- 4:     find group of applicable rules
- 5:     create a cube for it
- 6:     create first hypothesis in cube
- 7:     place cube in queue
- 8:   **end for**
- 9:   **for** specified number of pops **do**
- 10:     pop off best hypothesis of any cube in queue
- 11:     add it to the chart cell
- 12:     create its neighbors
- 13:   **end for**
- 14:   extend dotted rules over constituent labels
- 15: **end for**

- First stage: decoding without a language model (-LM decoding)
  - ➔ may be done exhaustively
  - ➔ optionally prune out low scoring hypotheses
- Second stage: add language model
  - ➔ limited to packed chart obtained in first stage
- Note: essentially, we do two-stage decoding for each span at a time

- Which spans should be more emphasized in search?
- Initial decoding stage can provide outside cost estimates



- Use min/max language model costs to obtain admissible heuristic (or at least something that will guide search better)

- Synchronous context free grammars
- Extracting rules from a syntactically parsed parallel corpus
- Bottom-up decoding
- Chart organization: dynamic programming, stacks, pruning
- Prefix tree for rules
- Dotted rules
- Cube pruning