# Syntax-Based Decoding 2

Philipp Koehn

14 November 2017

# flashback: syntax-based models

- Nonterminal rules

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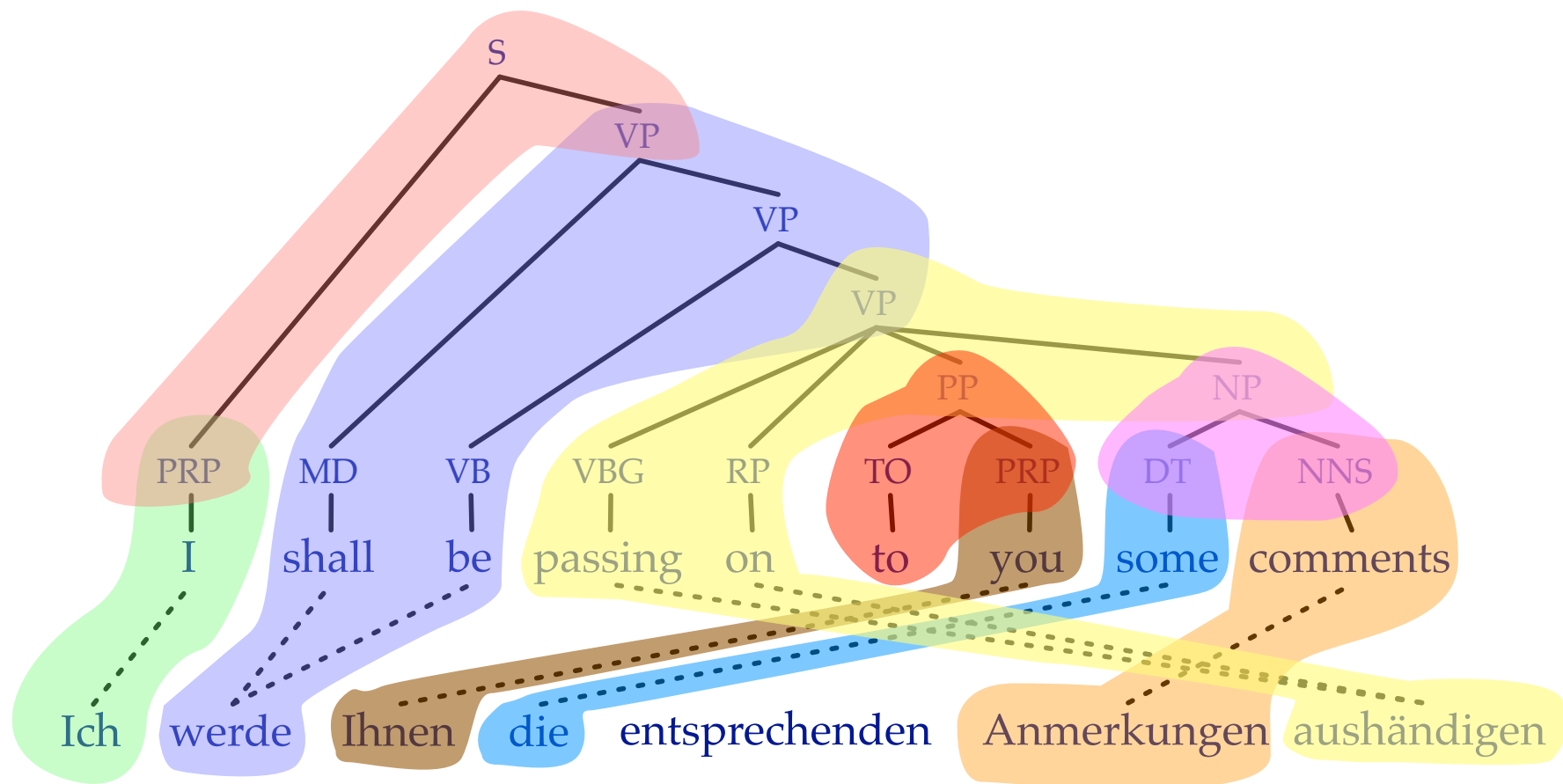
- Terminal rules

$$\text{N} \rightarrow \text{maison} \mid \text{house}$$

$$\text{NP} \rightarrow \text{la maison bleue} \mid \text{the blue house}$$

- Mixed rules

$$\text{NP} \rightarrow \text{la maison JJ}_1 \mid \text{the JJ}_1 \text{ house}$$
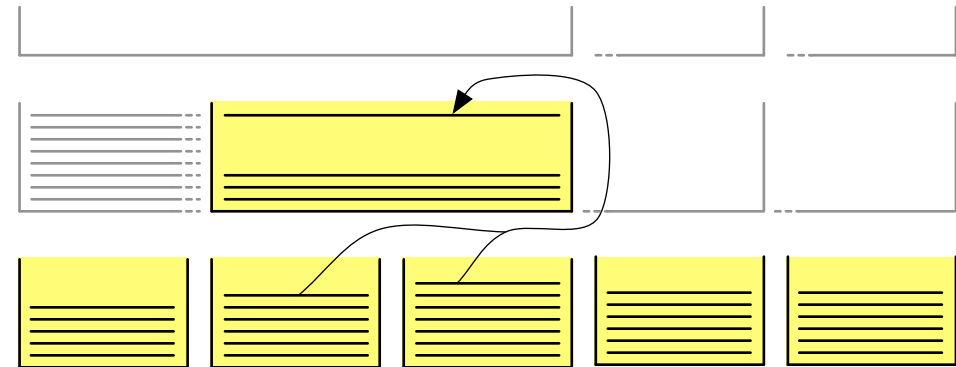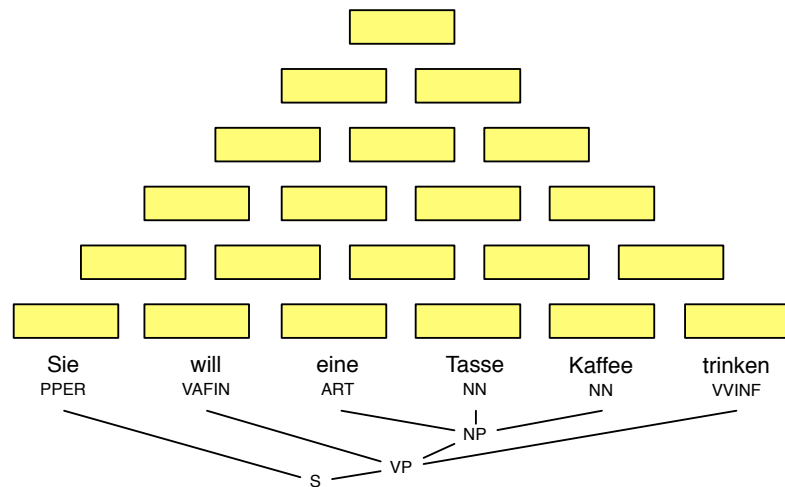
# Extracting Minimal Rules

Extracted rule: $S \rightarrow X_1 \, X_2 \mid PRP_1 \, VP_2$

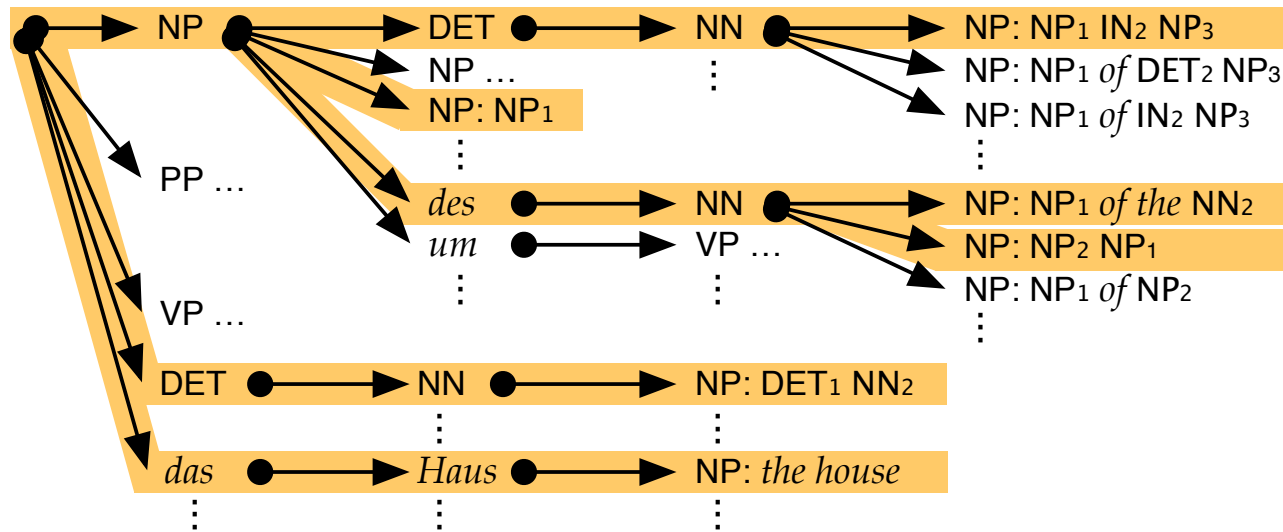DONE — note: one rule per alignable constituent

# flashback: decoding

# Chart Organization

- Chart consists of cells that cover contiguous spans over the input sentence

- For each span, a stack of (partial) translations is maintained

- Bottom-up: a higher stack is filled, once underlying stacks are complete

# Prefix Tree for Rules

## Highlighted Rules

$NP \rightarrow NP_1\ DET_2\ NN_3 \mid NP_1\ IN_2\ NN_3$

$NP \rightarrow NP_1 \mid NP_1$

$NP \rightarrow NP_1\ des\ NN_2 \mid NP_1\ of\ the\ NN_2$

$NP \rightarrow NP_1\ des\ NN_2 \mid NP_2\ NP_1$

$NP \rightarrow DET_1\ NN_2 \mid DET_1\ NN_2$
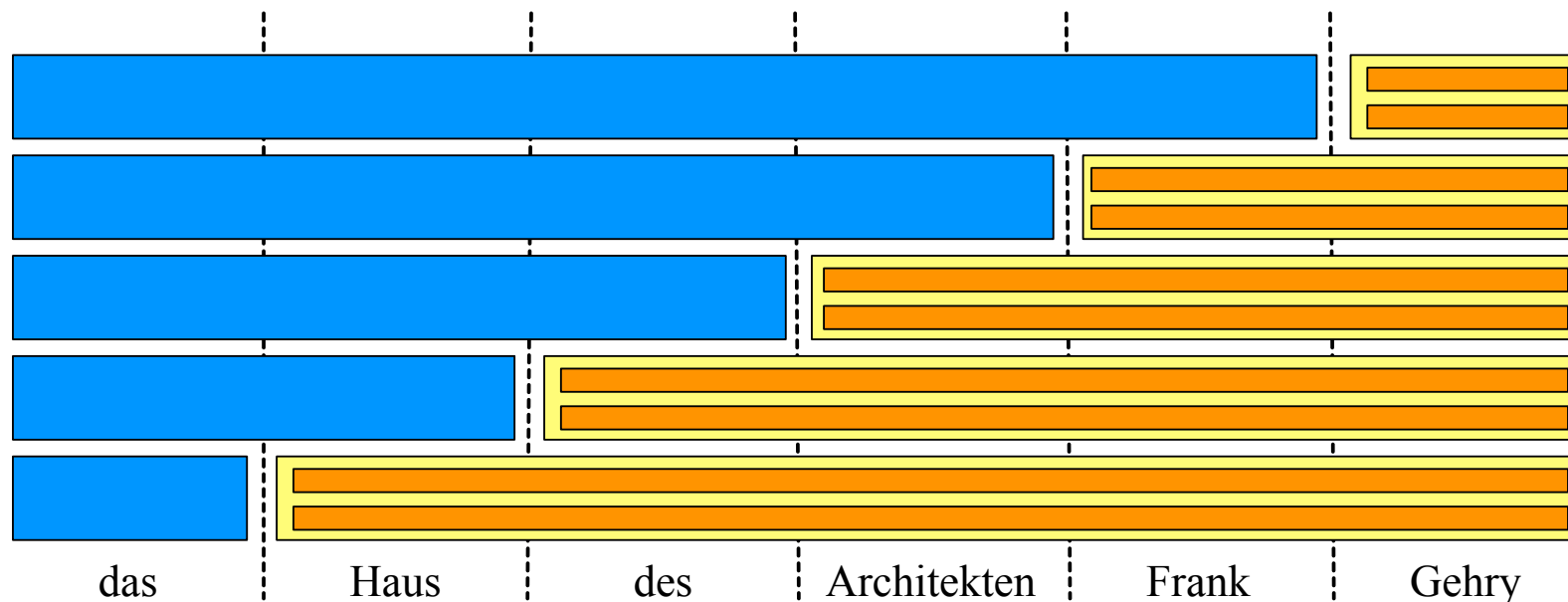
$NP \rightarrow das\ Haus \mid the\ house$

# CYK+ Parsing for SCFG

Extend lists of dotted rules with cell constituent labels

span's dotted rule list (with same start)
plus neighboring
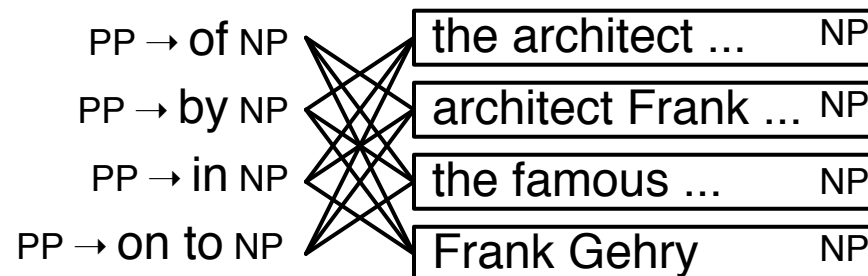span's constituent labels of hypotheses (with same end)

# pruning

# Where are we now?

- 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

# Rules with One Non-Terminal

Found applicable rules PP → des X | ... NP ...



- Non-terminal will be filled any of $h$ underlying matching hypotheses

- Choice of $t$ lexical translations

⇒ Complexity $O(ht)$

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

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



| a house |
|---|
| a building |
| the building |
| a new house |

NP → NP of NP
NP → NP by NP
NP → NP in NP
NP → NP on to NP

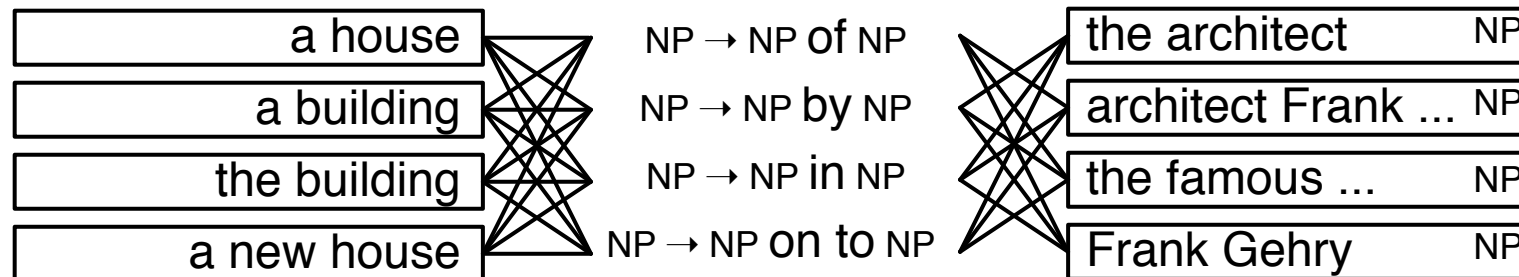| the architect | NP |
|---|---|
| architect Frank ... | NP |
| the famous ... | NP |
| Frank Gehry | NP |

- Two non-terminal will be filled any of $h$ underlying matching hypotheses each

- Choice of $t$ lexical translations

$\Rightarrow$ Complexity $O(h^2 t)$ — 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)

- Hypotheses created in cube: (0,0)

|  | 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)

- Hypotheses created in cube: (0,1), (1,0)

- Hypotheses in chart cell stack: (0,0)

- Hypotheses created in cube: (0,1)

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

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

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

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

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

# Queue of Cubes

- Several groups of rules will apply to a given span

- Each of them will have a cube

- We can create a queue of cubes

$\Rightarrow$ 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**

# recombination and pruning

# Dynamic Programming

Applying rule creates new hypothesis



apply rule:
NP → NP Kaffee | NP+P coffee

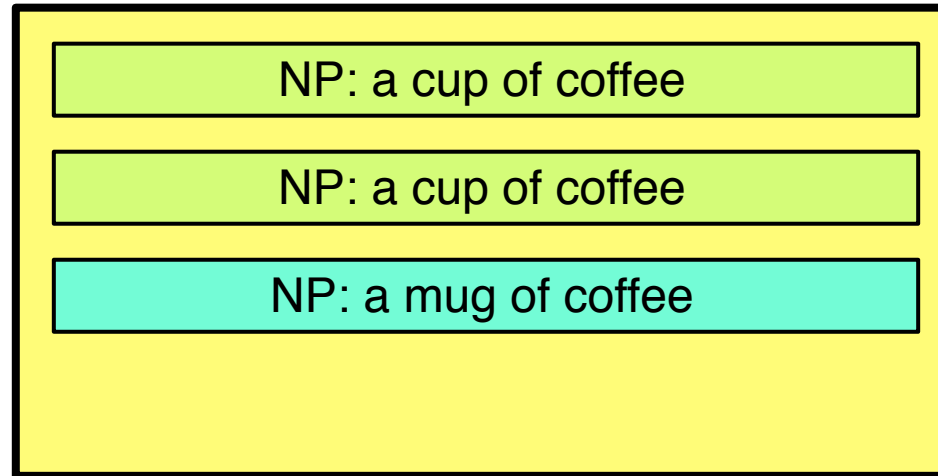# Dynamic Programming

Another hypothesis



Both hypotheses are indistiguishable in future search
$\rightarrow$ can be recombined

Recombinable?

# Recombinable States

Recombinable?

NP: **a** cup of **coffee**

NP: **a** cup of **coffee**

NP: **a** mug of **coffee**

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

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



S

(minister of Germany met with Condoleezza Rice)
the foreign ... ... in Frankfurt

NP

(minister)
the foreign ... ... of Germany

VP

(Condoleezza Rice)
met with ... ... in Frankfurt

**relevant history**          **un-scored words**

$p_{LM}$(met | of Germany)
$p_{LM}$(with | Germany met)

- Number of hypotheses in each chart cell explodes

$\Rightarrow$ need to discard bad hypotheses
  e.g., keep 100 best only

- Different stacks for different output constituent labels?

- Cost estimates

  – translation model cost known
  – language model cost for internal words known
    $\rightarrow$ estimates for initial words
  – outside cost estimate?
    (how useful will be a NP covering input words 3–5 later on?)

# scope 3 pruning

- Lexical rule $\rightarrow$ only once in sentence

$$NP \rightarrow \text{la maison bleue} \mid \text{the blue house}$$

- One non-terminal bounded by words $\rightarrow$ only once in sentence

$$NP \rightarrow \text{la } NN_1 \text{ bleue} \mid \text{the blue } NN_1$$

- One non-terminal at edge of rule $\rightarrow$ non-terminal can cover $O(n)$ words

$$NP \rightarrow \text{la } NN_1 \mid \text{the } NN_1$$

- Two non-terminals at edges $\rightarrow$ combined choices for both non-terminals $O(n^2)$

$$NP \rightarrow DET_1 \text{ maison } JJ_2 \mid DET_1 \text{ } JJ_2 \text{ house}$$

- 4 choice points $\rightarrow O(n^4)$ application contexts

- Too many choice points $\rightarrow$ rule applied to many times

- Having only one non-terminal symbol X

- Restrictions to limit complexity

  - at most 2 nonterminal symbols
  - no neighboring non-terminals on the source side
  - span at most 15 words (counting gaps)

$\Rightarrow$ At most 2 choice points ("scope 2")

- Convert grammar to Chomsky Normal Form (CNF) — scope 3

- Only allow two types of rules
$$A \rightarrow word$$
$$A \rightarrow B\ C$$

  (Note: for our rules, we would allow additional terminals)

- Convert rules
with more non-terminals

$$A \rightarrow X\ Y\ Z$$
$$\Downarrow$$
$$A \rightarrow X\ Q$$
$$Q \rightarrow Y\ Z$$

  ($Q$ is a new non-terminal, specific to this rule)

- But:

  - increases the number of non-terminals ("grammar constant")
  - can be tricky for SCFG rules

- Remove all rules with scope $> 3$
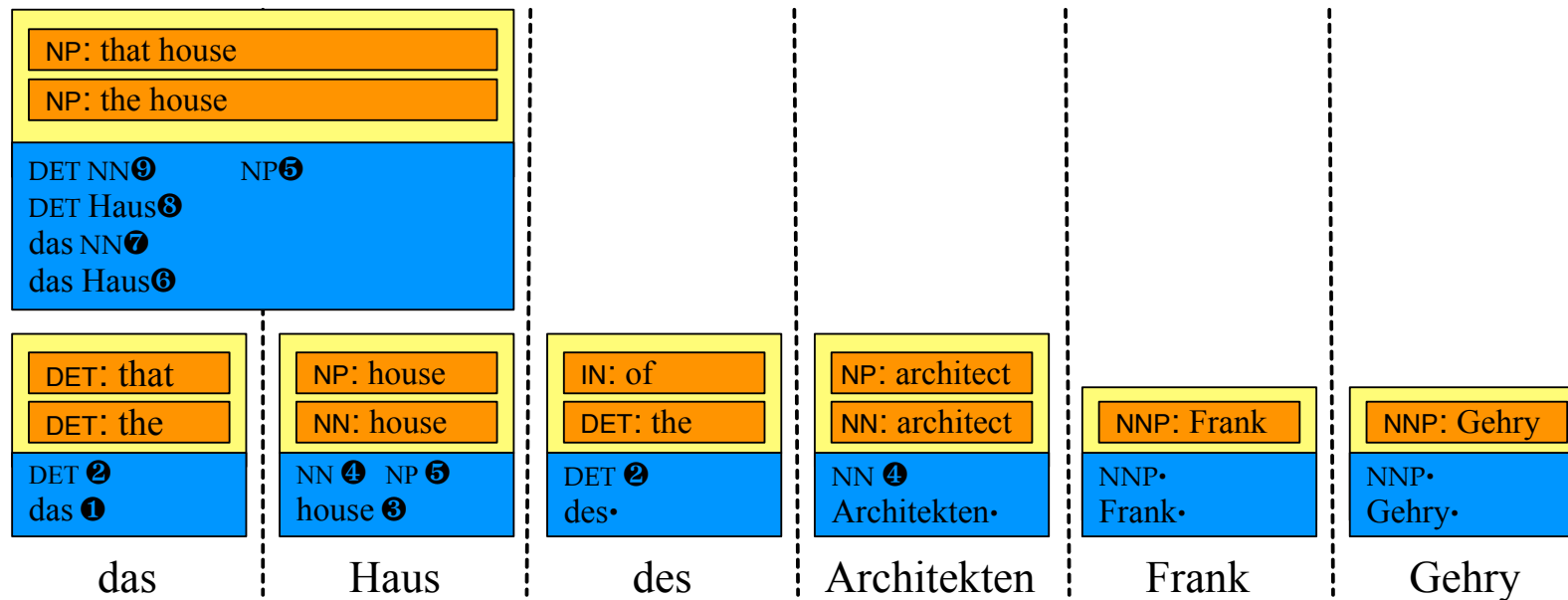
- Less restrictive than CNF
  e.g., allows:

$$A \rightarrow \text{DET}_1 \text{ maison } \text{JJ}_2 \text{ sur la } \text{NN}_3 \quad | \quad \text{DET}_1 \text{ JJ}_2 \text{ house on the } \text{NN}_3$$

  (2 choice points at edges)

- Better speed/quality trade-off than synchronous binarization
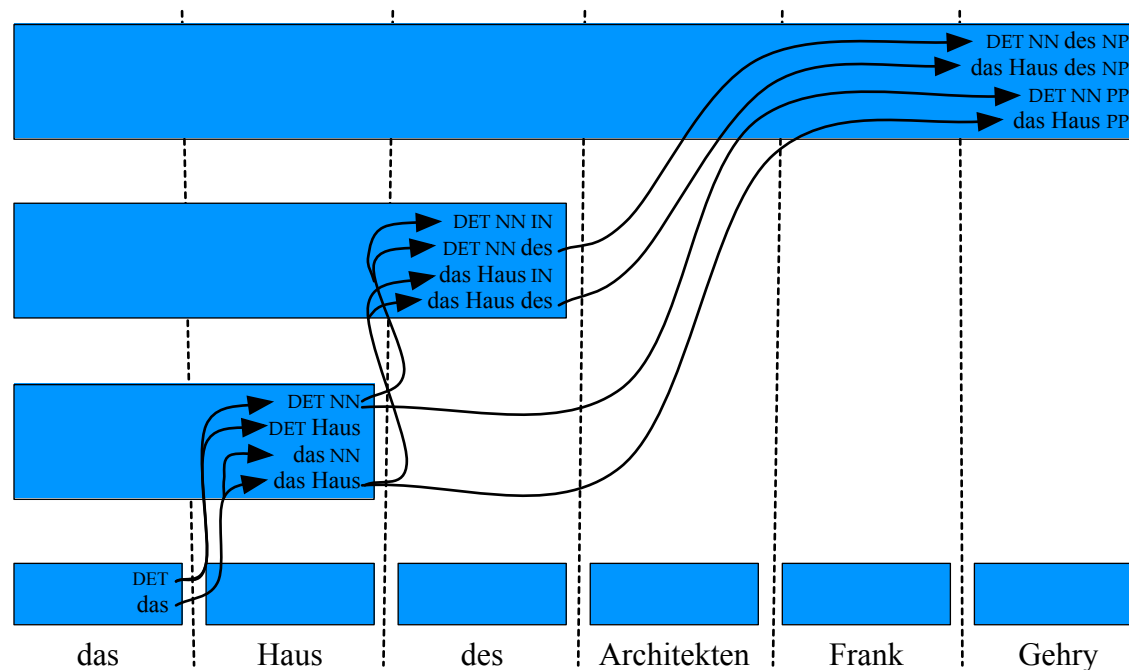
# recursive cky+

- Two charts: (1) hypothesis chart, (2) dotted rule chart



- Dotted rule chart allows dynamic programming of rules with same prefix

# Expansion of Dotted Rules
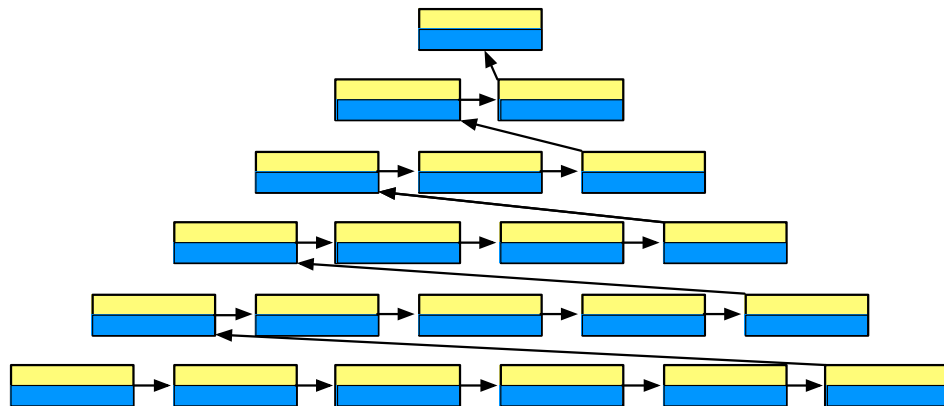
- Dotted rules are expanded recursively



- Dotted rules are stored with each chart cell

# Recursive CKY+

- Recursive CKY+ (Sennrich, 2014) removes need for dotted rule chart

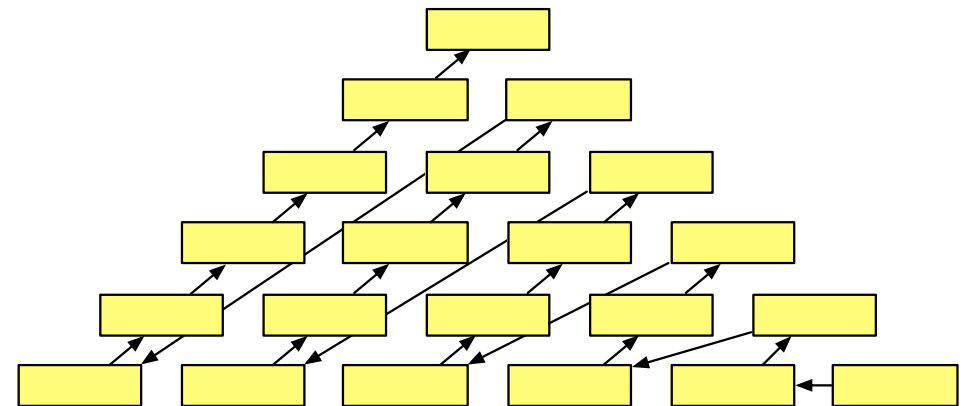- Chart traversal is re-arranged

CKY+

bottom-up, left-to-right

recursive CKY+

right-to-left, depth-first
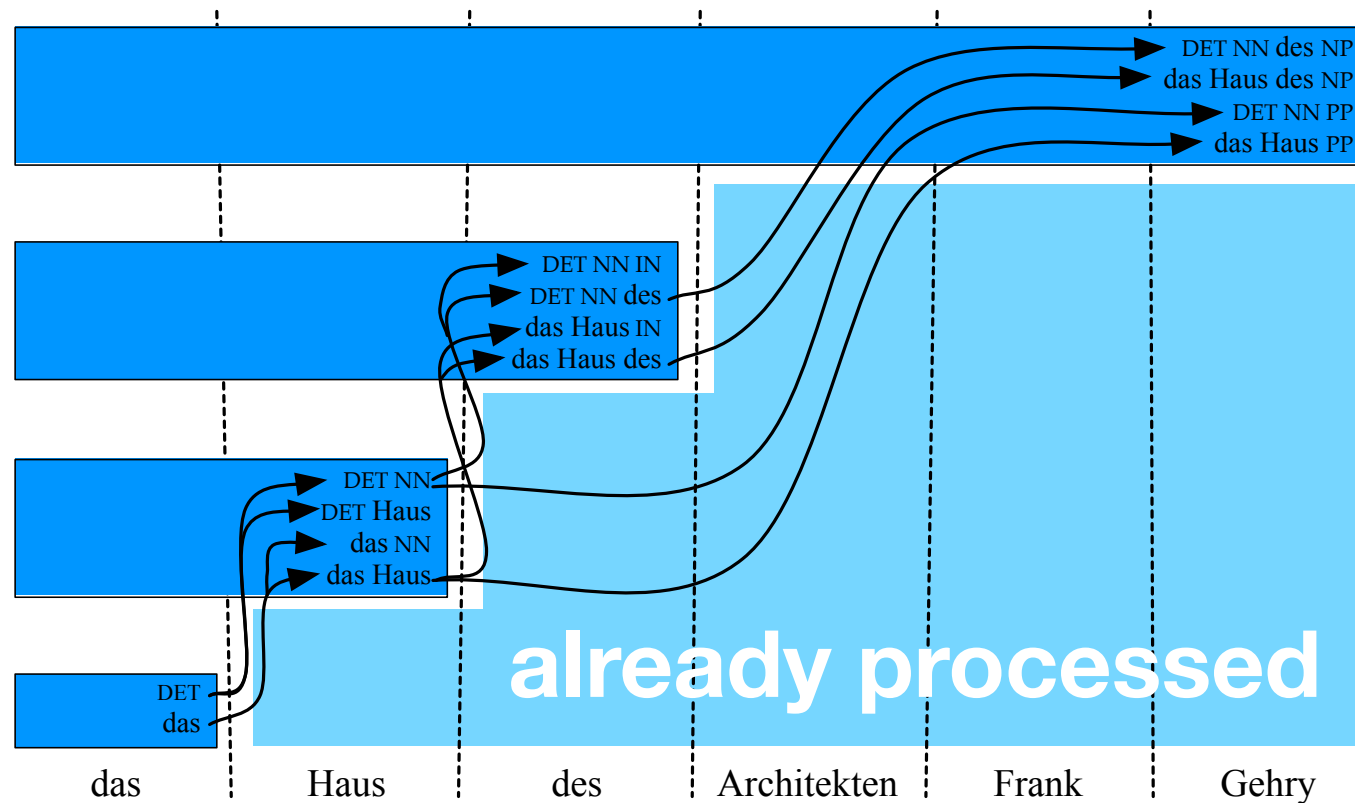


with dotted rule chart

without dotted rule chart

# Recursive CKY+

- Rule expansion by recursive function calls

- Rules can be immediately expanded, because all needed cells already processed
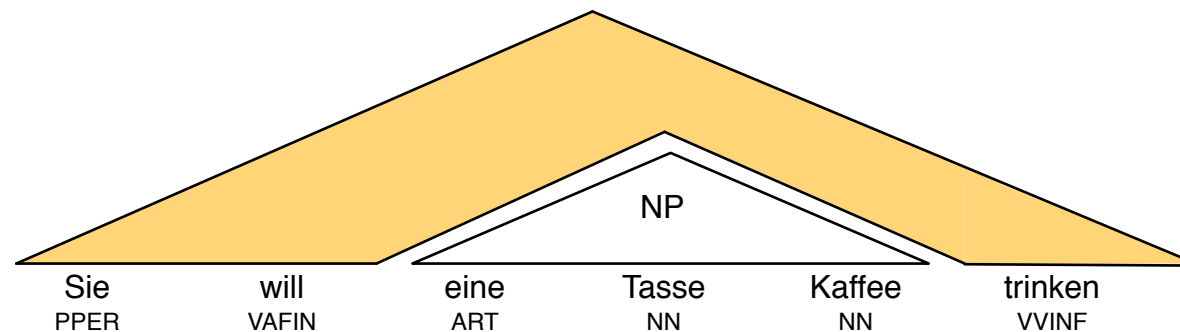
# search strategies

# Two-Stage Decoding

- First stage: decoding without a language model (-LM decoding)

  – may be done exhaustively
  – eliminate dead ends
  – 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

  – stage 1: find applicable rules
  – stage 2: cube pruning

- Decode with increasingly complex model

- Examples

    – reduced language model [Zhang and Gildea, 2008]
    – reduced set of non-terminals [DeNero et al., 2009]
    – language model on clustered word classes [Petrov et al., 2008]

# Outside Cost Estimation

- Which spans should be more emphasized in search?

- Initial decoding stage can provide outside cost estimates

```
                              NP

      Sie       will    eine     Tasse    Kaffee    trinken
      PPER      VAFIN   ART      NN       NN        VVINF
```

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

- What causes the high search error rate?

- Where does the best translation fall out the beam?

- How accurate are LM estimates?

- Are particular types of rules too quickly discarded?

- Are there systemic problems with cube pruning?