

# Statistical Machine Translation

-tree-based models-

**Artem Sokolov**

Computerlinguistik  
Universität Heidelberg  
Sommersemester 2015

material from P. Koehn, S. Riezler

- Traditional statistical models operate on sequences of words
- Linguistic theories tell us that a deeper structure exists
- Many translation problems can be best explained by pointing to syntax
  - ➔ reordering, e.g., verb movement in German–English translation
  - ➔ long distance agreement (e.g., subject-verb) in output
- ⇒ Translation models based on tree representation of language
  - ➔ significant ongoing research
  - ➔ state-of-the art for some language pairs
  - ➔ trend: as technology matures even more pairs are better handled by tree-based SMT

**Idea:** word groups should correspond to constituents of certain roles and functions

- Phrase structure
  - ➔ noun phrases: the big man, a house, ...
  - ➔ prepositional phrases: at 5 o'clock, in Edinburgh, ...
  - ➔ verb phrases: going out of business, eat chicken, ...
  - ➔ adjective phrases, angry with the high prices, faster than you, ...
- Weighted Context-free Grammars (CFG)
  - ➔  $G = \langle N, T, (P, \pi), (S, \sigma) \rangle$
  - ➔ non-terminal symbols  $N$ : phrase structure labels, part-of-speech tags
  - ➔ terminal symbols  $T$ : words
  - ➔ production rules  $P: N \rightarrow (N \cup T)^*$
  - ➔ weights of production rules:  $\pi : P \rightarrow K$  ( $K$  is a semiring)
  - ➔ start symbol  $S$
  - ➔ weights of start states:  $\sigma : S \rightarrow K$  ( $K$  is a semiring) example:  $NP \rightarrow$   
DET NN  
example:  $NP \rightarrow$  DET house

## Weighted Synchronous Context-free Grammars (SCFG)

- $G = \langle N, T^1, T^2, (P, \pi), (S, \sigma) \rangle$
- non-terminal symbols  $N$
- source and target terminal symbols  $T^1, T^2$
- production rules  $P: N \rightarrow (N \cup T^1)^* \times (\{\boxed{1}, \boxed{2}, \dots\} \cup T^2)^*$
- weights of production rules:  $\pi: P \rightarrow K$  ( $K$  is a semiring)
- start symbol  $S$
- weights of start states:  $\sigma: S \rightarrow K$  ( $K$  is a semiring)

- Nonterminal rules

$$\text{NP} \rightarrow \text{DET}_1 \text{ NN}_2 \text{ JJ}_3 \mid \text{DET}_1 \text{ JJ}_3 \text{ NN}_2$$

- Terminal rules

$$\begin{aligned} \text{N} &\rightarrow \text{maison} \mid \text{house} \\ \text{NP} &\rightarrow \text{la maison bleue} \mid \text{the blue house} \end{aligned}$$

- Mixed rules

$$\text{NP} \rightarrow \text{la maison JJ}_1 \mid \text{the JJ}_1 \text{ house}$$

- Translation by parsing
  - ➔ synchronous grammar has to parse entire input sentence
  - ➔ output tree is generated at the same time
  - ➔ process is broken up into a number of rule applications
- Each rule is weighted (definition)
- Total translation probability

$$\text{SCORE}(\text{TREE}, E, F) = \prod_i \text{RULE}_i$$

- Many ways to assign probabilities to rules (as there are many parses possible)

$$\Sigma = \{a, b, c\} \quad \Delta = \{x, y, z\} \quad V = \{A, B, C\} \quad S = \{C\} \quad \sigma(C) = 1.0$$

$$R = \left\{ \begin{array}{l} A \xrightarrow{0.5} \langle A a B, \boxed{2} \boxed{1} x \rangle, \\ A \xrightarrow{0.5} \langle b C, y \boxed{1} \rangle, \\ B \xrightarrow{0.6} \langle a b c, z z \rangle, \\ B \xrightarrow{0.4} \langle a b c, x y z \rangle, \\ C \xrightarrow{0.8} \langle A A, \boxed{1} \boxed{2} \rangle, \\ C \xrightarrow{0.2} \langle c, z \rangle \end{array} \right\}$$

Figure 2.4: Example of a weighted synchronous context-free grammar (WSCFG). Note that C is the start symbol.

Yield	$\alpha \rightarrow \langle \beta, \gamma \rangle$	Weight
$\langle C, \boxed{1} \rangle \Rightarrow$	$C \rightarrow \langle A A, \boxed{1} \boxed{2} \rangle$	$1.0 \times 0.8$
$\langle A A, \boxed{1} \boxed{2} \rangle \Rightarrow$	$A \rightarrow \langle A a B, \boxed{2} \boxed{1} x \rangle$	$\times 0.5$
$\langle A a B A, \boxed{2} \boxed{1} x \boxed{3} \rangle \Rightarrow$	$A \rightarrow \langle b C, y \boxed{1} \rangle$	$\times 0.5$
$\langle b C a B A, \boxed{2} y \boxed{1} x \boxed{3} \rangle \Rightarrow$	$C \rightarrow \langle c, z \rangle$	$\times 0.2$
$\langle b c a B A, \boxed{1} y z x \boxed{2} \rangle \Rightarrow$	$B \rightarrow \langle a b c, x y z \rangle$	$\times 0.4$
$\langle b c a a b c A, x y z y z x \boxed{1} \rangle \Rightarrow$	$A \rightarrow \langle b C, y \boxed{1} \rangle$	$\times 0.5$
$\langle b c a a b c b C, x y z y z x y \boxed{1} \rangle \Rightarrow$	$C \rightarrow \langle c, z \rangle$	$\times 0.2$
$\langle b c a a b c b c, x y z y z x y z \rangle$		$= 0.0016$

Figure 2.5: Example synchronous derivation using the WSCFG shown in Figure 2.4.

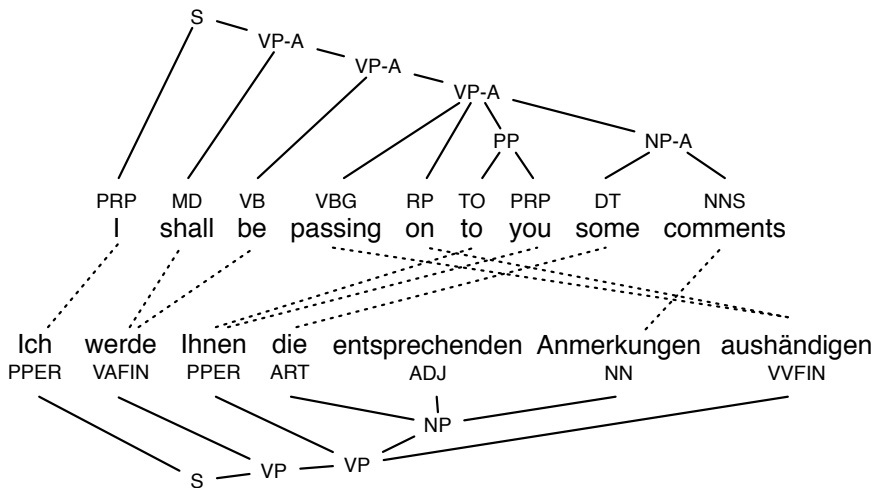
- so far assumed no particular syntax theory
- can use real syntactic annotation
- will need to store internal structure in the rule

## **Benefits:**

- input language syntax puts some constraints on the extracted rules
- output language can have a better-formed (syntactically) structure

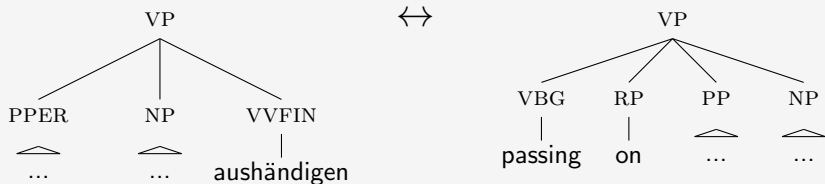


# Synchronous Tree-Substitution Grammars



Phrase structure grammar trees with word alignment  
(German–English sentence pair.)

- Subtree alignment



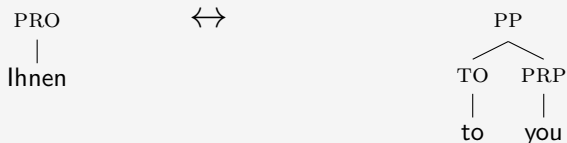
- Synchronous grammar rule

$$VP \rightarrow PPER_1 NP_2 \text{ aushändigen} \mid \text{passing on } PP_1 NP_2$$

- Note:

- ➔ one word *aushändigen* mapped to two words *passing on*
- ➔ effortless capture of reordering
- ➔ but: fully non-terminal rule not possible  
(one-to-one mapping constraint for nonterminals)

- Subtree alignment



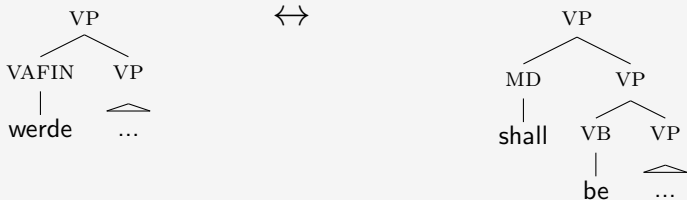
- Synchronous grammar rule (stripping out English internal structure)

$PRO/PP \rightarrow \text{Ihnen} \mid \text{to you}$

- Rule with internal structure

$PRO/PP \rightarrow \text{Ihnen} \mid \begin{array}{l} \text{TO} \quad \text{PRP} \\ | \quad | \\ \text{to} \quad \text{you} \end{array}$

- Translation of German *werde* to English *shall* be



- Translation rule needs to include mapping of VP
- ⇒ Complex rule

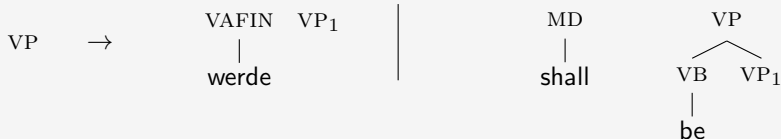


- Stripping out internal structure

$$VP \rightarrow \text{werde } VP_1 \mid \text{shall be } VP_1$$

$\Rightarrow$  synchronous context free grammar

- Maintaining internal structure



$\Rightarrow$  synchronous tree substitution grammar

- Extracting rules from a word-aligned parallel corpus
- First: Hierarchical phrase-based model
  - ➔ only **one** non-terminal symbol  $x$
  - ➔ no linguistic syntax, just a formally syntactic model
- Then: Synchronous phrase structure model
  - ➔ non-terminals for words and phrases: NP, VP, PP, ADJ, ...
  - ➔ corpus must be parsed with syntactic parser

- Suppose we want to learn a rule for *werde ... aushändigen*
- phrase-based SMT will probably fail here
  - ➔ the gap is too large, likely inconsistent
  - ➔ if extracted will contain all words in between (rarely applicable)

# Extracting Phrase Translation Rules

	Ich	werde	Ihnen	die	entsprechenden	Anmerkungen	aushändigen
I	■	■					
shall	■	■					
be	■	■					
passing							■
on							■
to			■				
you			■				
some				■			
comments						■	

→ shall be = werde



# Extracting Phrase Translation Rules

	Ich	werde	Ihnen	die	entsprechenden	Anmerkungen	aushändigen
I	black			light blue	light blue	light blue	
shall		black		light blue	light blue	light blue	
be		black		light blue	light blue	light blue	
passing				light blue	light blue	light blue	black
on				light blue	light blue	light blue	black
to			black	light blue	light blue	light blue	
you			black	light blue	light blue	light blue	
some	light blue	light blue	light blue	dark blue	blue	blue	
comments	light blue	light blue	light blue	blue	blue	dark blue	

some comments =  
die entsprechenden Anmerkungen

# Extracting Phrase Translation Rules

	Ich	werde	Ihnen	die	entsprechenden	Anmerkungen	aushändigen
I							
shall							
be							
passing							
on							
to							
you							
some							
comments							

werde Ihnen die entsprechenden  
Anmerkungen aushändigen  
= shall be passing on to you  
some comments

# Extracting Phrase Translation Rules

	Ich	werde	Ihnen	die	entsprechenden	Anmerkungen	aushändigen
I							
shall							
be							
passing							
on							
to							
you							
some							
comments							

subtracting  
subphrase

werde X aushändigen  
= shall be passing on X

- Recall: consistent phrase pairs

$(\bar{e}, \bar{f})$  consistent with  $A \Leftrightarrow$

$$\begin{aligned} & \forall e_i \in \bar{e} : (e_i, f_j) \in A \rightarrow f_j \in \bar{f} \\ \text{AND } & \forall f_j \in \bar{f} : (e_i, f_j) \in A \rightarrow e_i \in \bar{e} \\ \text{AND } & \exists e_i \in \bar{e}, f_j \in \bar{f} : (e_i, f_j) \in A \end{aligned}$$

- Let  $P$  be the set of all extracted phrase pairs  $(\bar{e}, \bar{f})$

- Extend recursively:

if  $(\bar{e}, \bar{f}) \in P$  AND  $(\bar{e}_{\text{SUB}}, \bar{f}_{\text{SUB}}) \in P$

AND  $\bar{e} = \bar{e}_{\text{PRE}} + \bar{e}_{\text{SUB}} + \bar{e}_{\text{POST}}$

AND  $\bar{f} = \bar{f}_{\text{PRE}} + \bar{f}_{\text{SUB}} + \bar{f}_{\text{POST}}$

AND  $\bar{e} \neq \bar{e}_{\text{SUB}}$  AND  $\bar{f} \neq \bar{f}_{\text{SUB}}$

add  $(e_{\text{PRE}} + X + e_{\text{POST}}, f_{\text{PRE}} + X + f_{\text{POST}})$  to  $P$

(note: any of  $e_{\text{PRE}}$ ,  $e_{\text{POST}}$ ,  $f_{\text{PRE}}$ , or  $f_{\text{POST}}$  may be empty)

- Set of hierarchical phrase pairs is the closure under this extension mechanism

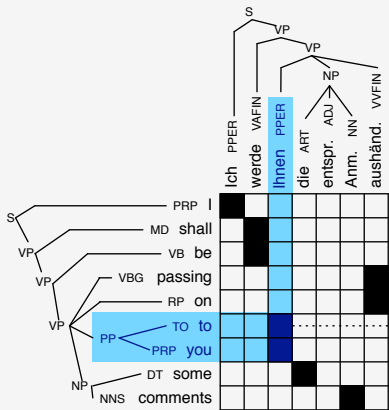
- Removal of multiple sub-phrases leads to rules with multiple non-terminals, such as:

$$Y \rightarrow X_1 X_2 \mid X_2 \textit{ of } X_1$$

- Typical restrictions to limit complexity [Chiang, 2005], to avoid exponential explosion
  - ➔ at most 2 nonterminal symbols
  - ➔ at least 1 but at most 5 words per language
  - ➔ span at most 15 words (counting gaps)
  - ➔ no 2 non-terminals are next to each other in both languages

Even without syntax tree-based models often gain about 1-2 BLEU points over phrase-based systems.

# Learning Syntactic Translation Rules



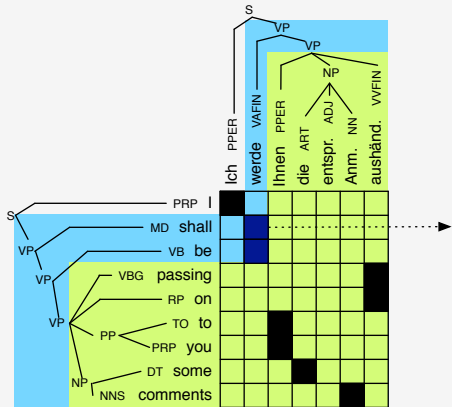
PRO =  
 |  
 Ihnen

PP  
 / \  
 TO PRP  
 | |  
 to you

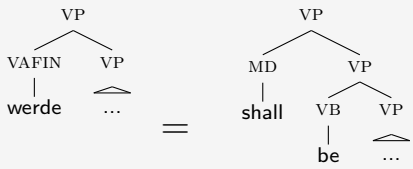
- Hierarchical: rule can cover any span  
⇔ syntactic rules must cover constituents in the tree ⇒ 1 node on top
- Hierarchical: gaps may cover any span  
⇔ gaps must cover constituents in the tree
- Moving up the tree introduces non-terminals
- Much less rules are extracted (all things being equal)







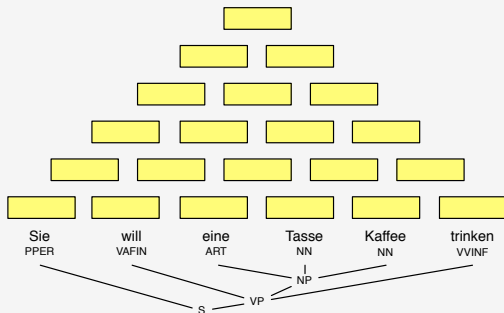
Rule with this phrase pair requires syntactic context



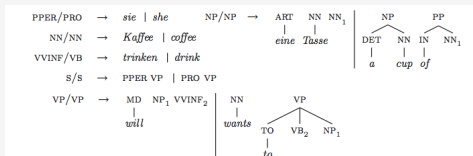
- Extract all rules from corpus
- Score based on counts
  - ➔ joint rule probability:  $p(\text{LHS}, \text{RHS}_f, \text{RHS}_e)$
  - ➔ rule application probability:  $p(\text{RHS}_f, \text{RHS}_e | \text{LHS})$
  - ➔ direct translation probability:  $p(\text{RHS}_e | \text{RHS}_f, \text{LHS})$
  - ➔ noisy channel translation probability:  $p(\text{RHS}_f | \text{RHS}_e, \text{LHS})$
  - ➔ lexical translation probability:  $\prod_{e_i \in \text{RHS}_e} p(e_i | \text{RHS}_f, a)$

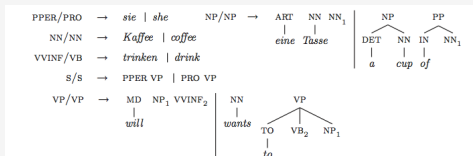
Inspired by monolingual syntactic chart parsing:

During decoding of the source sentence,  
a chart with translations for the  $O(n^2)$  spans has to be filled

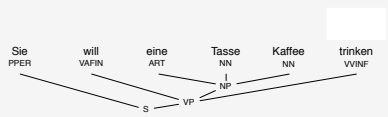


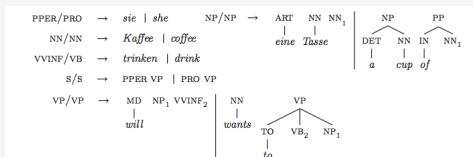
**note:** constrains limit the branching factor



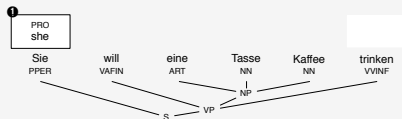


3





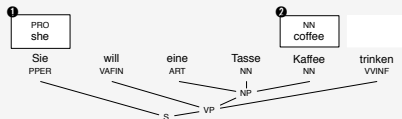
0



Purely lexical rule: filling a span with a translation (a constituent in the chart)

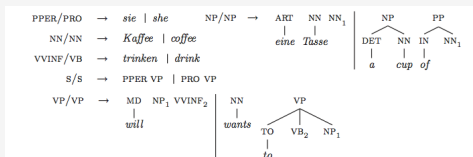
PPER/PRO	→	<i>sie</i>   <i>she</i>	NP/NP	→	ART	NN	NN <sub>1</sub>		NP	PP		
NN/NN	→	<i>Kaffee</i>   <i>coffee</i>			<i>eine</i>	<i>Tasse</i>			DET	NN	IN	NN <sub>1</sub>
VVINF/VB	→	<i>trinken</i>   <i>drink</i>							<i>a</i>	<i>cup</i>	<i>of</i>	
S/S	→	PPER VP		PRO VP								
VP/VP	→	MD	NP <sub>1</sub>	VVINF <sub>2</sub>	NN	VP			TO	VB <sub>2</sub>	NP <sub>1</sub>	
			<i>will</i>		<i>wants</i>				<i>to</i>			

0

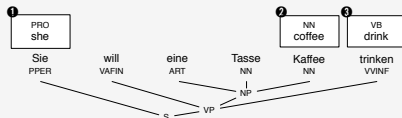


Purely lexical rule: filling a span with a translation (a constituent in the chart)





0

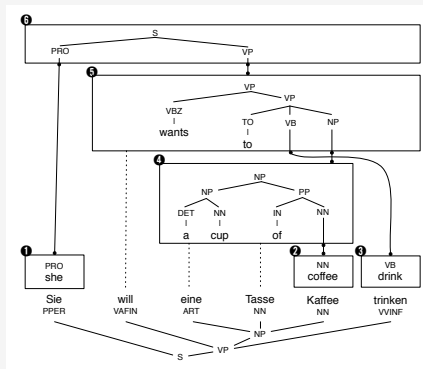


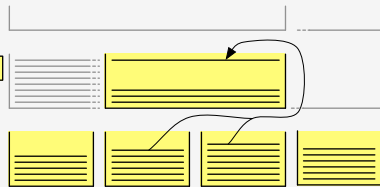
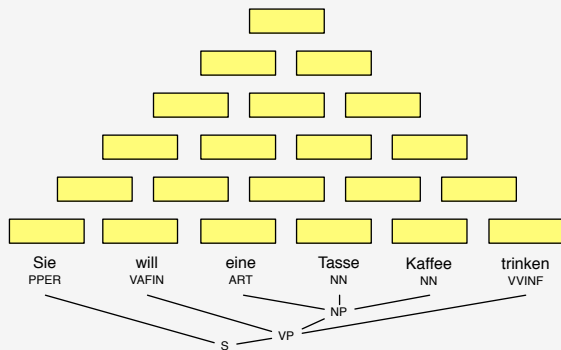
Purely lexical rule: filling a span with a translation (a constituent in the chart)





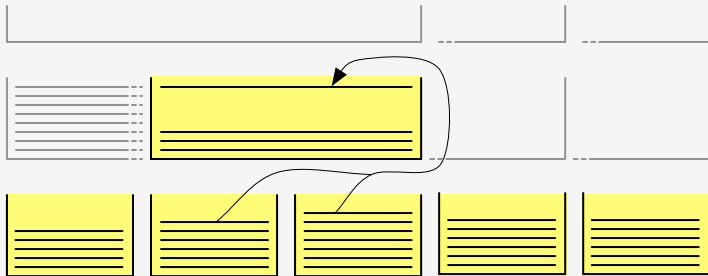
PPER/PRO	→	<i>sie</i>   <i>she</i>	NP/NP	→	ART	NN	NN <sub>1</sub>		NP	PP		
NN/NN	→	<i>Kaffee</i>   <i>coffee</i>							DET	NN	IN	NN <sub>1</sub>
VVINF/VB	→	<i>trinken</i>   <i>drink</i>			<i>eine</i>	<i>Tasse</i>						
S/S	→	PPER VP		PRO VP								
					<i>a</i>	<i>cup</i>	<i>of</i>					
VP/VP	→	MD	NP <sub>1</sub>	VVINF <sub>2</sub>		NN	VP		TO	VB <sub>2</sub>	NP <sub>1</sub>	
		<i>will</i>			<i>wants</i>				<i>to</i>			





- Chart consists of cells that cover contiguous spans over the input sentence
- Each cell contains a set of hypotheses
- Hypothesis = translation of span with target-side constituent

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

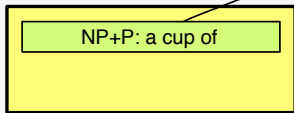
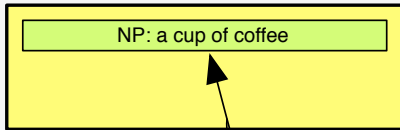


**Input:** Foreign sentence  $\mathbf{f} = f_1, \dots, f_{l_f}$ , with syntax tree

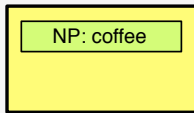
**Output:** English translation  $\mathbf{e}$

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



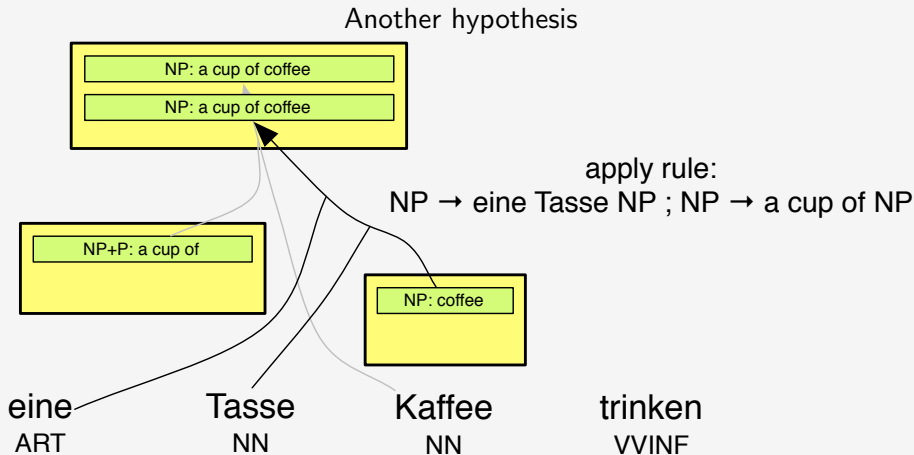
eine  
ART

Tasse  
NN

Kaffee  
NN

trinken  
VVINF





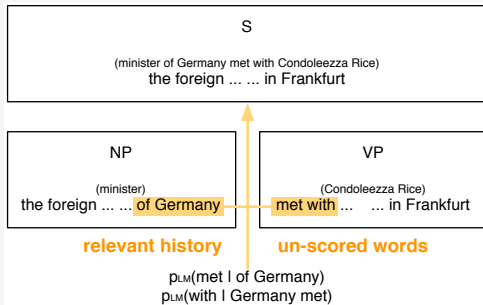
Both hypotheses are indistinguishable in future search  
 → can be recombined

Hypotheses have to match in

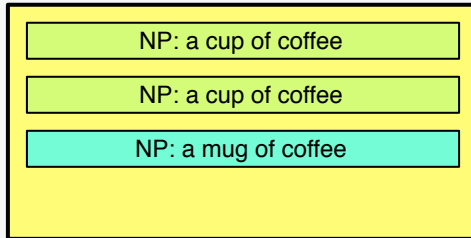
- span of input words covered
- output constituent label
- first  $n-1$  output words (not properly scored, since they lack context)
- last  $n-1$  output words (still affect scoring of subsequently added words, just like in phrase-based decoding)

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

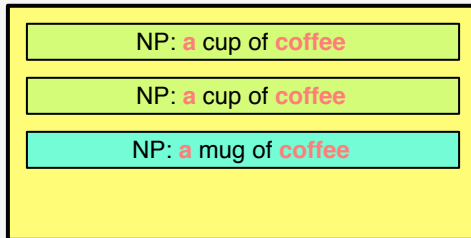
When merging hypotheses, internal language model contexts are absorbed



Recombinable?



Recombinable?



Yes, iff max. 2-gram language model is used