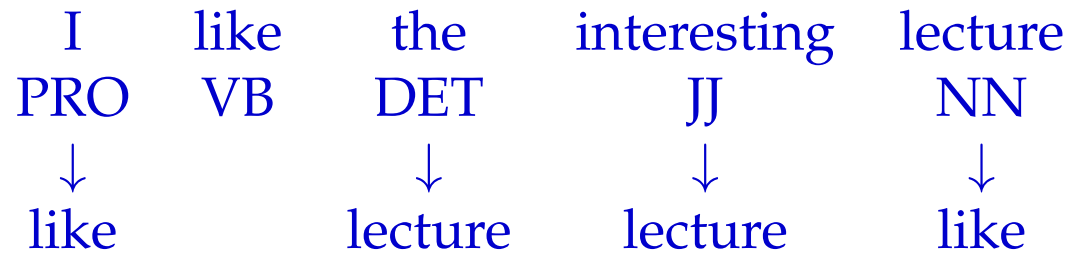# Syntax-Based Models

Philipp Koehn

7 November 2017

# what is syntax?

# Tree-Based Models

- Traditional statistical models operate on sequences of words

- 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
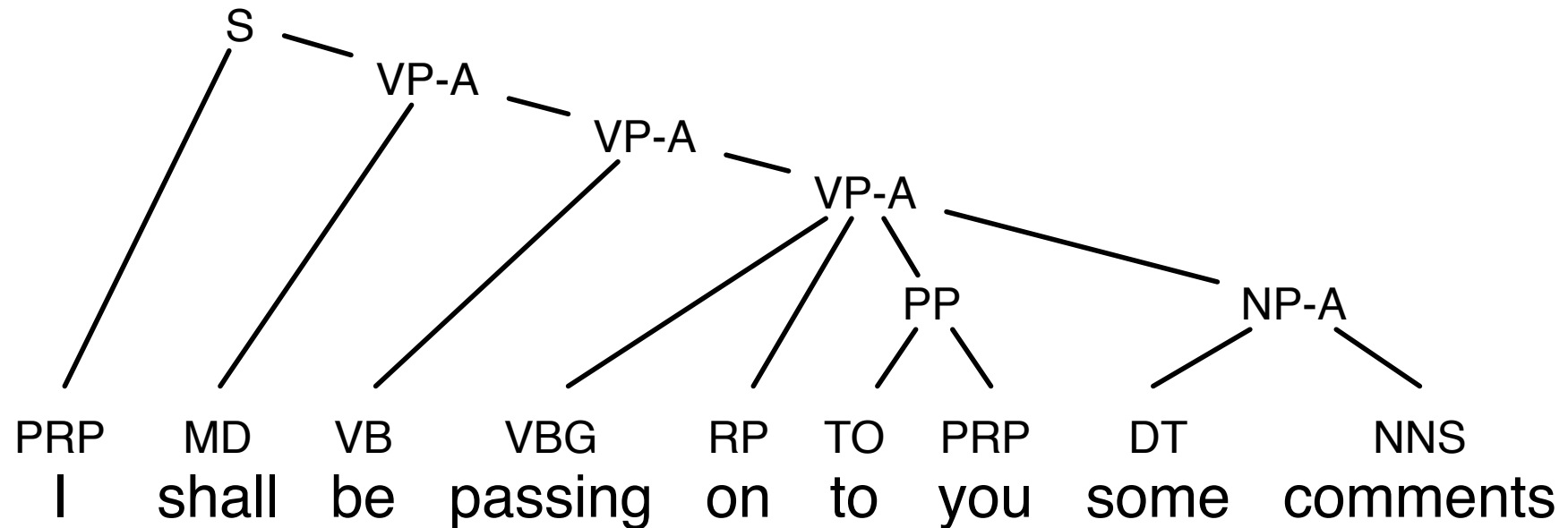  - state-of-the art for some language pairs

# Dependency Structure

$$
\begin{array}{ccccc}
\text{I} & \text{like} & \text{the} & \text{interesting} & \text{lecture} \\
\text{PRO} & \text{VB} & \text{DET} & \text{JJ} & \text{NN} \\
\downarrow & & \downarrow & \downarrow & \downarrow \\
\text{like} & & \text{lecture} & \text{lecture} & \text{like}
\end{array}
$$

- Center of a sentence is the verb

- Its dependents are its arguments (e.g., subject noun)

- These may have further dependents (adjective of noun)

- 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, ...

- Context-free Grammars (CFG)

  – non-terminal symbols: phrase structure labels, part-of-speech tags
  – terminal symbols: words
  – production rules: NT → [NT,T]+
    example: NP → DET NN

# Phrase Structure Grammar

Phrase structure grammar tree for an English sentence
(as produced Collins' parser)

# syntactic transfer

- English rule

$$NP \to DET\ JJ\ NN$$

- French rule

$$NP \to DET\ NN\ JJ$$

- Synchronous rule (indices indicate alignment):

$$NP \to DET_1\ NN_2\ JJ_3 \mid DET_1\ JJ_3\ NN_2$$

# Synchronous Grammar Rules

- Nonterminal rules

$$NP \rightarrow DET_1\ NN_2\ JJ_3\ |\ DET_1\ JJ_3\ NN_2$$

- Terminal rules

$$N \rightarrow maison\ |\ house$$

$$NP \rightarrow la\ maison\ bleue\ |\ the\ blue\ house$$

- Mixed rules

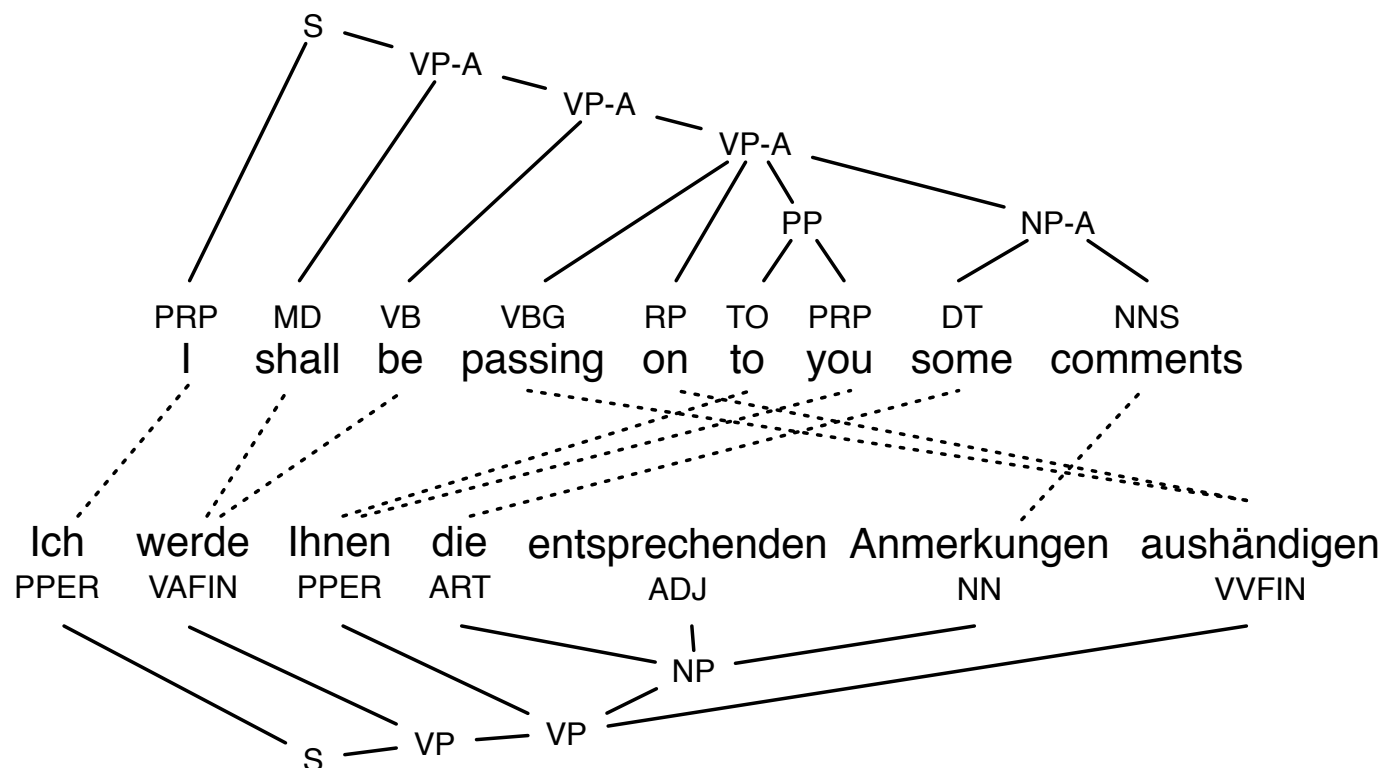$$NP \rightarrow la\ maison\ JJ_1\ |\ the\ JJ_1\ house$$

- 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

- Translation probability

$$\text{SCORE}(\text{TREE}, \text{E}, \text{F}) = \prod_i \text{RULE}_i$$
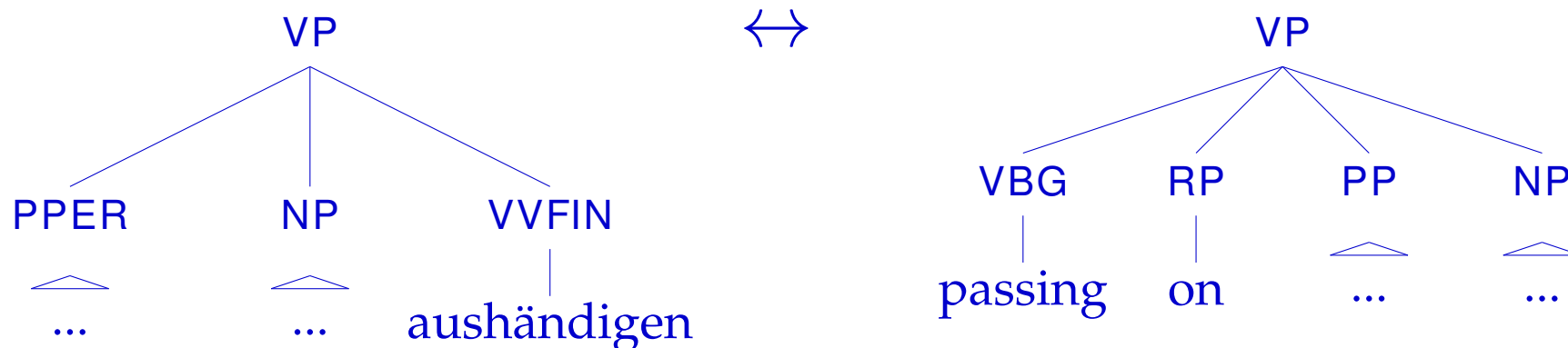
- Many ways to assign probabilities to rules

# Aligned Tree Pair



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

- Subtree alignment



- Synchronous grammar rule

$$VP \rightarrow PPER_1\ NP_2\ aushändigen\ \mid\ passing\ on\ PP_1\ NP_2$$

- Note:

  - one word aushändigen mapped to two words passing on ok
  - but: fully non-terminal rule not possible
    (one-to-one mapping constraint for nonterminals)

# Another Rule

- Subtree alignment

$$PRO \quad\quad \leftrightarrow \quad\quad PP$$

Ihnen

TO     PRP

to      you

- Synchronous grammar rule (stripping out English internal structure)

$$PRO/PP \rightarrow Ihnen \mid to \; you$$

- Rule with internal structure

$$PRO/PP \rightarrow \quad Ihnen \quad\Big| \quad TO \quad PRP$$

to      you

# Another Rule

- Translation of German werde to English shall be

$$VP \quad \leftrightarrow \quad VP$$

VAFIN   VP

werde   ...

MD   VP

shall   VB   VP

be   ...

- Translation rule needs to include mapping of VP

$\Rightarrow$ Complex rule

$$VP \rightarrow \begin{array}{c} VAFIN \quad VP_1 \\ | \\ werde \end{array} \quad \bigg| \quad \begin{array}{c} MD \quad VP \\ | \\ shall \quad VB \quad VP_1 \\ | \\ be \end{array}$$

- Stripping out internal structure

$$\text{VP} \rightarrow \text{werde VP}_1 \mid \text{shall be VP}_1$$

$\Rightarrow$ synchronous context free grammar

- Maintaining internal structure


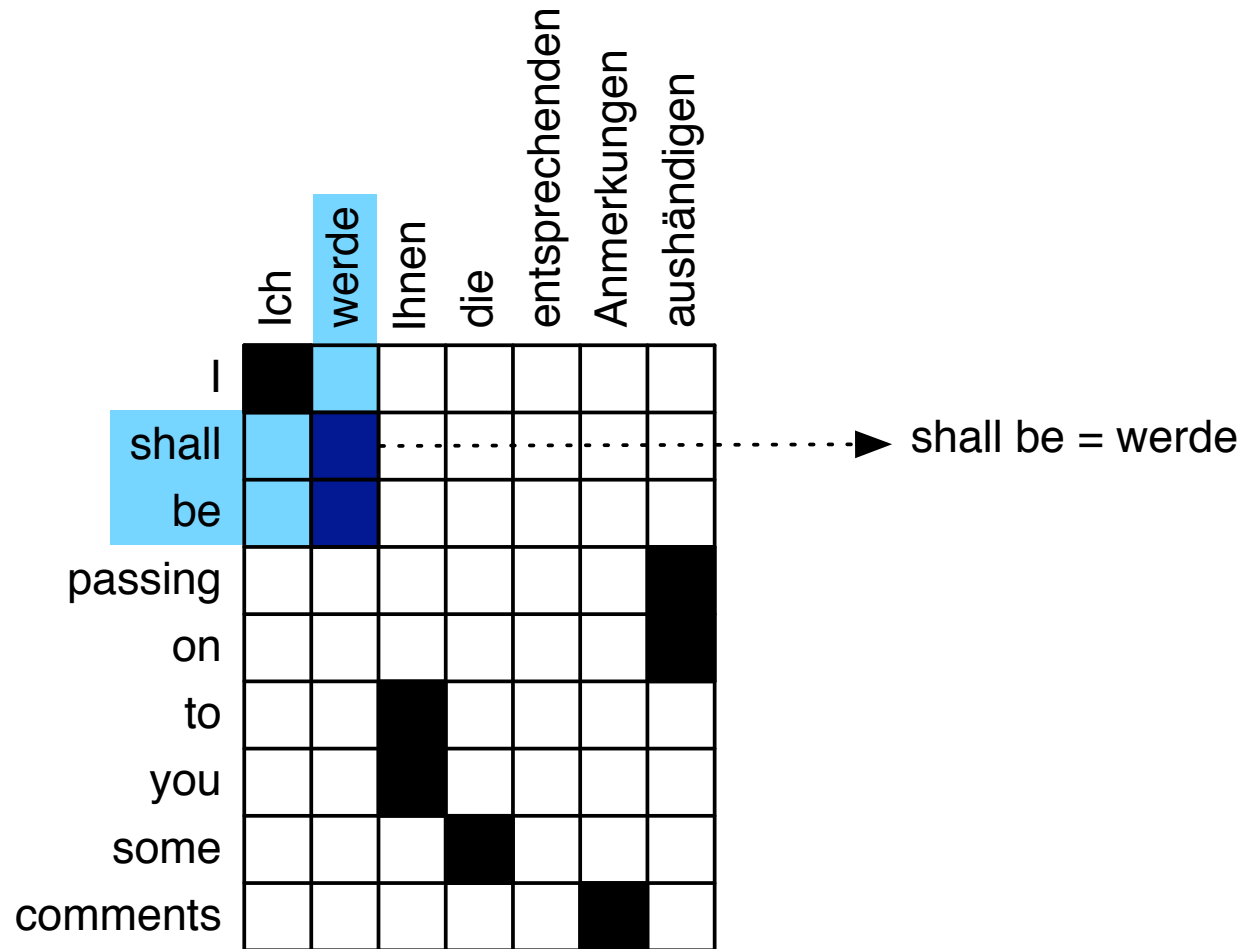
$\Rightarrow$ synchronous tree substitution grammar

# learning

Machine Translation: Syntax-Based Models

- 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 also be parsed with syntactic parser
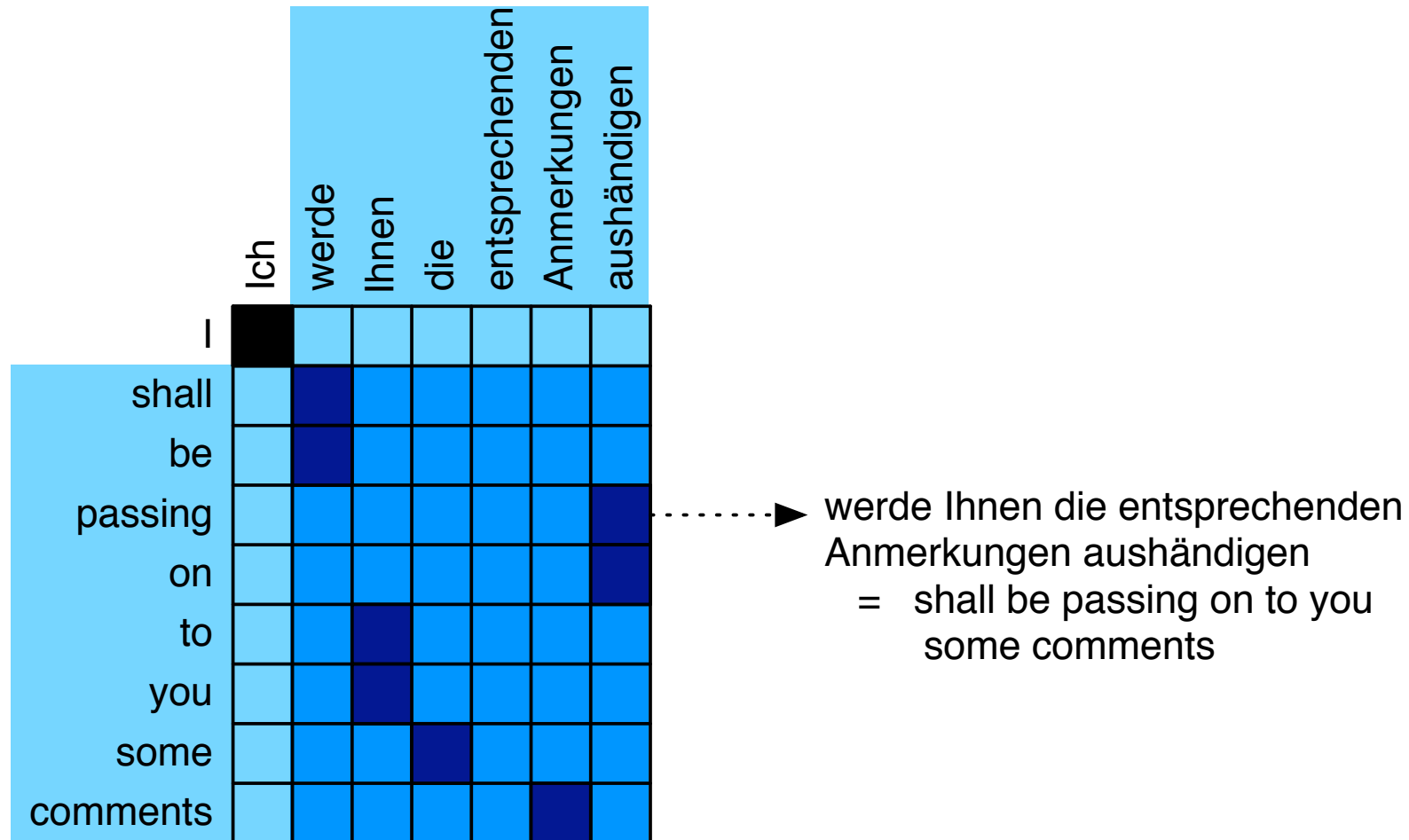
# Extracting Phrase Translation Rules



shall be = werde

some comments =
die entsprechenden Anmerkungen

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

subtracting subphrase

werde X aushändigen
= shall be passing on X

- Recall: consistent phrase pairs

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

$$\forall e_i \in \bar{e} : (e_i, f_j) \in A \rightarrow f_j \in \bar{f}$$

$$\text{AND } \forall f_j \in \bar{f} : (e_i, f_j) \in A \rightarrow e_i \in \bar{e}$$

$$\text{AND } \exists e_i \in \bar{e}, f_j \in \bar{f} : (e_i, f_j) \in A$$

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

# Formal Definition

- Extend recursively:

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

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

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

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

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

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

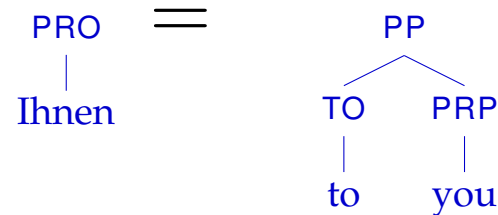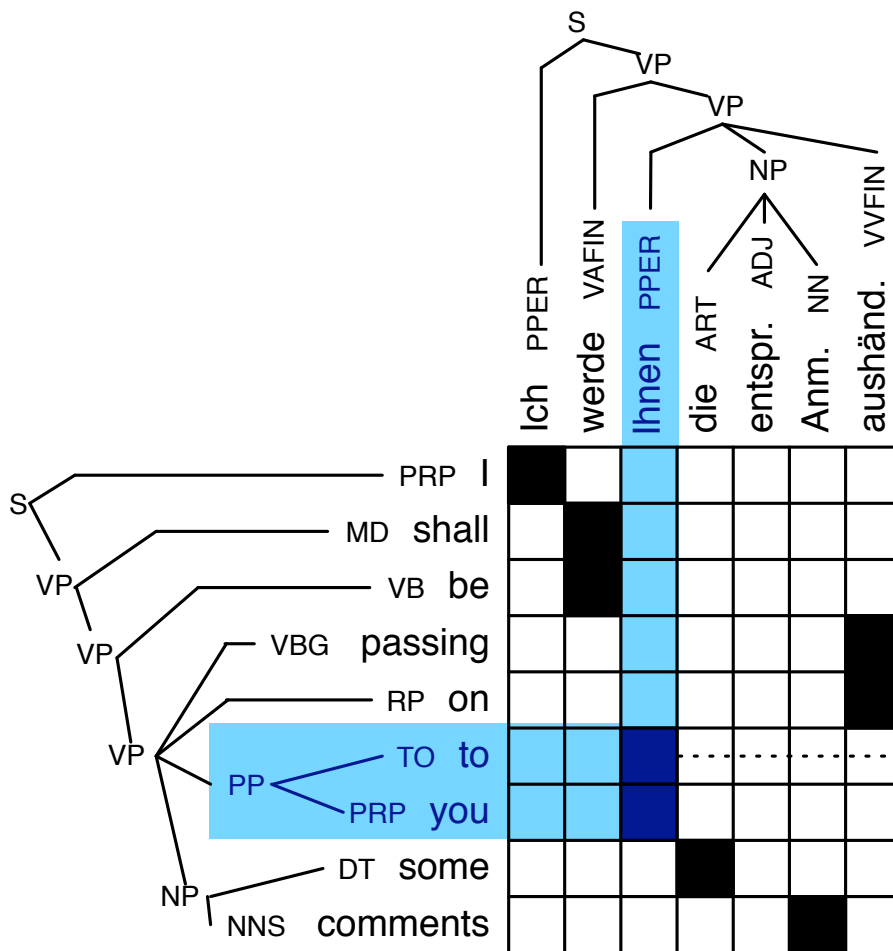- Set of hierarchical phrase pairs is the closure under this extension mechanism

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

$$Y \rightarrow X_1 \ X_2 \ \mid \ X_2 \ of \ X_1$$

- Typical restrictions to limit complexity (Chiang, 2005)

  - at most 2 nonterminal symbols
  - no neighboring non-terminals on the source side
  - at least 1 but at most 5 words per language
  - span at most 15 words (counting gaps)

# Constraints on Syntactic Rules

- Same word alignment constraints as hierarchical models

- Hierarchical: rule can cover any span
  $\Leftrightarrow$ syntactic rules must cover constituents in the tree

- Hierarchical: gaps may cover any span
  $\Leftrightarrow$ gaps must cover constituents in the tree
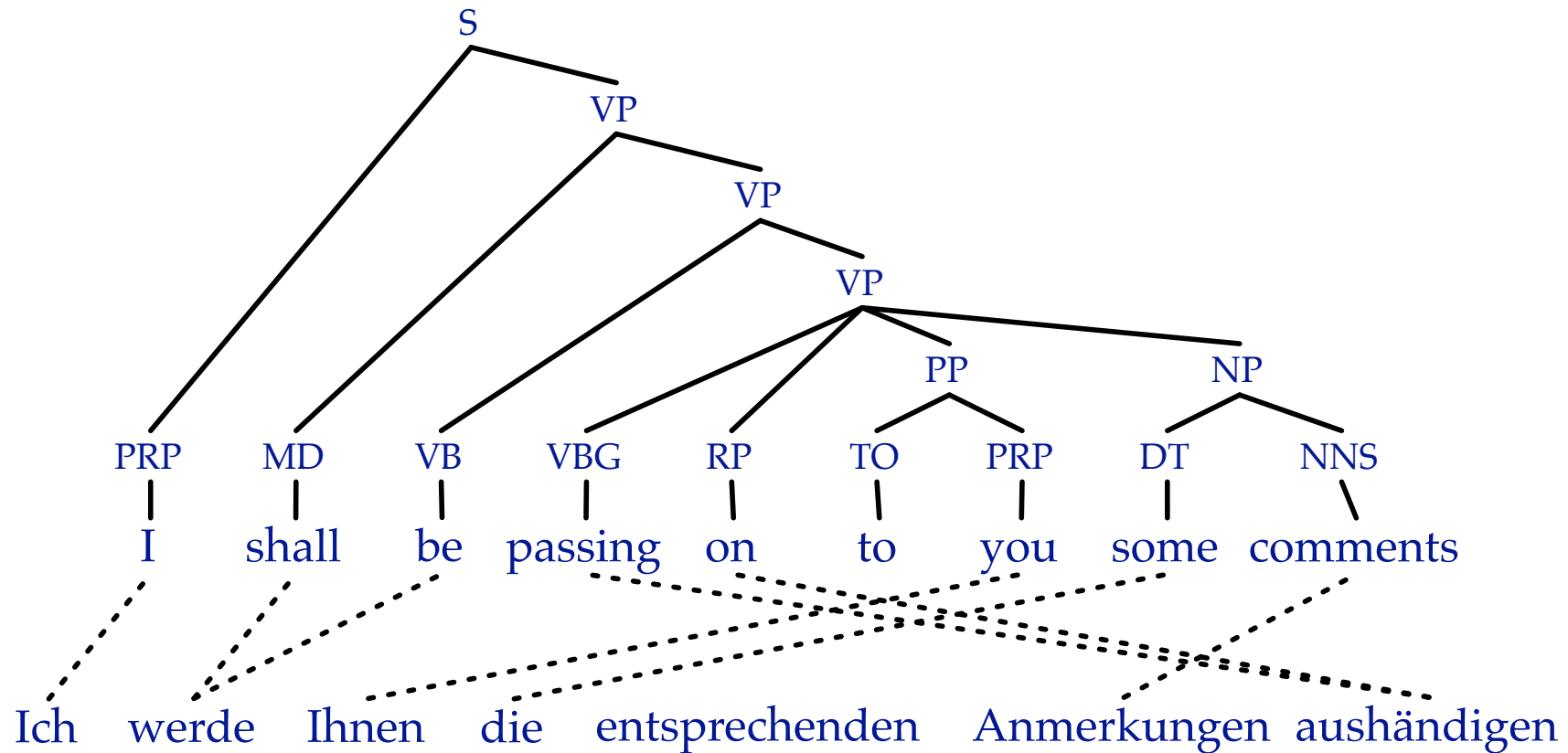
- Much less rules are extracted (all things being equal)

# Impossible Rules



English span not a constituent
no rule extracted

Rule with this phrase pair requires syntactic context

# Too Many Rules Extractable

- Huge number of rules can be extracted
  (every alignable node may or may not be part of a rule $\rightarrow$ exponential number of rules)

- Need to limit which rules to extract

- Option 1: similar restriction as for hierarchical model
  (maximum span size, maximum number of terminals and non-terminals, etc.)

- Option 2: only extract minimal rules ("GHKM" rules)

# refinements

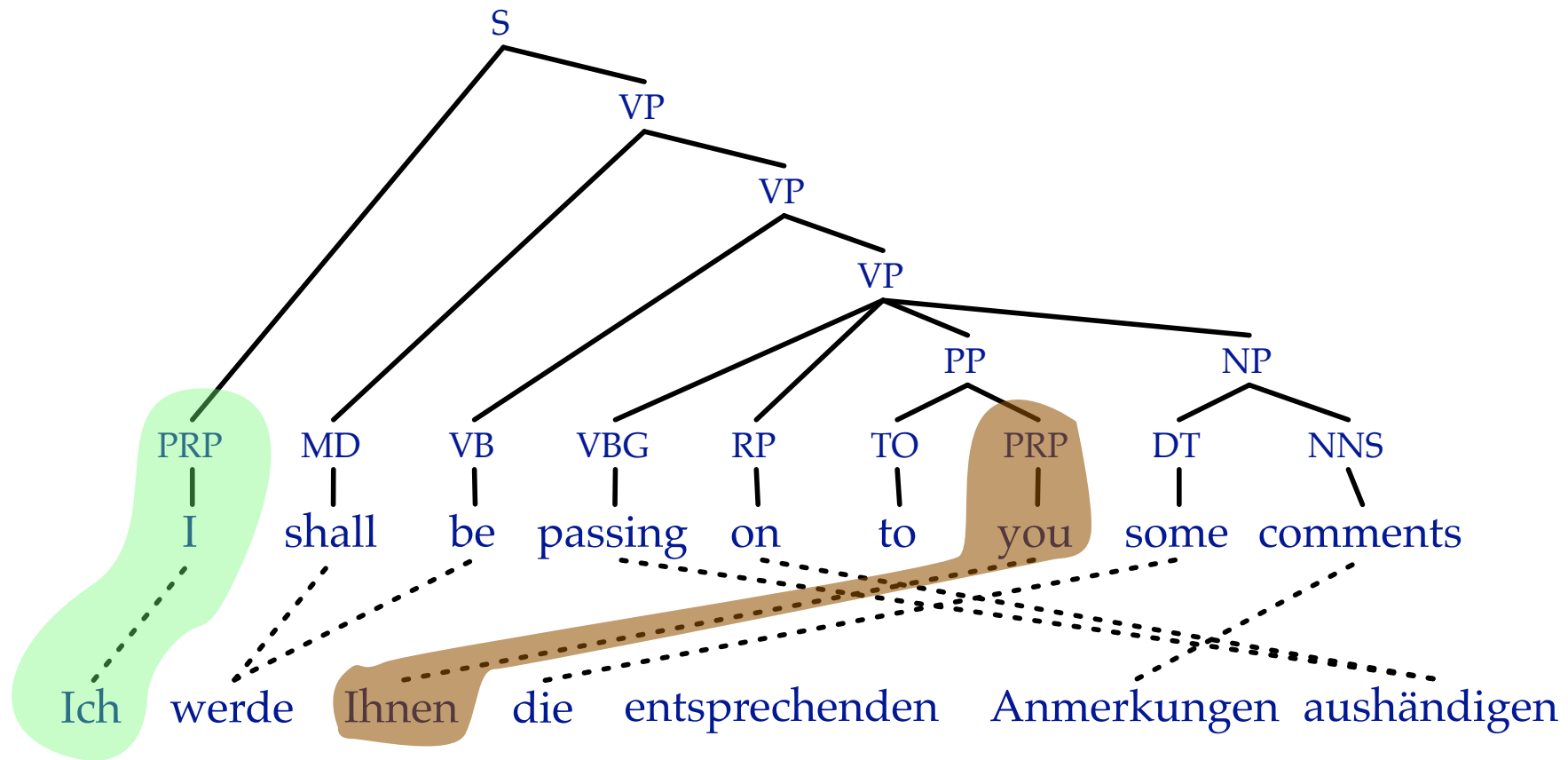Extract: set of smallest rules required to explain the sentence pair

# Lexical Rule

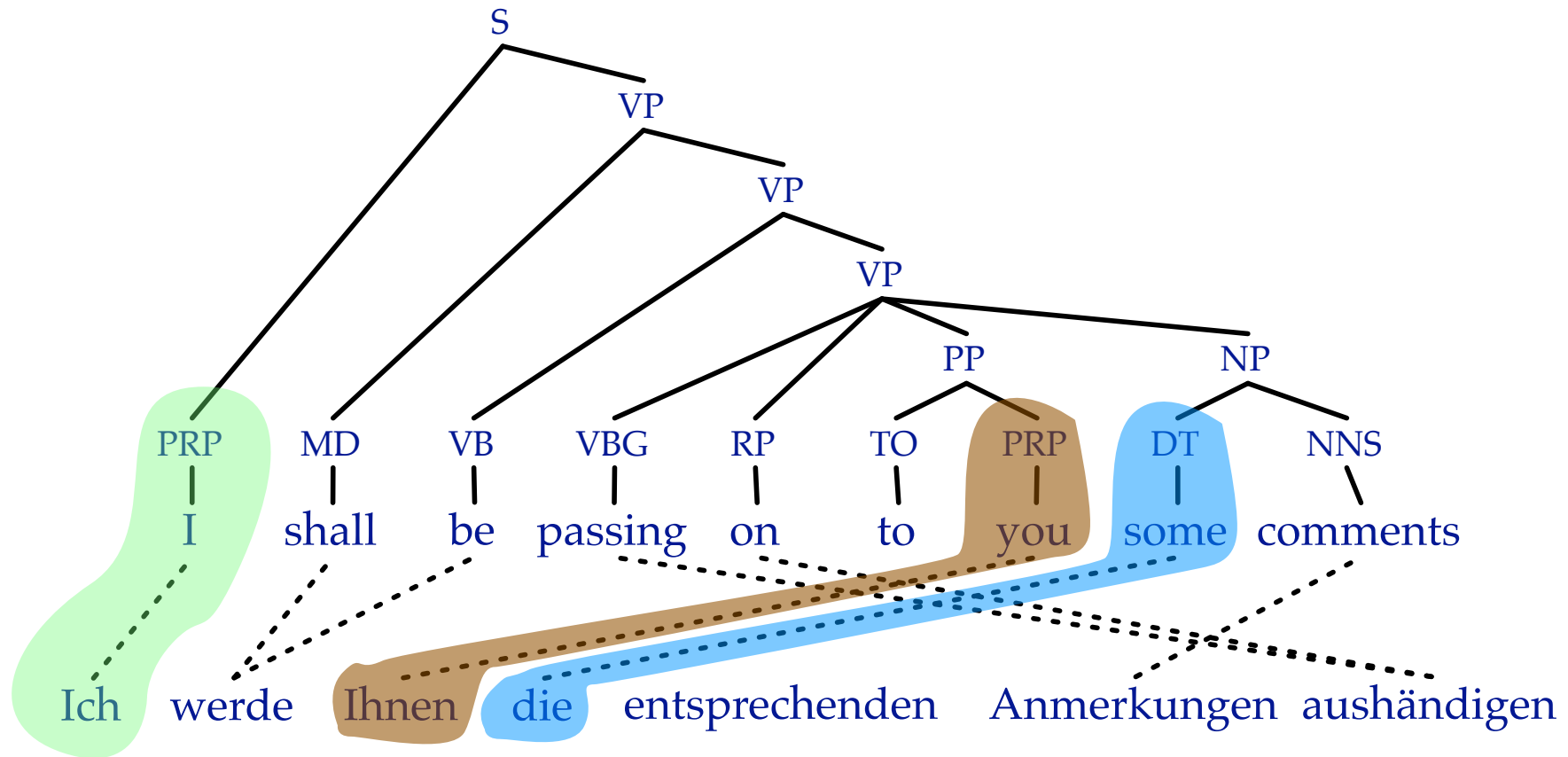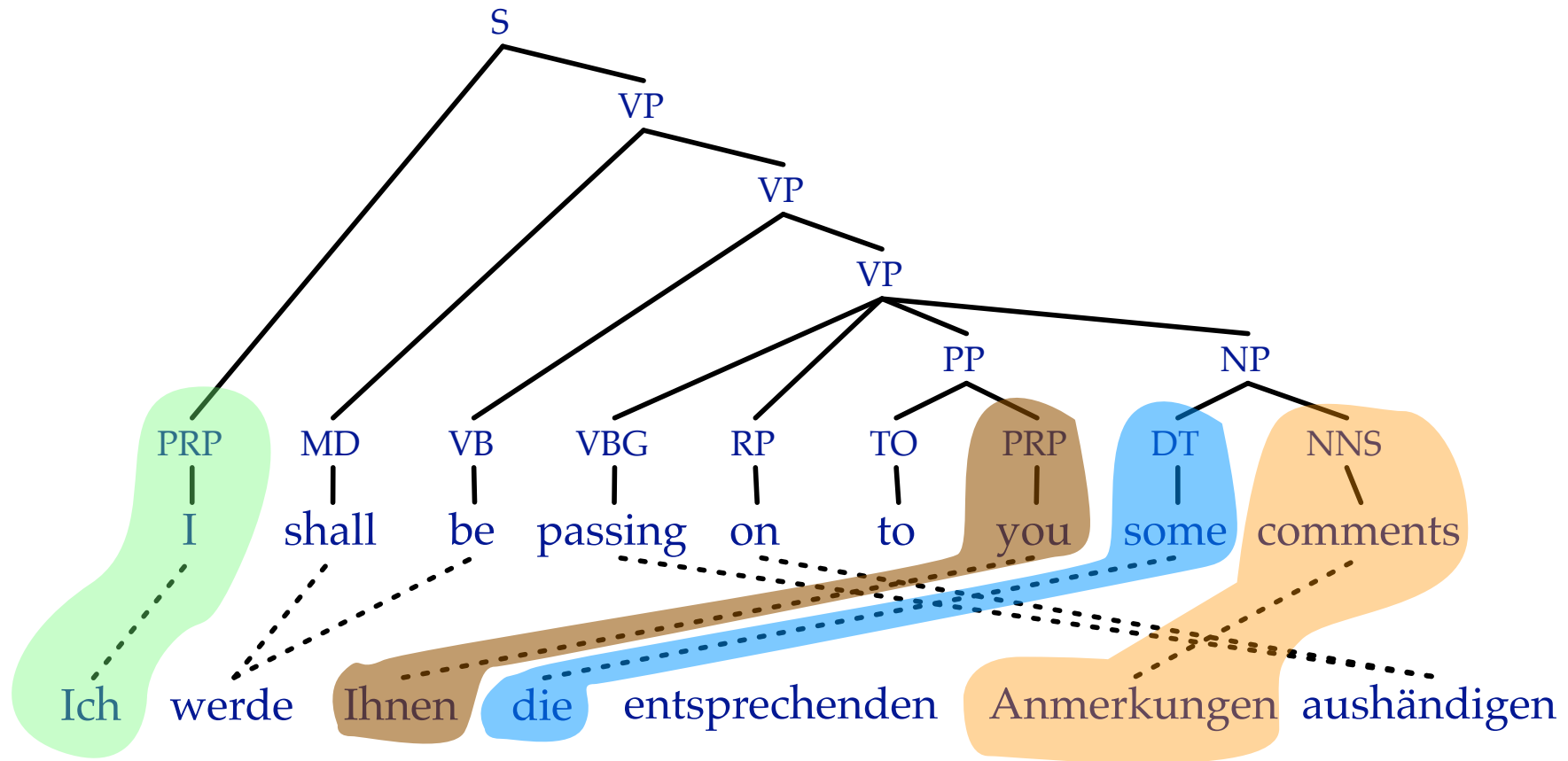

Extracted rule: PRP → Ich | I
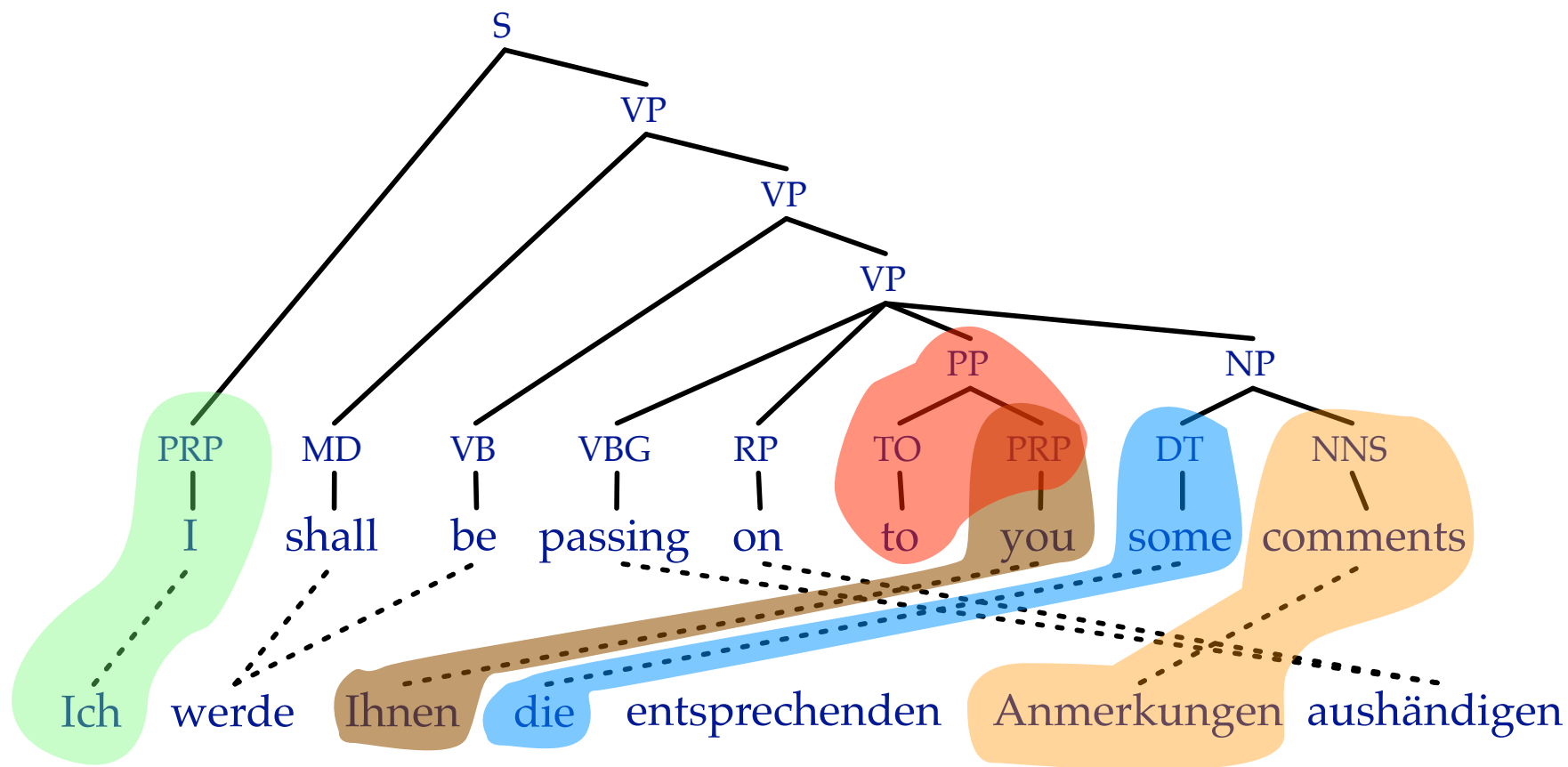
# Lexical Rule

Extracted rule: PRP → Ihnen | you

# Lexical Rule



Extracted rule: DT → die | some

# Lexical Rule
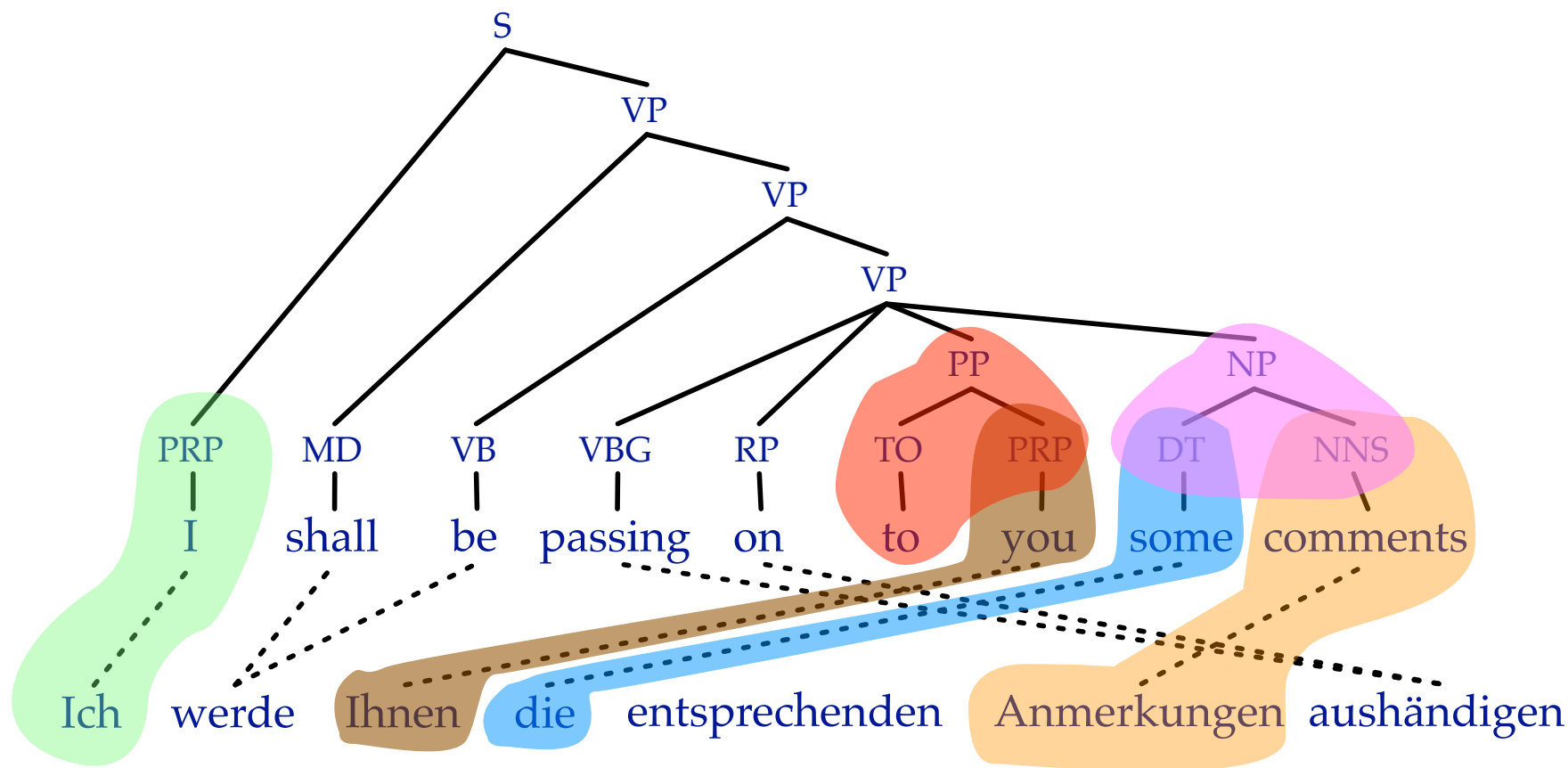


Extracted rule: NNS → Anmerkungen | comments

# Insertion Rule
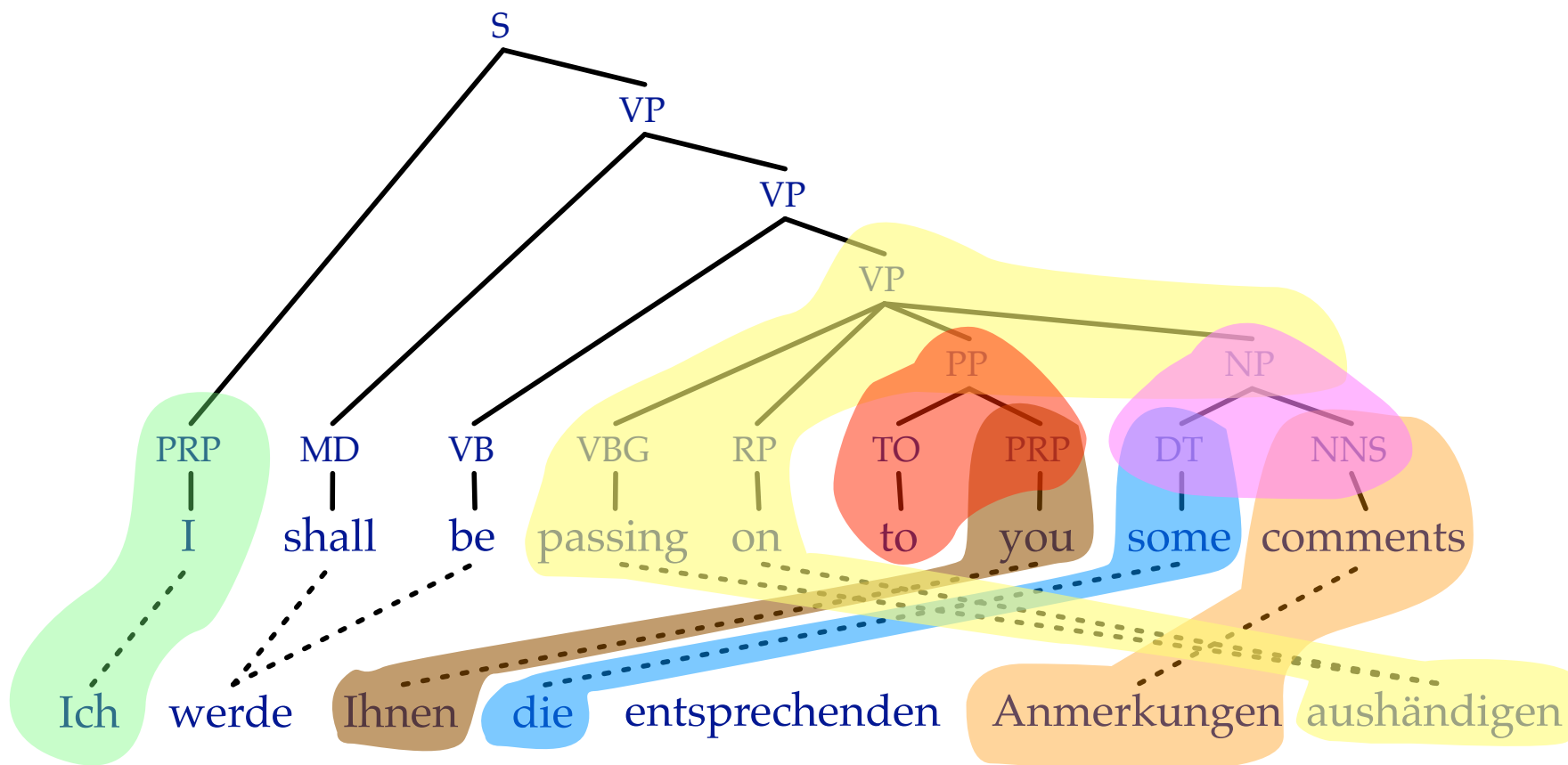
Extracted rule: PP → X | to PRP

# Non-Lexical Rule

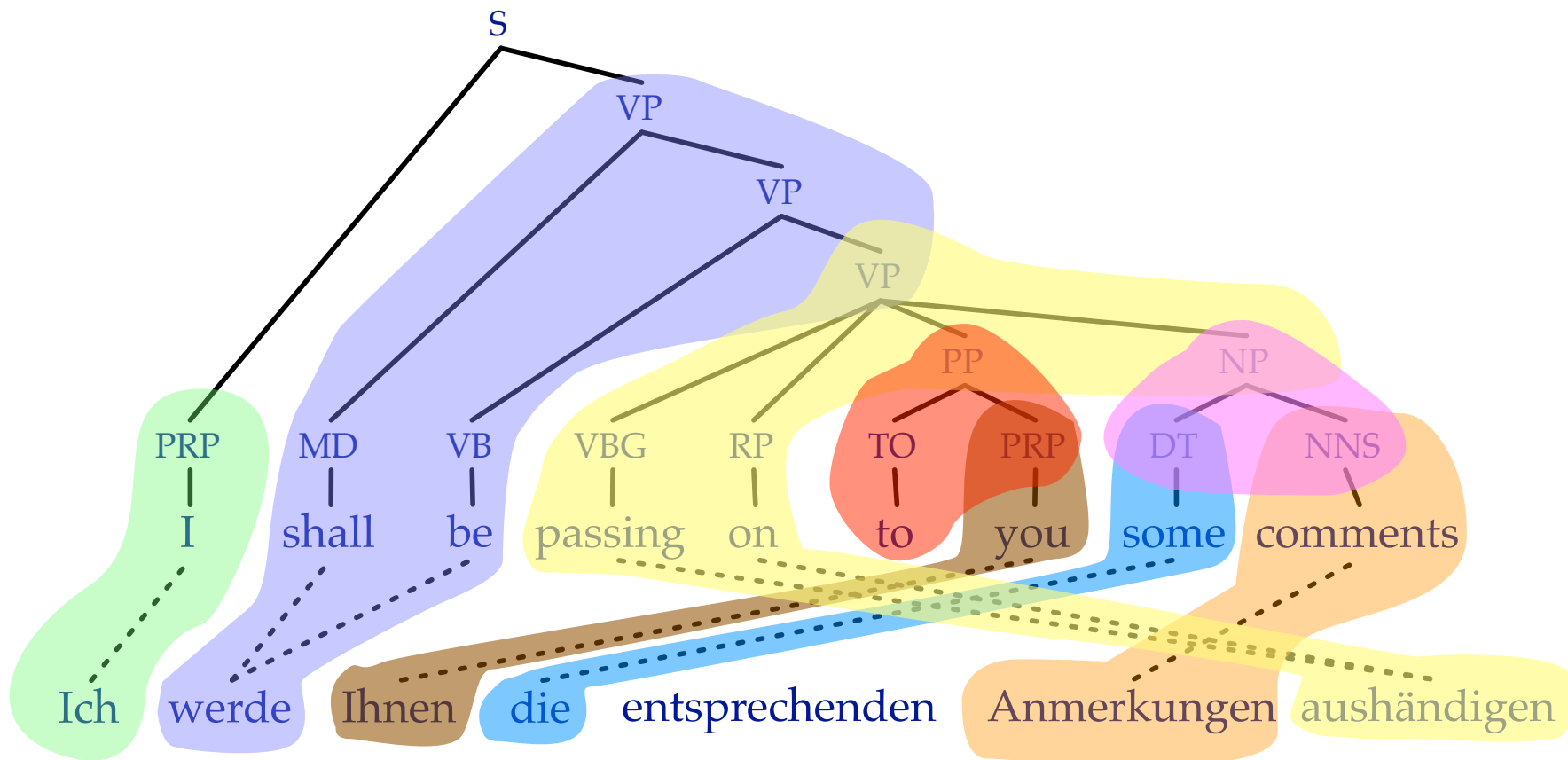Extracted rule: NP $\rightarrow$ X$_1$ X$_2$ | DT$_1$ NNS$_2$

# Lexical Rule with Syntactic Context



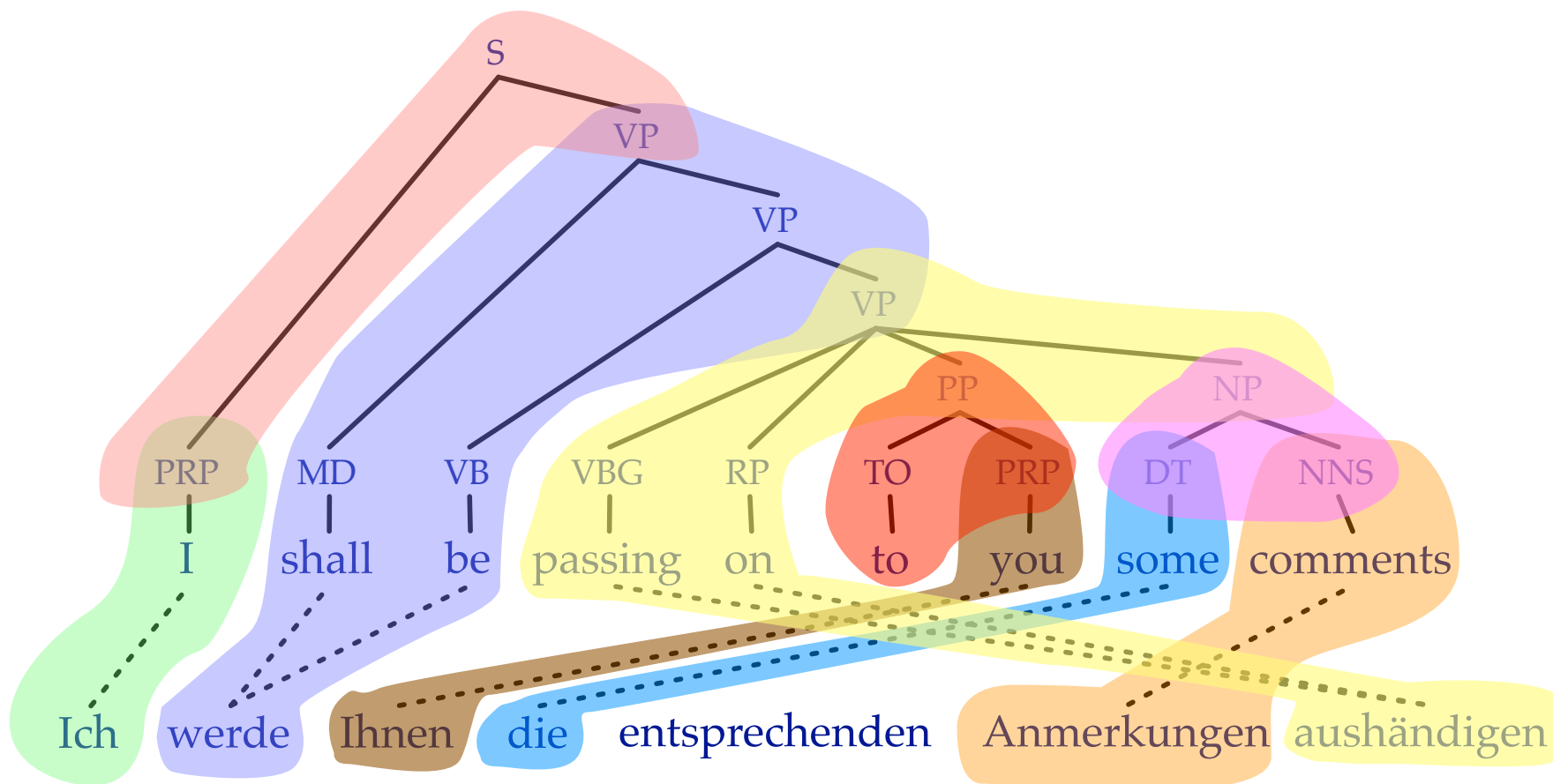Extracted rule: VP → $X_1$ $X_2$ aushändigen | passing on $PP_1$ $NP_2$

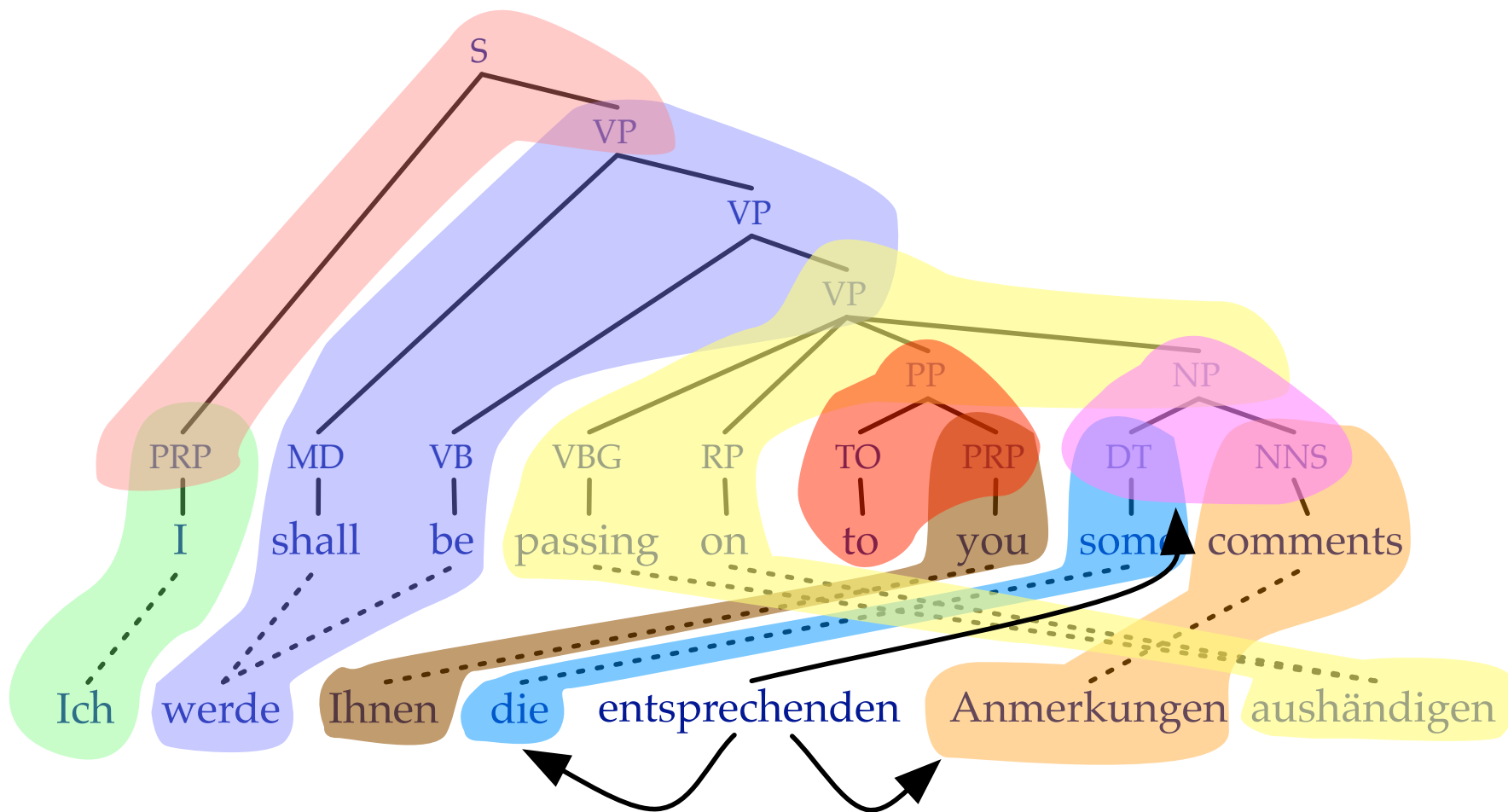Extracted rule: VP → werde X | shall be VP (ignoring internal structure)

Extracted rule: $S \rightarrow X_1 \; X_2 \mid PRP_1 \; VP_2$
DONE — note: one rule per alignable constituent

# Unaligned Source Words

Attach to neighboring words or higher nodes → additional rules
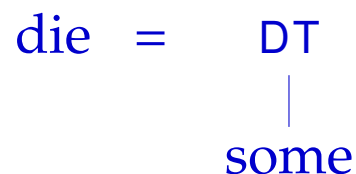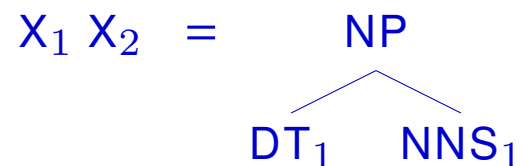
# Too Few Phrasal Rules?

- Lexical rules will be 1-to-1 mappings (unless word alignment requires otherwise)

- But: phrasal rules very beneficial in phrase-based models

- Solutions

  - combine rules that contain a maximum number of symbols
    (as in hierarchical models, recall: "Option 1")

  - compose minimal rules to cover a maximum number of non-leaf nodes

# Composed Rules

- Current rules

$$X_1\ X_2\ =\ NP$$
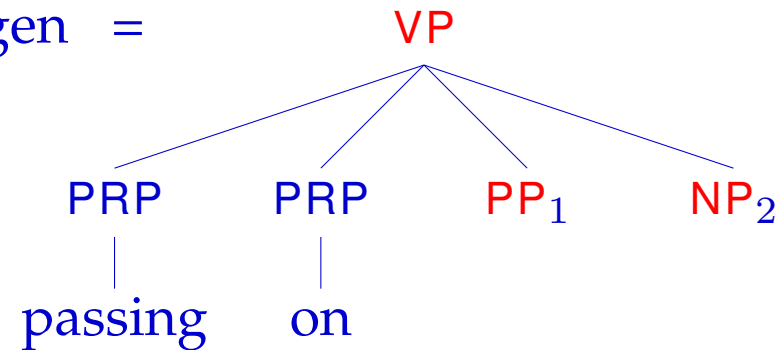
NP branches to $DT_1$ and $NNS_1$

die = DT — some

entsprechenden Anmerkungen = NNS — comments

- Composed rule

die entsprechenden Anmerkungen = NP

NP branches to DT (some) and NNS (comments)

(1 non-leaf node: NP)

# Composed Rules

- Minimal rule:     $X_1\ X_2$ aushändigen  =

```
                        VP
           /      |      |      \
        PRP     PRP    PP_1    NP_2
         |       |
      passing    on
```

  3 non-leaf nodes:
  VP, PP, NP

- Composed rule:     Ihnen $X_1$ aushändigen  =

```
                         VP
           /       |      |       \
        PRP      PRP     PP       NP_1
         |        |      /  \
      passing     on   TO   PRP
                       |     |
                       to   you
```

  3 non-leaf nodes:
  VP, PP and NP

- Impossible rule

$$X \quad = \quad MD \quad VB$$

| | | | |
| werde | | shall | be |

- Create new non-terminal label: MD+VB

⇒ New rule

$$X \quad = \quad MD+VB$$

werde

MD    VB

shall    be

# Zollmann Venugopal Relaxation

- If span consists of two constituents , join them: X+Y
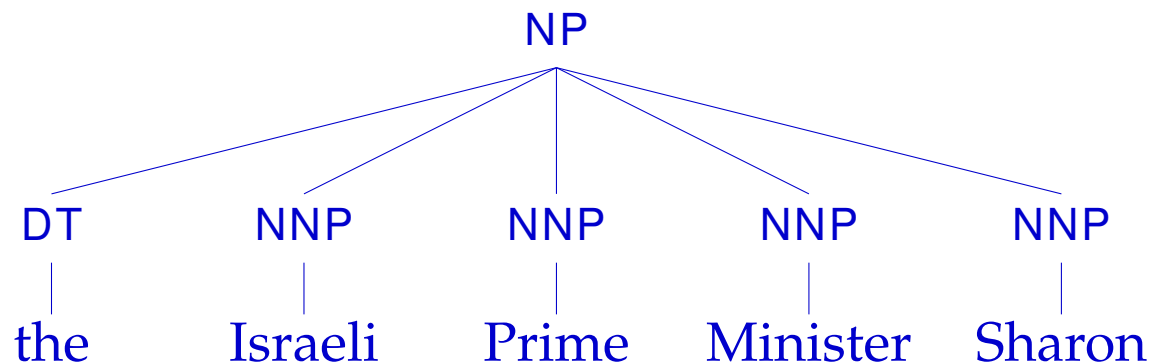
- If span conststs of three constituents, join them: X+Y+Z

- If span covers constituents with the same parent x and include
  - every but the first child Y, label as X\Y
  - every but the last child Y, label as X/Y

- For all other cases, label as FAIL

$\Rightarrow$ More rules can be extracted, but number of non-terminals blows up
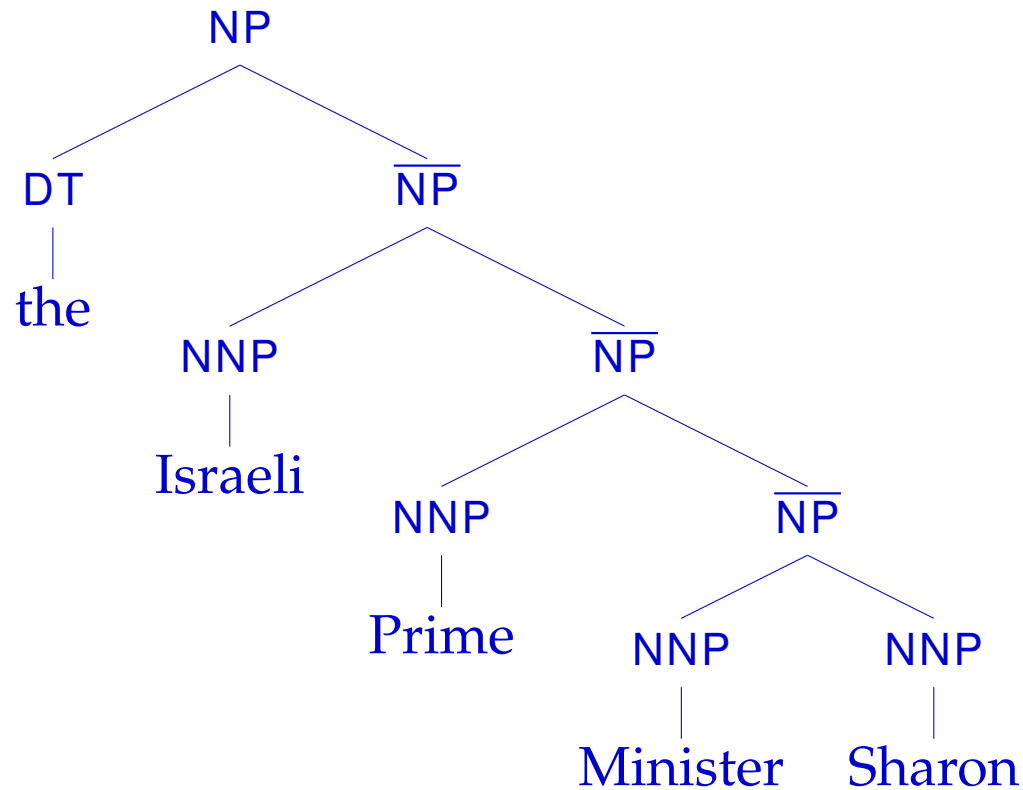
# Special Problem: Flat Structures

- Flat structures severely limit rule extraction

```
                          NP
         ┌──────────┬──────┼──────┬──────────┐
        DT         NNP    NNP    NNP        NNP
         │          │      │      │          │
        the       Israeli Prime Minister   Sharon
```

- Can only extract rules for individual words or entire phrase

More rules can be extracted

Left-binarization or right-binarization?

- Extract all rules from corpus

- Score based on counts

  - joint rule probability: $p(\mathsf{LHS}, \mathsf{RHS}_f, \mathsf{RHS}_e)$
  - rule application probability: $p(\mathsf{RHS}_f, \mathsf{RHS}_e | \mathsf{LHS})$
  - direct translation probability: $p(\mathsf{RHS}_e | \mathsf{RHS}_f, \mathsf{LHS})$
  - noisy channel translation probability: $p(\mathsf{RHS}_f | \mathsf{RHS}_e, \mathsf{LHS})$
  - lexical translation probability: $\prod_{e_i \in \mathsf{RHS}_e} p(e_i | \mathsf{RHS}_f, a)$

# next lecture: decoding