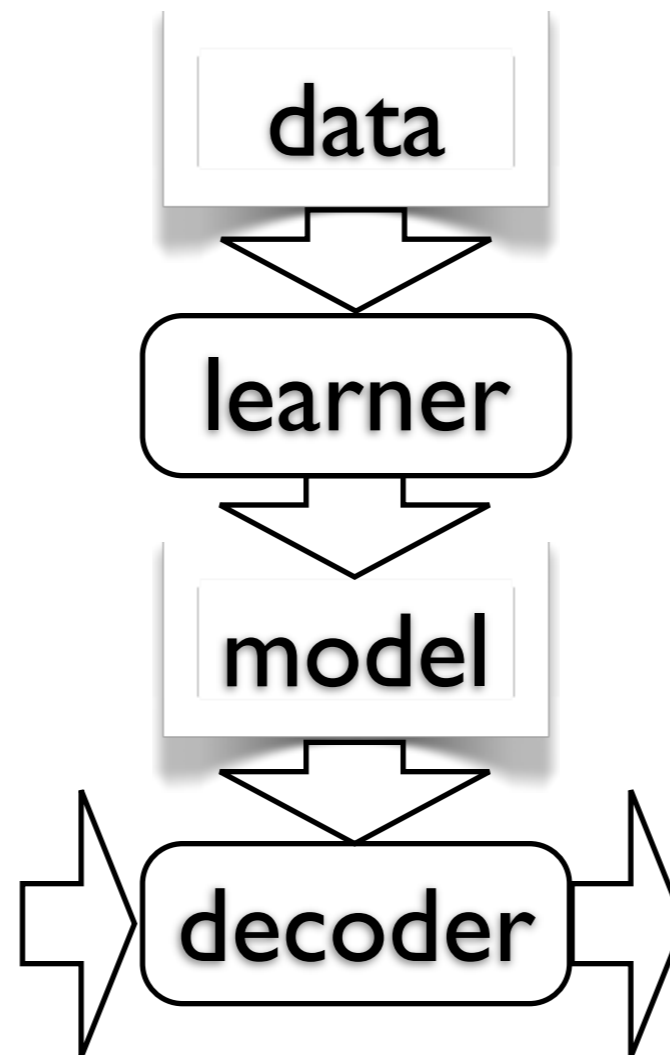# Structures in Statistical Machine Translation

Taro Watanabe @ NICT
taro . watanabe @ nict . go . jp

# Machine Translation

data

↓

learner

↓

model

↓

decoder

黑山头口岸联检部门将原来要二至三天办完的出入境手续改为一天办完。

The United Inspection Department of Heishantou Port has shortened the procedures for leaving and entering the territory from originally 2 - 3 days to 1 day.

- We learn parameters from data assuming a "model"

- Decode by the learned parameters

# Channel Model

$X$ ➡️ | Process | ➡️ $Y$

# Channel Model
## + noise



encoder $\Rightarrow Y \Rightarrow$ channel $\Rightarrow X \Rightarrow$ decoder

$$
\begin{aligned}
\hat{y} &= \underset{y}{\operatorname{argmax}} \, Pr(y|x) \\
&= \underset{y}{\operatorname{argmax}} \, \frac{Pr(x|y)Pr(y)}{Pr(x)} \\
&= \underset{y}{\operatorname{argmax}} \, Pr(x|y)Pr(y)
\end{aligned}
$$

f = source

e = target

$$
\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \, Pr(\mathbf{f}|\mathbf{e})Pr(\mathbf{e})
$$

- Employed in: ASR, OCR, MT...

# Translation Model

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \; \boxed{Pr(\mathbf{f}|\mathbf{e})} \; \boxed{Pr(\mathbf{e})}$$

Translation Model    Language Model

(Brown et al., 1990)

- Translation Model: adequacy of translation

- Language Model: grammatical correctness, consistent style, fluency

# Language Model

$$Pr(\text{I do not know}) \quad = \quad ?$$
$$Pr(\text{I not do know}) \quad = \quad ?$$

- Likelihood of a string of English words

- Usually modeled by ngrams

$$W = w_1, w_2, w_3, \cdots w_N$$

$$
\begin{aligned}
p(W) \quad &= \quad p(w_1, w_2, w_3, \cdots, w_N) \\
&= \quad p(w_1)p(w_2|w_1)p(w_3|w_1, w_2)\cdots \\
&\qquad p(w_N|w_1, w_2, w_3, \cdots, w_{N-1})
\end{aligned}
$$

# ngram Language Model

- Markov assumption: only n-word history is memorized

- Bigram:

$$p(\text{I do not know}) \quad = \quad p(\text{I})p(\text{do}|\text{I})p(\text{not}|\text{do})p(\text{know}|\text{not})$$
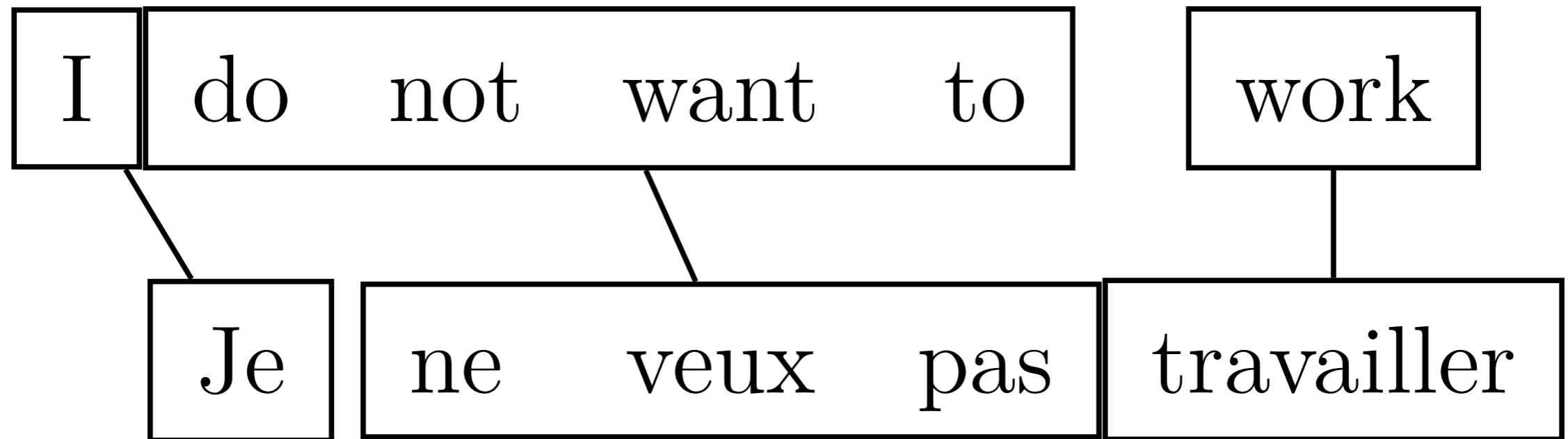
- Training: Maximum likelihood estimate + smoothing (Good-Turing, Witten-Bell, Kneser-Ney etc.)
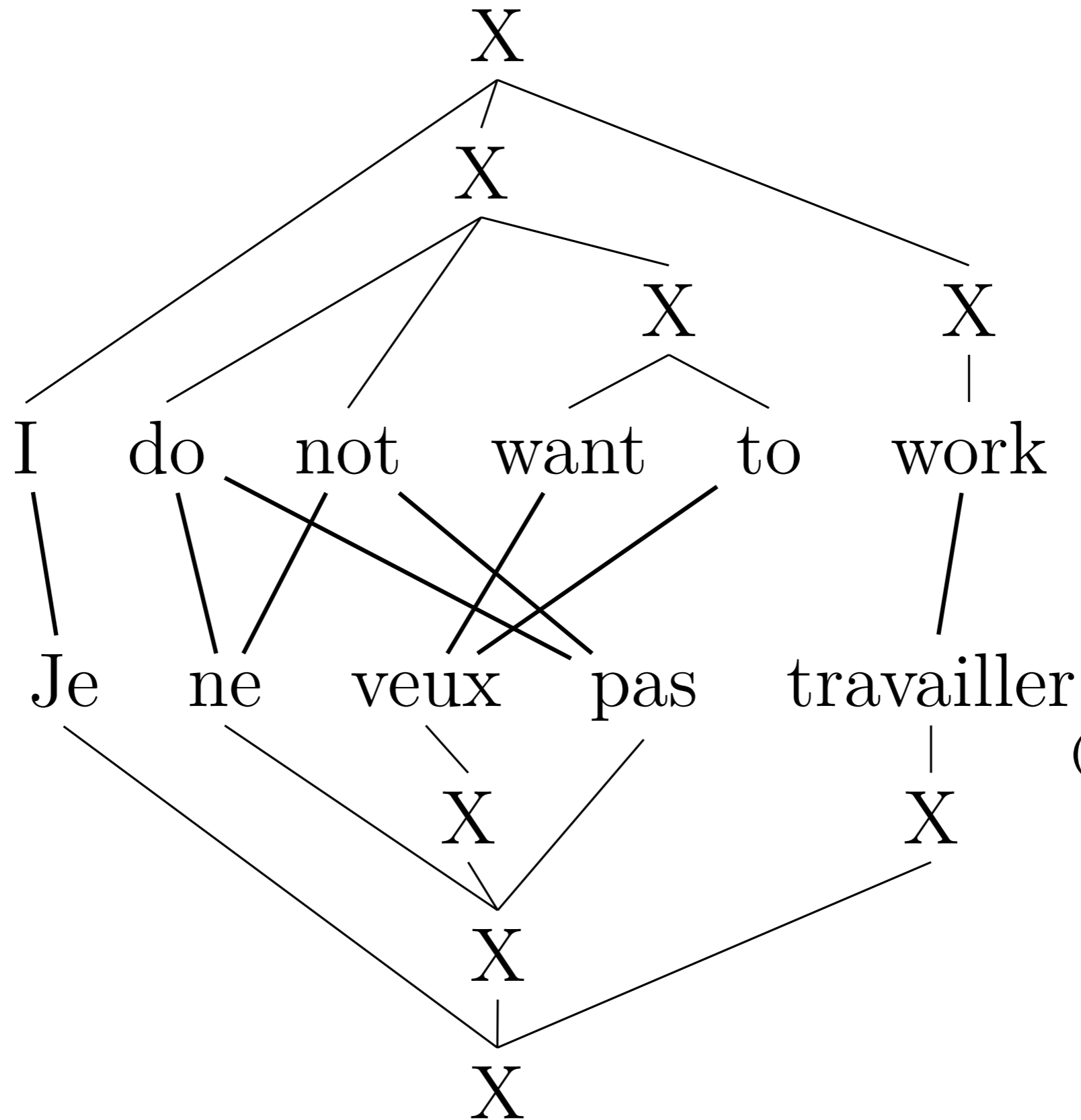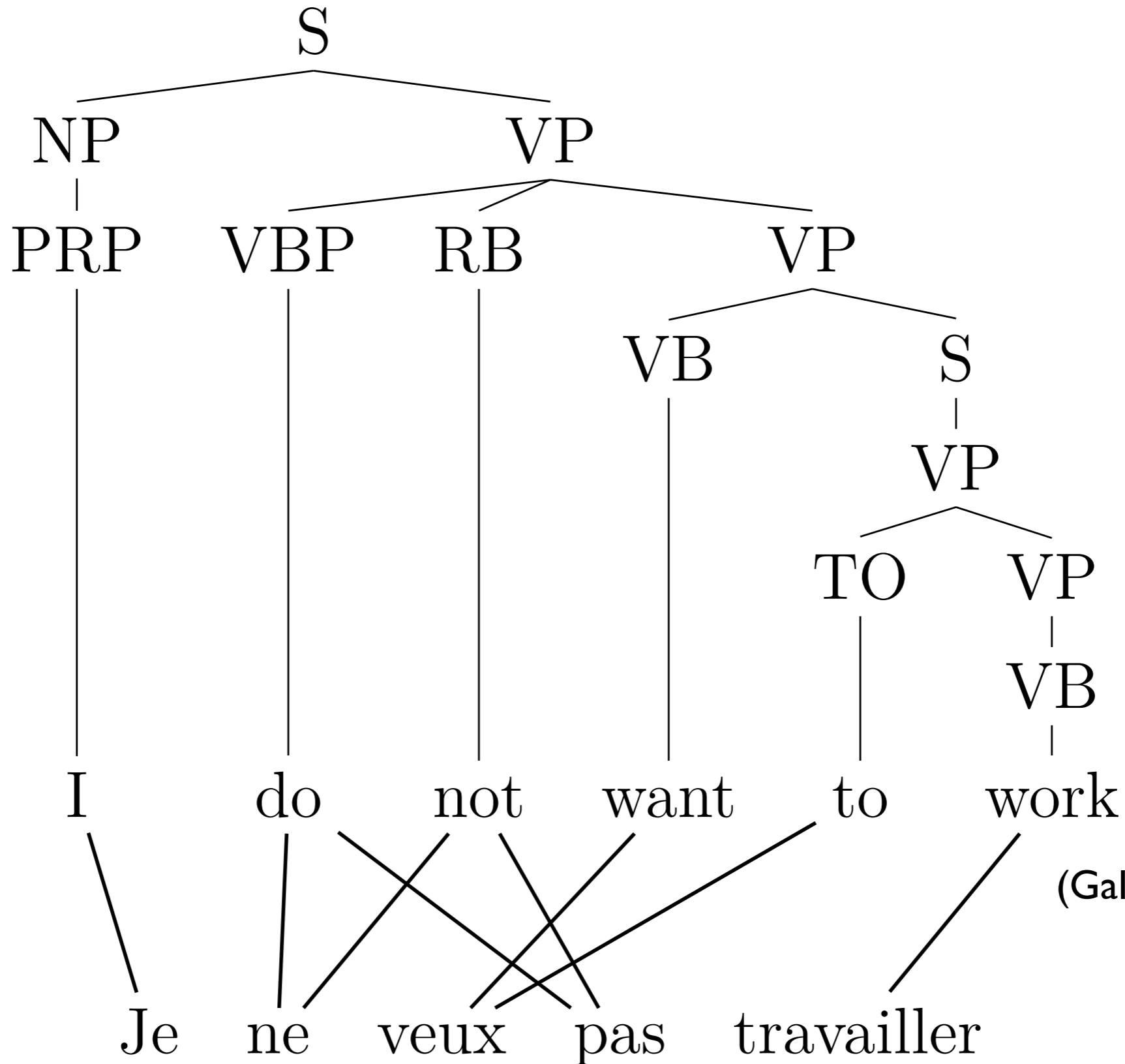
# Word-based MT

I  do  not  want  to  work

Je  ne  veux  pas  travailler

(Brown et al., 1993)

# Phrase-based MT



(Koehn et al., 2003)

# Hierarchical PBMT



(Chiang, 2007)

# Syntax-based MT



(Galley et al., 2004)
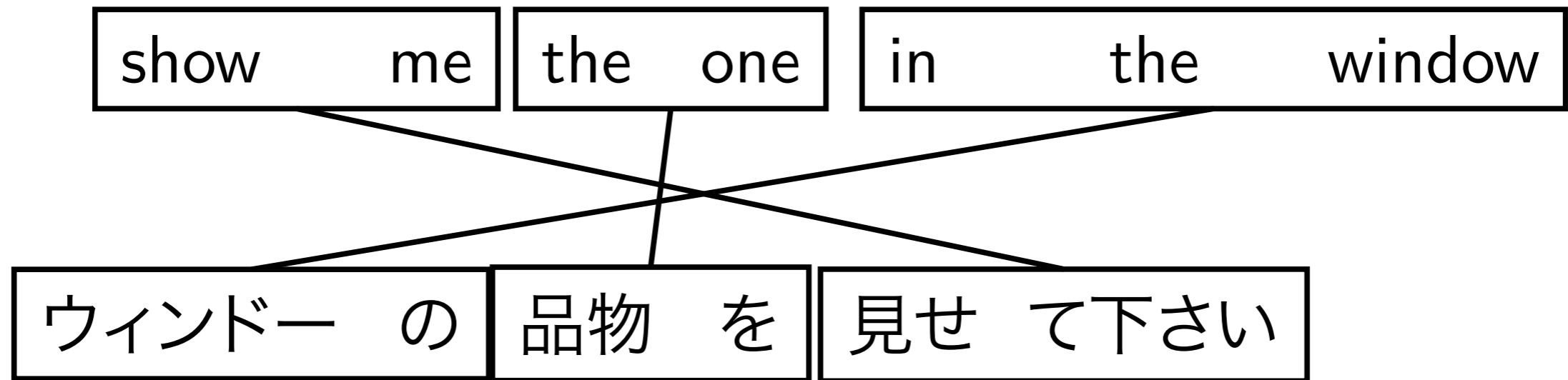
# Structures in SMT

- Tutorial

  - **Phrase-based MT**

  - Tree-based MT

- Syntactic Structures in System Combination

# Why Phrases?

- Use phrases as a unit of translations

    - Directly handle many-to-many word correspondence + local reordering

    - Allow local context + non-compositional phrases

- Employed in many systems, including Google, NICT(VoiceTra, TexTra) and open-source, Moses (http://www.statmt.org/moses/)

# Phrase-based Model

| show        me | the   one | in        the        window |
|:---|:---|:---|

| ウィンドー   の | 品物   を | 見せ て下さい |
|:---|:---|:---|

- Generative story:
  - f is segmented into phrases
  - Each phrase is translated
  - Translated phrases are reordered

# Phrase-based Model

$$\hat{\mathbf{e}} = \operatorname*{argmax}_{\mathbf{e}} \frac{\exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})\right)}{\sum_{\mathbf{e}', \phi'} \exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}', \phi', \mathbf{f})\right)}$$

$$= \operatorname*{argmax}_{\mathbf{e}} \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- Maximization of a log-linear combination of multiple feature functions h(e, Φ, f)

- Φ: phrasal partition of f and e

- w: weight of feature functions

# Questions

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \, \mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- **Training: How to learn phrases and parameters (Φ and h)?**

- Decoding (or search): How to find the best translation (argmax)?

- Tuning (or optimization): How to learn the scaling of features (w)?

# Training

- Learn phrase pairs from $\mathcal{D} = \langle \mathcal{F}, \mathcal{E} \rangle$

- A standard heuristic approach (Koehn et al., 2003)

  - Compute word alignment

  - Extract phrase pairs

  - Score phrases

# Word alignment



(Example from Huang and Chiang, 2007)

# Extract Phrase Pairs



- From word alignment, extract a phrase pair consistent with word alignment

# Exhaustive Extraction



- Exhaustively extract phrases from f, e

# Features from Phrases

$$\log p_\phi(\bar{\mathbf{f}}|\bar{\mathbf{e}}) \quad = \quad \log \frac{\text{count}(\bar{\mathbf{e}}, \bar{\mathbf{f}})}{\sum_{\bar{\mathbf{f}}'} \text{count}(\bar{\mathbf{e}}, \bar{\mathbf{f}}')}$$

$$\log p_\phi(\bar{\mathbf{e}}|\bar{\mathbf{f}}) \quad = \quad \log \frac{\text{count}(\bar{\mathbf{e}}, \bar{\mathbf{f}})}{\sum_{\bar{\mathbf{e}}'} \text{count}(\bar{\mathbf{e}}', \bar{\mathbf{f}})}$$

- Collect all the phrase pairs from the data

- Maximum likelihood estimates by relative frequencies
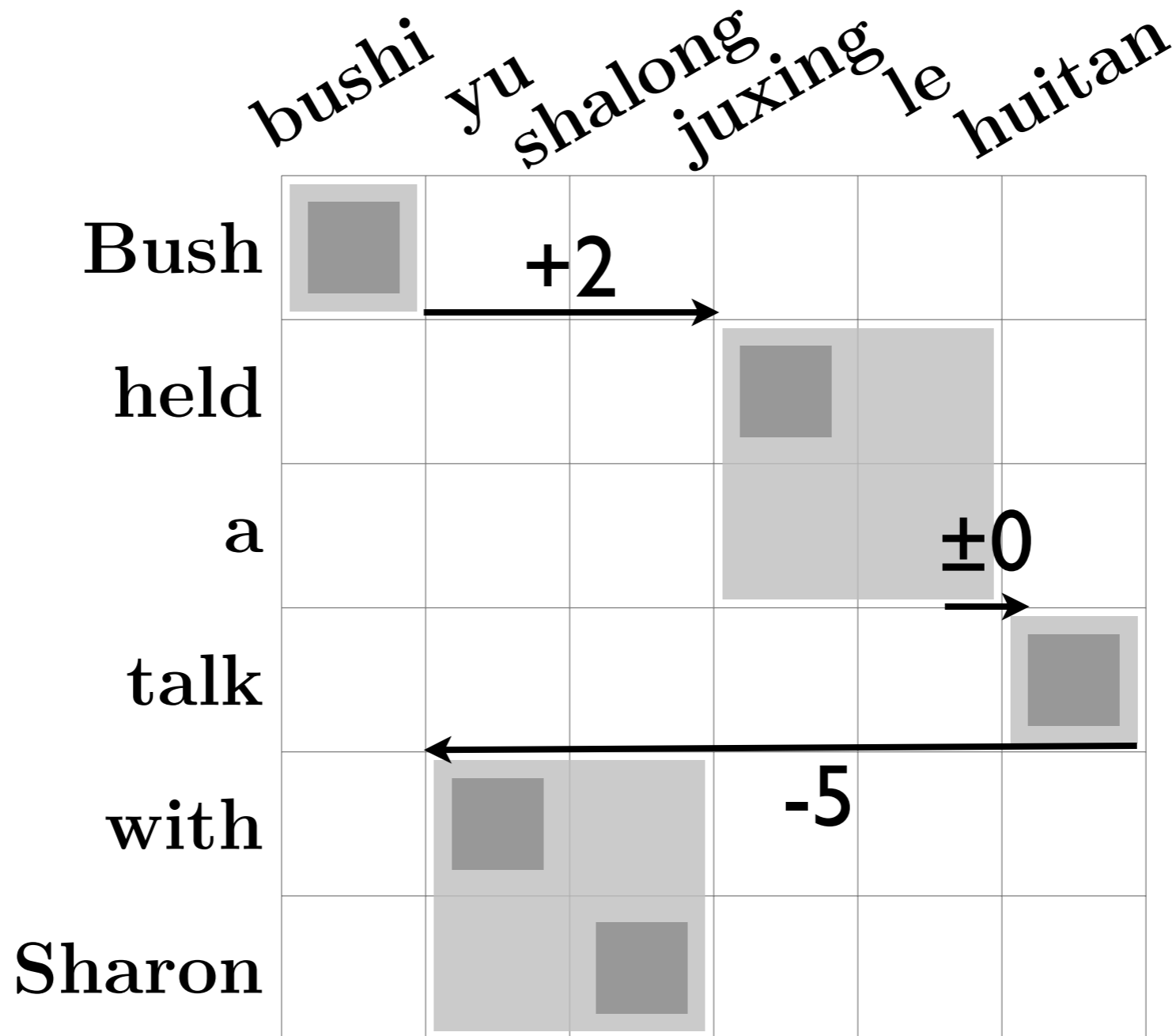
- Employ scores in two directions

# Features from Alignment

$$\log p_{lex}(\bar{\mathbf{f}}|\bar{\mathbf{e}}, \bar{\mathbf{a}}) = \log \prod_i^{|\bar{\mathbf{e}}|} \frac{1}{|\{j|(i,j) \in \bar{\mathbf{a}}\}|} \sum_{\forall(i,j) \in \bar{\mathbf{a}}} t(e_i|f_j)$$

$$\log p_{lex}(\bar{\mathbf{e}}|\bar{\mathbf{f}}, \bar{\mathbf{a}}) = \log \prod_j^{|\bar{\mathbf{f}}|} \frac{1}{|\{i|(j,i) \in \bar{\mathbf{a}}\}|} \sum_{\forall(j,i) \in \bar{\mathbf{a}}} t(f_j|e_i)$$

- Lexical weighing which scores by word translation probabilities

- Idea: counts for rare phrase pairs are unreliable

  - Smoothing effect by decomposing into word pairs

# Features for Distortion



- Distance-based distortion modeling

$$d(\mathbf{f}, \phi, \mathbf{e}) = |+2| + |0| + |-5| = 7$$

# Features for Reordering



- Fine grained reordering features: $\log p_o(o \in \{m, s, d\} \,|\, \bar{\mathbf{f}}, \bar{\mathbf{e}})$

- Either monotone, swap, discontinuous

# Other Features

- log of ngram language model(s)

- word count: bias for ngram language model(s)

- phrase count: shorter or longer phrases

# Direct Training

- Instead of word alignment + extraction pipeline, directly learn phrase-pairs (Marcu and Wong, 2002)

- Bayesian approach + blocked Gibbs sampling to learn parameters (Blunsom et al., 2009)

- Exhaustively memorize longer phrases (Neubig et al., 2011)

# Questions

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \, \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- Training: How to learn phrases and parameters (Φ and h)?

- **Decoding (or search): How to find the best translation (argmax)?**

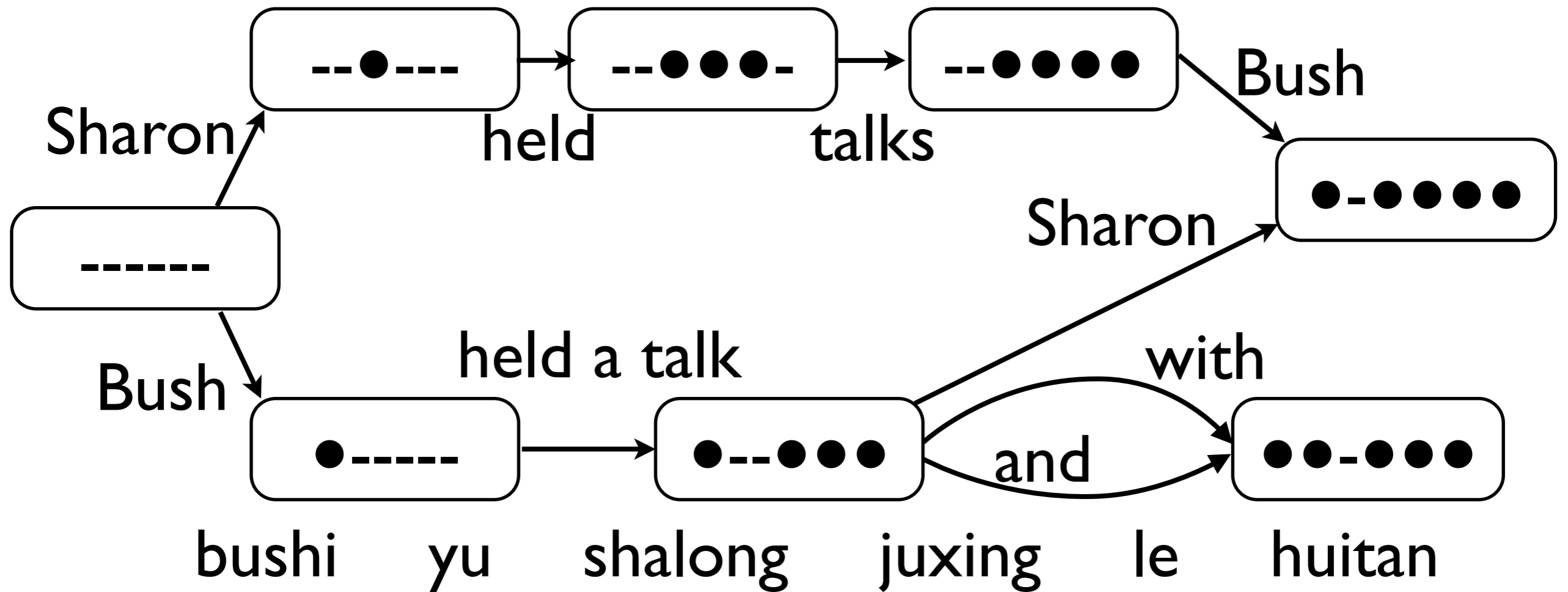- Tuning (or optimization): How to learn the scaling of features (w)?

# Decoding

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \frac{\exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})\right)}{\sum_{\mathbf{e}', \phi'} \exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}', \phi', \mathbf{f})\right)}$$

$$= \underset{\mathbf{e}}{\operatorname{argmax}} \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- Given an input sentence f and phrasal model h and w, seek e with the highest score

- Potential errors:

  - Search error: we cannot find the best scored hypothesis

  - Translation error: highest scored hypothesis is bad

# Enumerate Phrase Pairs

| Bush and | | | held a talk | |
| | with Sharon | | talked | |
| | with | | held | meeting |
| Bush | and | Sharon | hold | talks |

bushi    yu    shalong    juxing    le    huitan

- Given a input sentence f, we can enumerate all possible phrases that match with the source side

- Choose the best phrase pair + ordering

# Phrase-based Search Space



```
          Sharon        --●---        held        --●●●-        talks        --●●●●        Bush
                                                                                          ●-●●●● ●
      ------                                                                    Sharon
          Bush                                                    with
                  ●------    held a talk    ●--●●●              ●●-●●● ●
                                              and
      bushi      yu      shalong      juxing      le      huitan
```
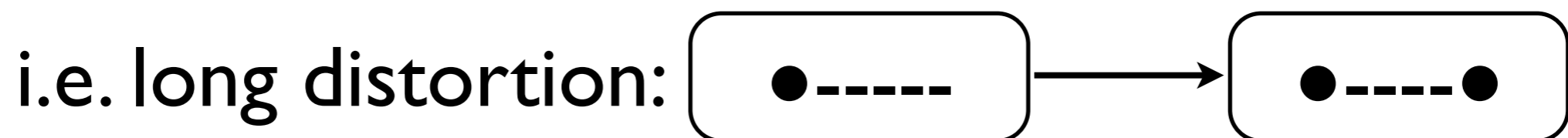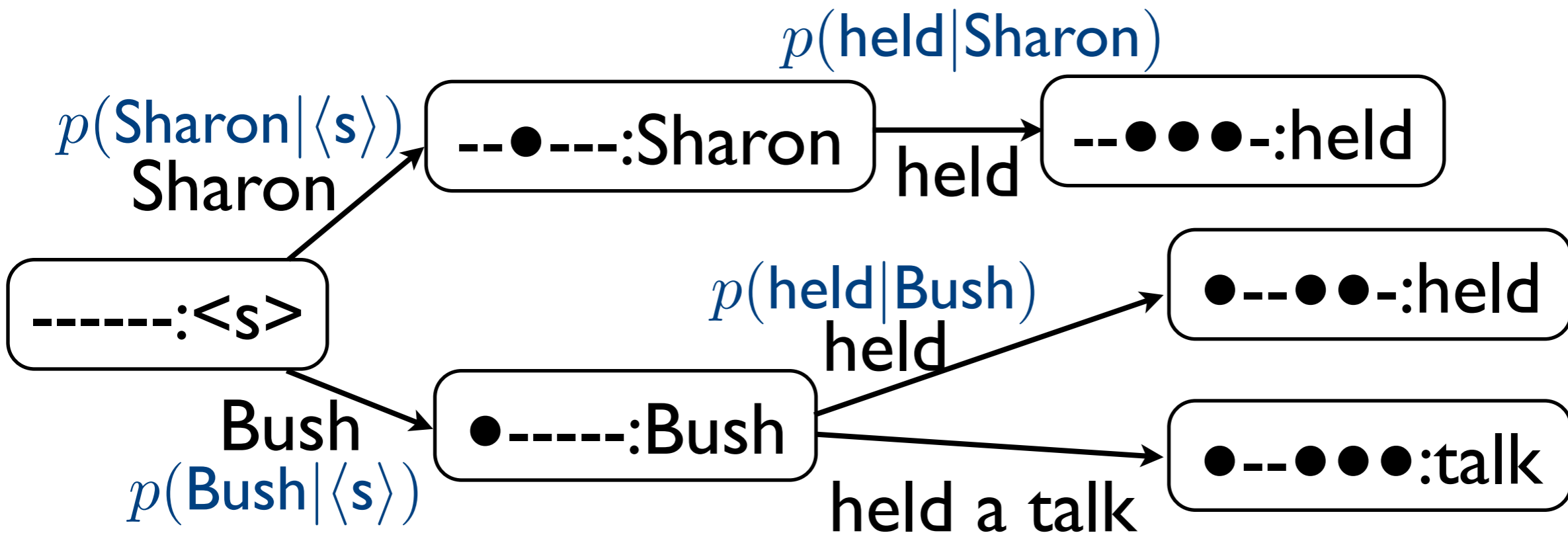
- Node: bit-vector representing covered source words

- Edge: phrasal translations, strictly left-to-right

- Search space: $O(2^n)$, Time: $O(2^n n^2)$ (Why?)

# Traveling Salesman Problem

- NP-hard problem: visit each city only once

- MT as a Traveling Salesman Problem (Knight, 1999)

  - Each source word corresponds to a city

  - A Dynamic Programming solution:

    - State: visited cities (bit-vector)

    - Search space: $O(2^n)$

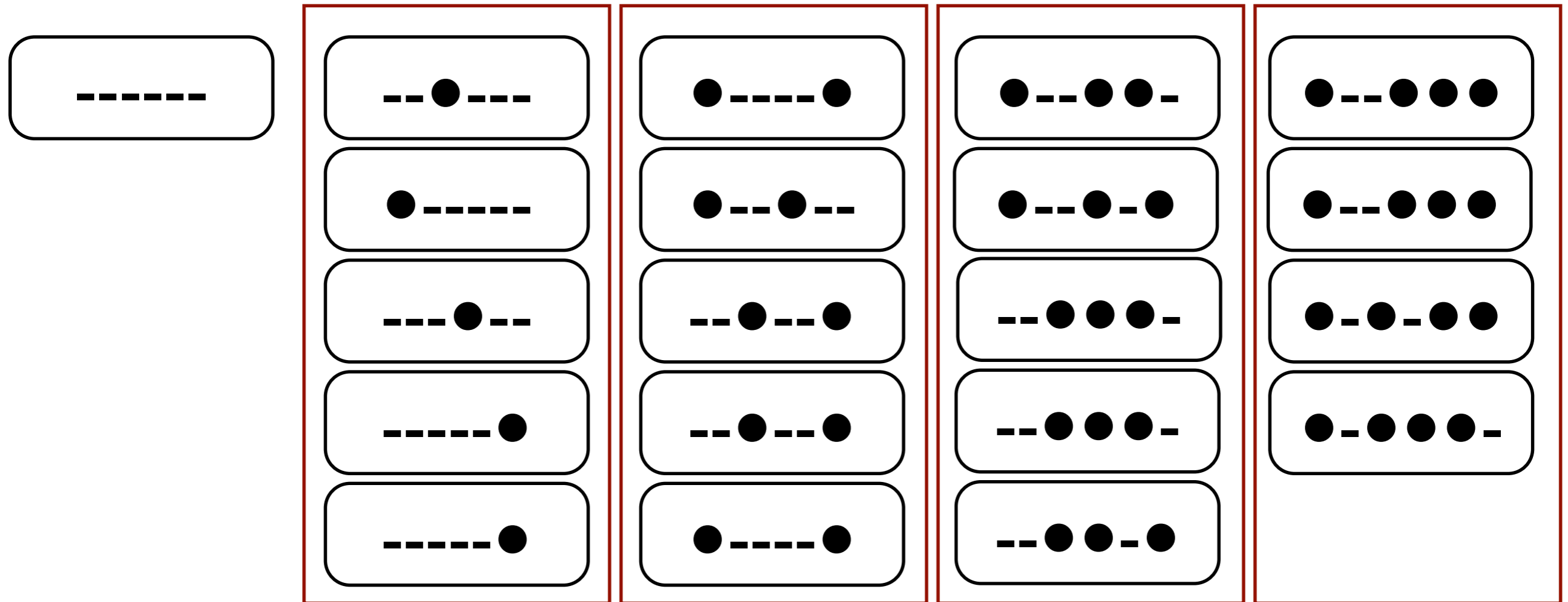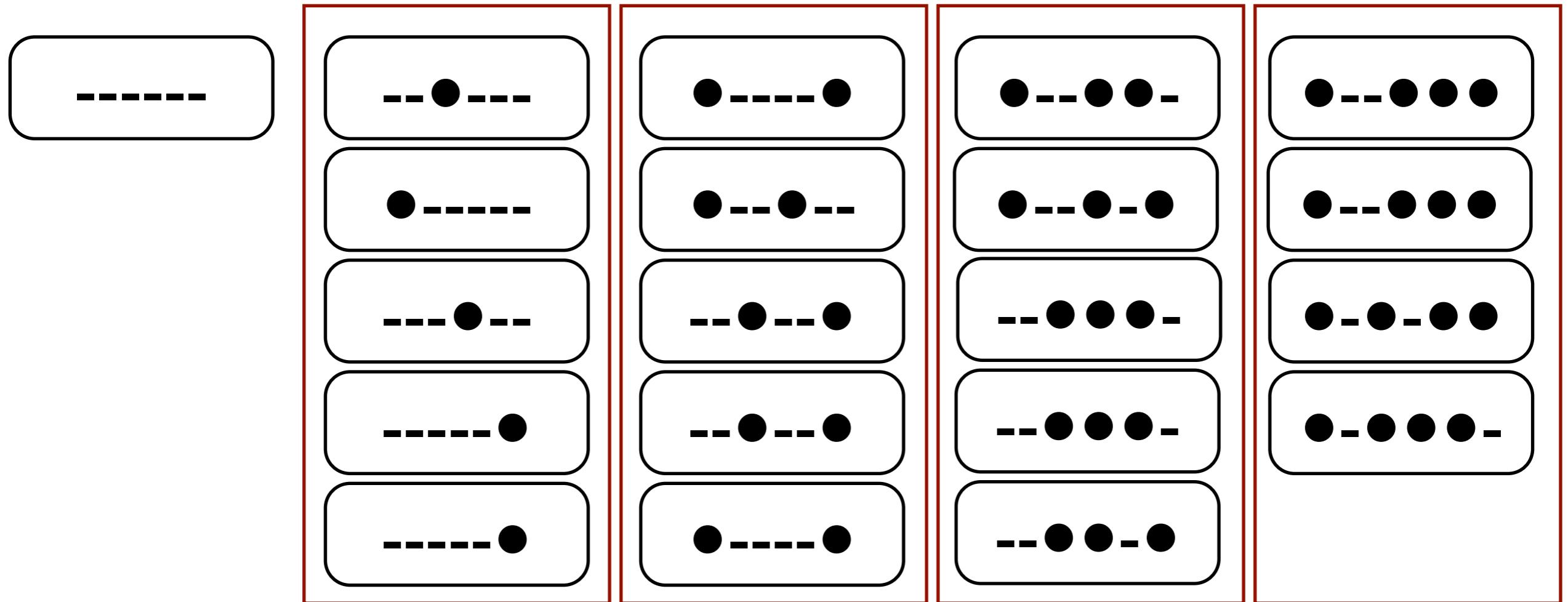  - Distortion limit to reduce search space

    i.e. long distortion:

# Non-local features

$p(\text{held}|\text{Sharon})$

$p(\text{Sharon}|\langle s\rangle)$

Sharon

--●----:Sharon   held→   --●●●-:held

------:<s>

$p(\text{held}|\text{Bush})$
held

●--●●-:held

Bush

●-----:Bush

$p(\text{Bush}|\langle s\rangle)$

held a talk

●--●●●:talk

$p(\text{held}|\text{Bush})p(\text{a}|\text{held})p(\text{talk}|\text{a})$

- Features that requires scoring out of phrases: bigram language model

- Additional state representation required for "future scoring": 1-word for bigram LM

- Space: $O(2^n V^{m-1})$, Time: $O(2^n V^{m-1} n^2)$ for m-gram LM

# Phrase-based Decoding



- Re-organize the search space by the cardinality (= # of covered source words)

- Expand hypotheses from the smallest cardinality first

# Pruning



- Prune hypotheses in a bin sharing the same cardinality

- Expand survived hypotheses only

# Questions

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}}\, \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- Training: How to learn phrases and parameters (Φ and h)?

- Decoding (or search): How to find the best translation (argmax)?

- **Tuning (or optimization): How to learn the scaling of features (w)?**

# Tuning

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \frac{\exp\left(\mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})\right)}{\sum_{\mathbf{e}', \phi'} \exp\left(\mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}', \phi', \mathbf{f})\right)}$$

$$= \underset{\mathbf{e}}{\operatorname{argmax}} \mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, \phi, \mathbf{f})$$

- Three popular objectives (in SMT) for tuning w

  - (Direct) Error Minimization (Och, 2003)

  - Maximum Entropy (Och and Ney, 2002)

  - Large Margin (Watanabe et al., 2007; Chiang et al., 2008; Hopkins and May, 2011)

# (Direct) Minimum Error

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \sum_{s=1}^{S} l(\operatorname*{argmax}_{\mathbf{e}} \mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, \mathbf{f}_s), \mathbf{e}_s)$$

- MERT (Minimum ERror Training)

- Standard in SMT (but not in other NLP areas, such as tagging etc.)

  - We can incorporate arbitrary error functions, l

  - "Summation" can be replaced by document-wise BLEU specific summation
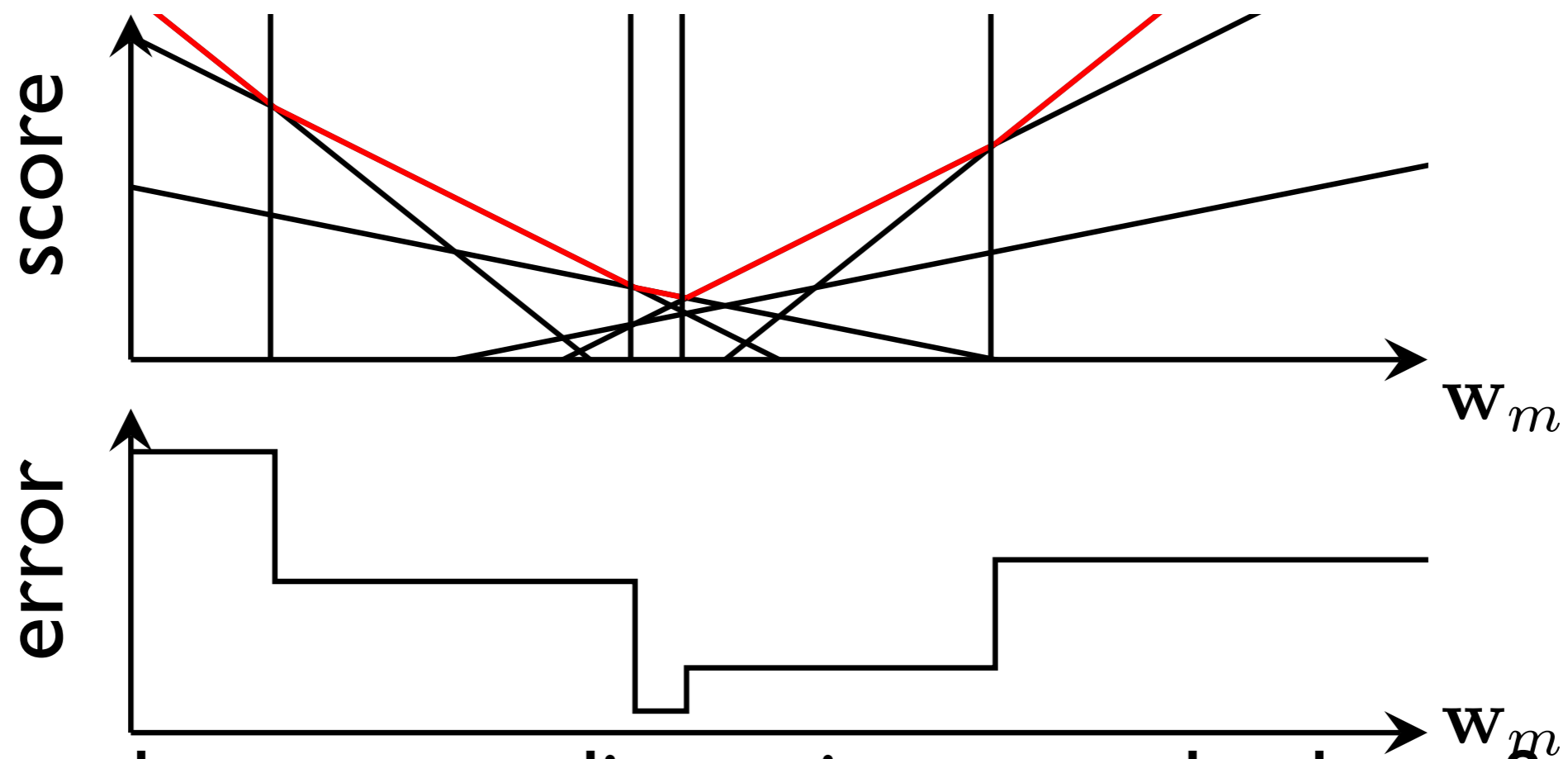
  - 10+ real valued features

# n-best Approximation

1: **procedure** MERT($\{(\mathbf{e}_s, \mathbf{f}_s)\}_{s=1}^{S}$)
2:     **for** $n = 1...N$ **do**
3:         Decode and generate nbest list using $\mathbf{w}$
4:         Merge nbest list
5:         **for** $k = 1...K$ **do**
6:             **for** each parameter $m = 1...M$ **do**
7:                 Solve one dimensional optimization
8:             **end for**
9:             update $\mathbf{w}$
10:        **end for**
11:    **end for**
12: **end procedure**

- N iterations, with each iteration, n-bests are generated and merged

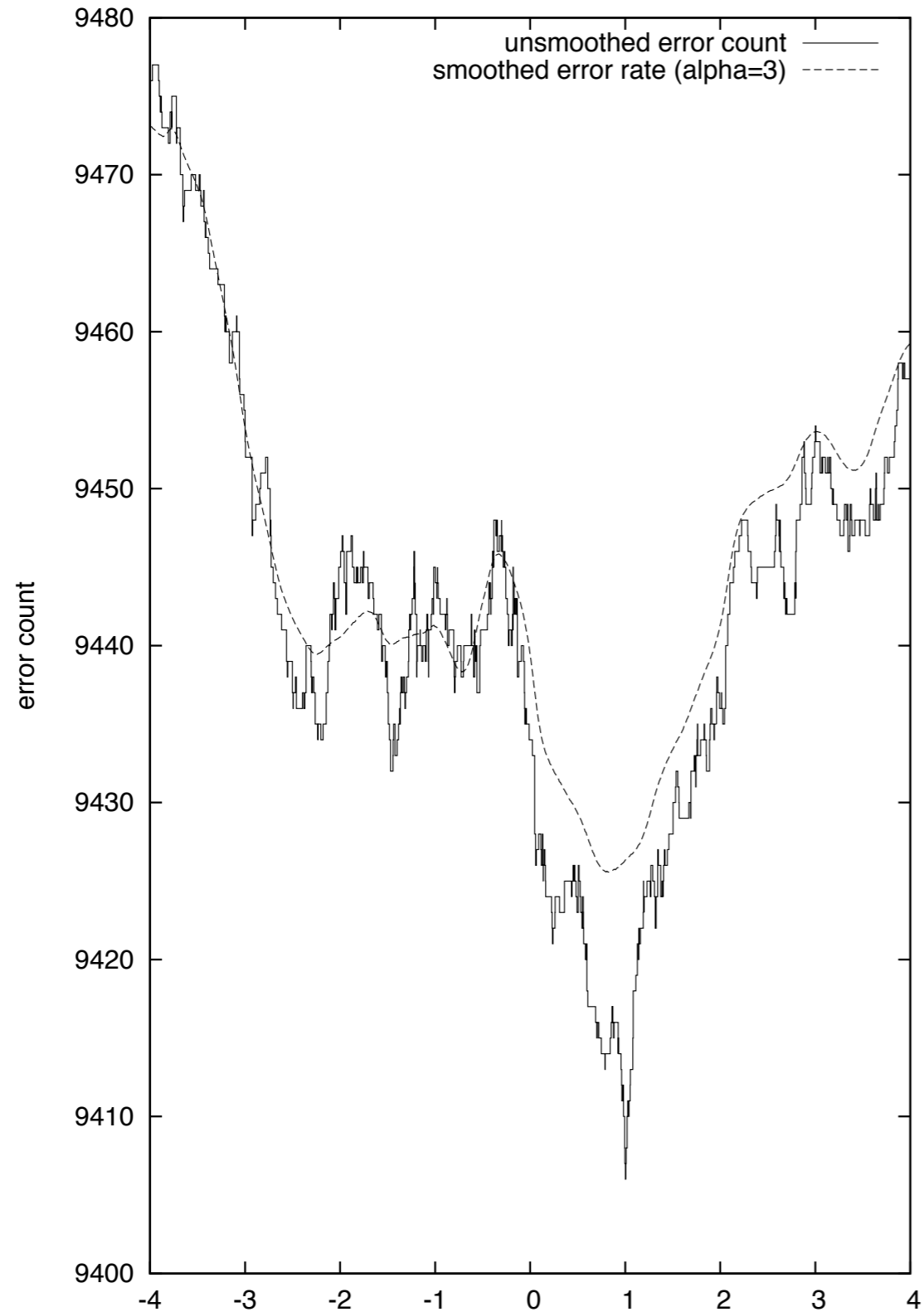  - K iterations, with each iteration, M dimensions are tried (M = # of features), and w is updated

38

# Efficient Line Search

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\operatorname{argmax}} \; \mathbf{w}_m^\top \cdot \underbrace{\mathbf{h}_m(\mathbf{e}, \mathbf{f}_s)}_{\text{slope}} + \underbrace{\mathbf{w}_{m_-}^\top \cdot \mathbf{h}_{m_-}(\mathbf{e}, \mathbf{f}_s)}_{\text{constant}}$$



- If we choose one dimension m, and others fixed, we can treat each hypothesis e as a "line"

- Compute convex hull of a set of "lines"

39

# Error Surface



(Och, 2003)

# MERT in Practice

- Many random starting points (Macherey et al., 2008; Moore and Quirk, 2008)

- Many random directions (Macherey et al., 2008)

- Error count smoothing (Cer et al., 2008)

- Regularization (Hayashi et al., 2009)

- Multi-dimensional search by efficiently computing convex hull (Galley and Quirk, 2011)

- MERT at least 3 times, and report average BLEU (Clark et al., 2011)

# Maximum Entropy

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \frac{\lambda}{2}\|\mathbf{w}\|^2 - \sum_{s=1}^{S} \log \frac{\displaystyle\sum_{\mathbf{e}*\in\mathsf{ORACLE}(\mathbf{f}_s)} \exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}^*, \mathbf{f}_s)\right)}{\displaystyle\sum_{\mathbf{e}'\in\mathsf{GEN}(\mathbf{f}_s)} \exp\left(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}', \mathbf{f}_s)\right)}$$

- Minimize the negative log-likelihood of generating good translations (Och and Ney, 2002)

- ORACLE is a subset of GEN, a set of hypotheses with minimum loss

- Optimized by L-BFGS or SGD

- Potentially large # of features as in NLP tasks

# Why Not MaxEnt?

| error criterion used in training | mWER [%] | mPER [%] | BLEU [%] | NIST | # words |
|:---:|:---:|:---:|:---:|:---:|:---:|
| confidence intervals | +/- 2.7 | +/- 1.9 | +/- 0.8 | +/- 0.12 | - |
| MMI | **68.0** | *51.0* | 11.3 | 5.76 | 21933 |
| mWER | *68.3* | *50.2* | 13.5 | 6.28 | 22914 |
| smoothed-mWER | *68.2* | *50.2* | 13.2 | 6.27 | 22902 |
| mPER | *70.2* | *49.8* | 15.2 | *6.71* | 24399 |
| smoothed-mPER | *70.0* | **49.7** | 15.2 | *6.69* | 24198 |
| BLEU | 76.1 | 53.2 | **17.2** | 6.66 | 28002 |
| NIST | 73.3 | *51.5* | *16.4* | **6.80** | 26602 |

- In Och and Ney (2002), they used

  - WER to select oracle translations

  - n-best merging approach to approximate summation as in MERT

43

# Large Margin

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \frac{\lambda}{2}||\mathbf{w}||^2 + \sum_{s=1}^{S}\sum_{\mathbf{e}_s^*}\sum_{\mathbf{e}_s'} \xi_{s,\mathbf{e}_s^*,\mathbf{e}_s'}$$

$$\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}_s^*, \mathbf{f}_s) - \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}_s', \mathbf{f}_s) \geq l(\mathbf{e}_s', \mathbf{e}_s^*) - \xi_{s,\mathbf{e}_s^*,\mathbf{e}_s'}$$

$$\mathbf{e}_s^* \in \mathsf{ORACLE}(\mathbf{f}_s)$$

$$\mathbf{e}_s' \in \mathsf{GEN}(\mathbf{f}_s)$$

- Structured output learning approach

- Very hard to enumerate all possible e' and oracle translations e*

- Solution: online learning or n-best approximation

# Online Learning

**Require:** $\{(\mathbf{f}_s, \mathbf{e}_s)\}_{s=1}^{S}$

1: $\mathbf{w}^1 = \{0\}$

2: $t = 1$

3: **for** $1...N$ **do**

4:      $s \sim \text{random}(1, S)$

5:      $\hat{\mathbf{e}} \in \mathsf{GEN}(\mathbf{f}_s, \mathbf{w}^{t-1})$

6:      **if** $l(\hat{\mathbf{e}}, \mathbf{e}_s) \geq 0$ **then**

7:          $\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{h}(\mathbf{e}_s, \mathbf{f}_s) - \mathbf{h}(\hat{\mathbf{e}}, \mathbf{f}_s)$

8:          $t = t + 1$

9:      **end if**

10: **end for**

11: **return** $\mathbf{w}^t$ or $\frac{1}{N} \sum_{i=1}^{N} \mathbf{w}^j$

- Averaged perceptron (Liang et al., 2006)

- Scale to large data, but each iteration requires decoding + weight update

# Online Large Margin

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}'} \frac{\lambda}{2}||\mathbf{w}' - \mathbf{w}||^2 + \max\left(l_s - \mathbf{w}'^{\top} \cdot \Delta\mathbf{h}_s\right)$$

$$\hat{\mathbf{e}}_s = \operatorname*{argmax}_{e} \mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, \mathbf{f}_s)$$

$$l_s = l(\hat{\mathbf{e}}_s) - l(\mathbf{e}_s^*)$$

$$\Delta\mathbf{h}_s = \mathbf{h}(\hat{\mathbf{e}}_s, \mathbf{f}_s) - \mathbf{h}(\mathbf{e}^*, \mathbf{f}_s)$$

- line 7 is replaced by the solution of the above equation

- Still, requires decoding + update in each iteration

- Hard to determine when to stop (watch another dev data)

# Ranking Approach

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \frac{\lambda}{2} ||\mathbf{w}||^2 + \sum_{s=1}^{S} \sum_{\mathbf{e}''_s} \sum_{\mathbf{e}'_s} \xi_{s,\mathbf{e}''_s,\mathbf{e}'_s}$$

$$-\log\left(1 + \exp(-\mathbf{w}^\top \cdot \Delta\mathbf{h}_{\mathbf{e}''_s,\mathbf{e}'_s})\right) \geq -\xi_{s,\mathbf{e}''_s,\mathbf{e}'_s}$$

$$\mathbf{e}''_s, \mathbf{e}'_s \in \mathsf{GEN}(\mathbf{f}_s)$$

$$l(\mathbf{e}'_s, \mathbf{e}''_s) > 0$$

$$\Delta\mathbf{h}_{\mathbf{e}''_s,\mathbf{e}'_s} = \mathbf{h}(\mathbf{e}''_s, \mathbf{f}_s) - \mathbf{h}(\mathbf{e}'_s, \mathbf{f}_s)$$

- An n-best approximation approach (Hopkins and May, 2011)

- Pair-wise comparison of all the hypotheses

- logistic-loss (or 0-1 loss): use an off-the-shelf binary classifier

# Results



- Reranking is competitive to MERT and MIRA, and scales to large # of features

# Answered?

- Grammar-less model (but very strong)

- Fast decoding

- Why MERT? (Good for non-binary, numerical features)

# Structures in SMT

- Tutorial
  - Phrase-based MT
- **Tree-based MT**
- Syntactic Structures in System Combination

# Tree-based MT

- Backgrounds

  - CFG, parsing, hypergraph, deductive system semirings

- Tree-based SMT

  - Synchronous-CFG

  - String-to-Tree, Tree-to-String

# Backgrounds: CFG

$$S \rightarrow NP\ VP$$
$$NP \rightarrow NNP$$
$$NP \rightarrow NP\ PP$$
$$NP \rightarrow DP\ NN$$
$$NP \rightarrow DT\ NN$$
$$VP \rightarrow VBD\ NP$$
$$NNP \rightarrow Bush$$
$$VBD \rightarrow held$$
$$\vdots$$



- parsing = intersection of CFG with a string (regular grammar)

# Parsing: CKY



$X_{i,j}$

$X \rightarrow Y \ Z$

$Y_{i,k}$   $Z_{k,j}$

Bush   held   a   talk   with   Sharon

- O(n³) : For each length n, for each position i, for each rule X → Y Z, for each split point k

- (Bottom-up) topological order

# Hypergraph

$$\text{S}_{0,6}$$

$$\text{NP}_{0,1} \qquad \text{VP}_{1,6} \qquad e \quad = \quad \langle \underbrace{\text{VP}_{1,6}}_{h(e)}, \underbrace{\{\text{VBD}_{1,2}, \text{NP}_{2,6}\}}_{T(e)} \rangle$$

$$\uparrow$$

$$\text{NNP}_{0,1} \quad \text{VBD}_{1,2} \quad \text{NP}_{2,6}$$

$$\uparrow \qquad \uparrow$$

Bush     held

$$\text{VP}_{1,6}$$

$$\wedge$$

$$\text{VBD}_{1,2} \quad \text{NP}_{2,6}$$

(Klein and Manning, 2001)

- Generalization of graphs:
  - h(e): head node of hyperedge e
  - T(e): tail node(s) of hyperedge e, arity = |T(e)|
  - hyperedge = instantiated rule
- Represented as and-or graphs

# Deductive System

$$\text{VBD}_{1,2} \quad \text{NP}_{2,6}$$

$$\text{VP}_{1,6}$$

$$\overbrace{\frac{\text{VBD}_{1,2} \ \text{NP}_{2,6}}{\underbrace{\text{VP}_{1,6}}_{\text{consequent}}}}^{\text{antecedents}} \text{VP}_{[i,j]} \to \text{VBZ}_{[j,k]} \ \text{NP}_{[i,k]}$$

(Shieber et al., 1995)

- Parsing algorithm as a deductive system

- We start from initial items (axioms) until we reach a goal item

- If antecedents are proved, its consequent is proved

- deduction = hyperedge

# Packed Forest

$$VP_{1,6}$$

$$VBD_{1,2}$$  $$NP_{2,6}$$

$$NP_{2,4}$$  $$PP_{4,6}$$

$$\frac{VBD_{1,2} \ \dfrac{NP_{2,4} \ PP_{4,6}}{NP_{2,6}}}{VP_{1,6}}$$

$$\frac{VBD_{1,2} \ NP_{2,4} \ PP_{4,6}}{VP_{1,6}}$$

(Klein and Manning, 2001; Huang and Chiang, 2005)

- A polynomial space encoding of exponentially many parses by sharing common sub-derivations

- Single derivation = tree

# Summary of Formalisms

| hypergraph | AND/OR graph | CFG | deductive system |
|---|---|---|---|
| vertex | OR-node | symbol | item |
| source-vertex | leaf OR-node | terminal | axiom |
| target-vertex | root OR-node | start symbol | goal item |
| hyperedge | AND-node | production | instantiated deduction |

$$\langle v, \{u_1, u_2\} \rangle \qquad \qquad v \to u_1\ u_2 \qquad \qquad \frac{u_1\ u_2}{v}$$

57

# Weights and Semirings

$$\text{VP} \quad \overset{w_1}{\rightarrow} \quad \text{VBD NP}$$

$$\text{NP} \quad \overset{w_2}{\rightarrow} \quad \text{NP PP}$$

$\text{VP}_{1,6} : w_1 \otimes c \otimes d$

$\text{VBD}_{1,2} : c \quad \text{NP}_{2,6} : d$

$\text{NP}_{2,6} : w_2 \otimes a \otimes b$

$\text{NP}_{2,4} : a \quad \text{PP}_{4,6} : b$

$$\frac{\text{VBD}_{1,2} : c \ \text{NP}_{2,6} : d}{\text{VP}_{1,6} : w_1 \otimes c \otimes d} : w_1$$

$$\frac{\text{NP}_{2,4} : a \ \text{PP}_{4,6} : b}{\text{NP}_{2,6} : w_2 \otimes a \otimes b} : w_2$$

- Associate weights as in WFST

- $\otimes$ : extension (multiplicative), $\oplus$ : summary (additive)

58

# Weights and Semirings



$$d(v) \quad = \quad (w(e_1, u_1, u_2) \otimes d(u_1) \otimes d(u_2))$$

$$\oplus \, (w(e_2, u_3, u_4) \otimes d(u_3) \otimes d(u_4))$$

- The weight of a hyperedge is dependent on antecedents (non-monotonic)

- The weight of a derivation is the product of hyperedge weights

- The weight of a vertex is the summary of (sub-)derivation weights

# Semirings

$$\mathbf{K} = \langle K, \oplus, \otimes, \mathbf{0}, \mathbf{1} \rangle$$

| semiring | K | $\oplus$ | $\otimes$ | 0 | 1 |
|---|---|---|---|---|---|
| Viterbi | [0,1] | max | × | 0 | 1 |
| Real | R | + | x | 0 | 1 |
| Log | R | logsumexp | + | $+\infty$ | 0 |
| Tropical | R | min | + | $+\infty$ | 0 |
| Expectation | <P,R> | <$p_1 \oplus p_2$, $r_1 \oplus r_2$> | <$p_1 \otimes p_2$, $p_1 \otimes r_2 \oplus p_2 \otimes r_1$> | <0,0> | <1,0> |

# Conclusion

- Review important concepts from "parsing"

  - CFG, parsing, hypergraph, deductive system, weights, semirings

# Tree-based MT

- Backgrounds

  - CFG, parsing, hypergraph, deductive system semirings

- **Tree-based SMT**

  - **Synchronous-CFG**

  - String-to-Tree, Tree-to-String

# Synchronous-CFG

$$X_{\boxed{1}} \qquad\qquad X_{\boxed{1}}$$

$$X_{\boxed{2}} \qquad X_{\boxed{3}} \qquad X_{\boxed{2}} \qquad X_{\boxed{3}}$$

布什 $X_{\boxed{4}}$ 举行 $X_{\boxed{5}}$ Bush held $X_{\boxed{5}}$ $\qquad X_{\boxed{4}}$

与 沙龙 $\qquad$ 了 会谈 $\qquad$ a talk with Sharon

(Chiang, 2007)

$$\hat{\mathbf{e}} \quad = \quad \underset{\mathbf{e}}{\arg\max} \frac{\exp\left(\mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, D, \mathbf{f})\right)}{\sum_{\mathbf{e}', D'} \exp\left(\mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}', D', \mathbf{f})\right)}$$

$$= \quad \underset{\mathbf{e}}{\arg\max}\, \mathbf{w}^{\top} \cdot \mathbf{h}(\mathbf{e}, D, \mathbf{f})$$

- D: a single derivation constructed by intersecting SCFG with input string

# Synchronous-CFG: Model

$$S \rightarrow \langle S_{\boxed{1}}\ X_{\boxed{2}}, S_{\boxed{1}}\ X_{\boxed{2}} \rangle$$

$$S \rightarrow \langle X_{\boxed{1}}, X_{\boxed{1}} \rangle$$

$$X \rightarrow \langle X_{\boxed{1}}\ 举行\ X_{\boxed{2}}, \text{hold}\ X_{\boxed{2}}\ X_{\boxed{1}} \rangle$$

$$X \rightarrow \langle 与\ 沙龙, \text{with Sharon} \rangle$$

$$VP \rightarrow \langle VBD_{\boxed{1}}\ NP_{\boxed{2}},\ NP_{\boxed{2}}\ VBD_{\boxed{1}} \rangle$$

$$NP \rightarrow \langle NP_{\boxed{1}}\ PP_{\boxed{2}},\ NP_{\boxed{1}}\ PP_{\boxed{2}} \rangle$$

$$VP \rightarrow \langle VBD_{\boxed{1}}\ NP_{\boxed{2}}\ PP_{\boxed{3}},\ NP_{\boxed{2}}\ PP_{\boxed{3}}\ VBD_{\boxed{1}} \rangle$$

- We use two categories, S and X (Chiang, 2007)

- Or, borrow linguistic categories from syntactic parse (Zollman and Venugopal, 2006)

# Rule Extraction

布什 与 沙龙举行 了 会谈

⟨held a talk with Sharon,
与 沙龙 举行 了 会谈⟩

⟨with Sharon, 与 沙龙⟩

⟨held, 举行⟩

$X \rightarrow \langle X_{1}\ X_{2}$ 了 会谈, $X_{2}$ a talk $X_{1} \rangle$

(Example from Huang and Chiang, 2007)

- As in phrase-based models, extract phrases then, use sub-phrases as non-terminals, aka Hiero (Chiang, 2007)

# Syntactic Categories

布什 与 沙龙举行 了 会谈

⟨held a talk with Sharon,
与 沙龙 举行 了 会谈⟩

⟨with Sharon, 与 沙龙⟩

⟨held, 举行⟩

VP → VBD a talk PP, PP VBD 了 会谈

VBD

VP

VBD+NP

NP

PP

Bush

held

a

talk

with

Sharon

- Borrow syntactic categories either from source/target side, aka SAMT (Zollman and Venugopal, 2006)

# Exhaustive Extraction

布什 与 沙龙举行 了 会谈

$X_{\boxed{1}}$ $X_{\boxed{2}}$ 了 会谈   $X_{\boxed{2}}$ a talk $X_{\boxed{1}}$

$X_{\boxed{1}}$ $X_{\boxed{2}}$ 会谈   $X_{\boxed{2}}$ a talk $X_{\boxed{1}}$

$X_{\boxed{1}}$ $X_{\boxed{2}}$ 会谈   $X_{\boxed{2}}$ talk $X_{\boxed{1}}$

$X_{\boxed{1}}$ 举行 $X_{\boxed{2}}$   held $X_{\boxed{2}}$ $X_{\boxed{1}}$

$X_{\boxed{1}}$ 举行 了 $X_{\boxed{2}}$   held a $X_{\boxed{2}}$ $X_{\boxed{1}}$

与 沙龙 $X_{\boxed{1}}$   $X_{\boxed{1}}$ with Sharon

与 $X_{\boxed{1}}$ $X_{\boxed{2}}$   $X_{\boxed{2}}$ with $X_{\boxed{1}}$

$S \rightarrow \langle S_{\boxed{1}} X_{\boxed{2}}, S_{\boxed{1}} X_{\boxed{2}} \rangle$

$S \rightarrow \langle X_{\boxed{1}}, X_{\boxed{1}} \rangle$

Rows (top to bottom): Bush, held, a, talk, with, Sharon

- Exhaustively extract rules as in phrase-based MT

- + glue rules

# Features from Rules

$$\log p_r(\bar{\alpha}|\bar{\beta}) = \log \frac{\text{count}(\bar{\beta}, \bar{\alpha})}{\sum_{\bar{\alpha}'} \text{count}(\bar{\beta}, \bar{\alpha}')}$$

$$\log p_r(\bar{\beta}|\bar{\alpha}) = \log \frac{\text{count}(\bar{\beta}, \bar{\alpha})}{\sum_{\bar{\beta}'} \text{count}(\bar{\beta}', \bar{\alpha})}$$

- Collect all the rules (α, β) from the data:
  - α = source side string, β = target side string
- Maximum likelihood estimates by relative frequencies
- Employ scores in two directions

# Remarks on Rules

- Too many rules extracted (Chiang, 2007):

  - at most two non-terminal symbols

  - at least one terminal between non-terminals in the source side

  - Span at most 15 words for "holes"

- Fractional counts (Chiang, 2007):

  - Each phrases counted in phrase-based MT

  - Fractional counts for rules sharing the same source/target span

# Other Features

- Lexical weights as used in phrase-based MT

- ngram language model(s)

- word count: bias for ngram language model(s)

- rule count: shorter or longer phrases

- glue-rule counts: bias for monotonic glue rules

# Synchronous-CFG: Parsing



$X \rightarrow \langle X_{\boxed{1}} \ X_{\boxed{2}} \ \text{le huitan}, \ X_{\boxed{2}} \ \text{a talk} \ X_{\boxed{1}} \rangle$

$X \rightarrow \langle X_{\boxed{1}} \ \text{juxing} \ X_{\boxed{2}}, \ \text{held} \ X_{\boxed{2}} \ X_{\boxed{1}} \rangle$

- Parse input sentence using the source side, and construct a translation forest by target side

# Synchronous-CFG: Parsing

- Translation by SCFG = monolingual parsing using the source side grammar

- Construct forest by the projected target side

- From forests, compute the best derivation (Huang and Chiang, 2005)

- Complexity: $O(n^3)$ as in monolingual CKY

# Non-Local Features

$$X \rightarrow \langle X_{\boxed{1}} \text{ juxing } X_{\boxed{2}},$$
$$\text{held } X_{\boxed{2}} \ X_{\boxed{1}} \rangle$$

$X_{1,6}$

held a talk with Sharon

held talks with Sharon

held a talk and Sharon

held meeting Sharon with

Update boundary words only

$X_{3,4}$

$X_{1,3}$

P(talk | a)   a talk

talks

meeting

meetings

with Sharon   P(Sharon | with)

and Sharon   P(Sharon | and)

Sharon with   P(with | Sharon)

Sharon and   P(and | Sharon)

- non-local features which requires out-of-span context, i.e. bigram LM

# Bigram Features

$$X \rightarrow \langle X_{\boxed{1}} \text{ juxing } X_{\boxed{2}},$$
$$\text{held } X_{\boxed{2}} \; X_{\boxed{1}} \rangle$$

held * Sharon

held * Sharon

held * Sharon

held * with

$X_{1,6}$

$X_{3,4}$

a * talk

talks

meeting

meetings

$X_{1,3}$

with * Sharon

and * Sharon

Sharon * with

Sharon * and

- We keep only bigram states: (Why 2 words?)

# Language Model Updates

- Each hypothesis keeps two contexts:

  - Prefix: ngrams to be scored with antecedents

  - Suffix: contexts for future ngrams (i.e. Phrase-based MT)

- Complexity: $O(n^3 V^{2(m-1)})$

- Very inefficient: we need to explicitly enumerate all the hypotheses in antecedents

# Forest Rescoring

- Translation by SCFG = monolingual parsing using the source side grammar

- Construct forest by the projected target side   + Rescore with non-local features

- From forests, compute the best derivation (Huang and Chiang, 2005)

- ~~Complexity: $O(n^3)$ as in monolingual CKY~~

# Cube Pruning

$$X \rightarrow \langle X_{\boxed{1}} \text{ juxing } X_{\boxed{2}},$$
$$\text{held } X_{\boxed{2}} X_{\boxed{1}} \rangle$$

|  |  | with * Sharon 1.5 | and * Sharon 1.7 | Sharon * with 2.6 | Sharon * and 3.2 |
|---|---|---|---|---|---|
| a * talk | 1.0 | 2.5 | 2.7 | 3.6 | 4.2 |
| talks | 1.3 | 2.8 | 3.0 | 3.9 | 4.5 |
| meeting | 2.2 | 3.7 | 3.9 | 4.8 | 5.4 |
| meetings | 2.6 | 4.1 | 4.3 | 5.2 | 5.8 |

- For each hyperedge, create a "cube" representing combinations of antecedents (Huang and Chiang, 2007)

# Cube Pruning

$$X \rightarrow \langle X_{\boxed{1}} \text{ juxing } X_{\boxed{2}},$$
$$\text{held } X_{\boxed{2}} \ X_{\boxed{1}} \rangle$$

|  |  | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
|  |  | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 2.5 +0.5 | 2.7 +1.0 | 3.6 +1.5 | 4.2 +1.5 |
| talks | 1.3 | 2.8 +0.3 | 3.0 +1.5 | 3.9 +2.0 | 4.5 +2.0 |
| meeting | 2.2 | 3.7 +0.5 | 3.9 +1.0 | 4.8 +1.5 | 5.4 +1.5 |
| meetings | 2.6 | 4.1 +0.3 | 4.3 +1.5 | 5.2 +2.0 | 5.8 +2.0 |

- Bigrams require contexts from antecedents: non-monotonic scoring

# Cube Pruning

queue: (0,0)

k-best:

| | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
| | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | | | |
| talks | 1.3 | | | | |
| meeting | 2.2 | | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue:

k-best: (0,0)

|  | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
|  | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | | | |
| talks | 1.3 | | | | |
| meeting | 2.2 | | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (0,1)(1,0)
k-best: (0,0)

| | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
| | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | | |
| talks | 1.3 | 3.1 | | | |
| meeting | 2.2 | | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (1,0)

k-best: (0,0)(0,1)

|  | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
|  | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | | |
| talks | 1.3 | 3.1 | | | |
| meeting | 2.2 | | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

82

# Cube Pruning

queue: (1,0) (0,2) (1,1)

k-best: (0,0) (0,1)

| | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
| | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | | |
| talks | 1.3 | 3.1 | 4.5 | | |
| meeting | 2.2 | 4.2 | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (0,2) (1,1)

k-best: (0,0) (0,1) (1,0)

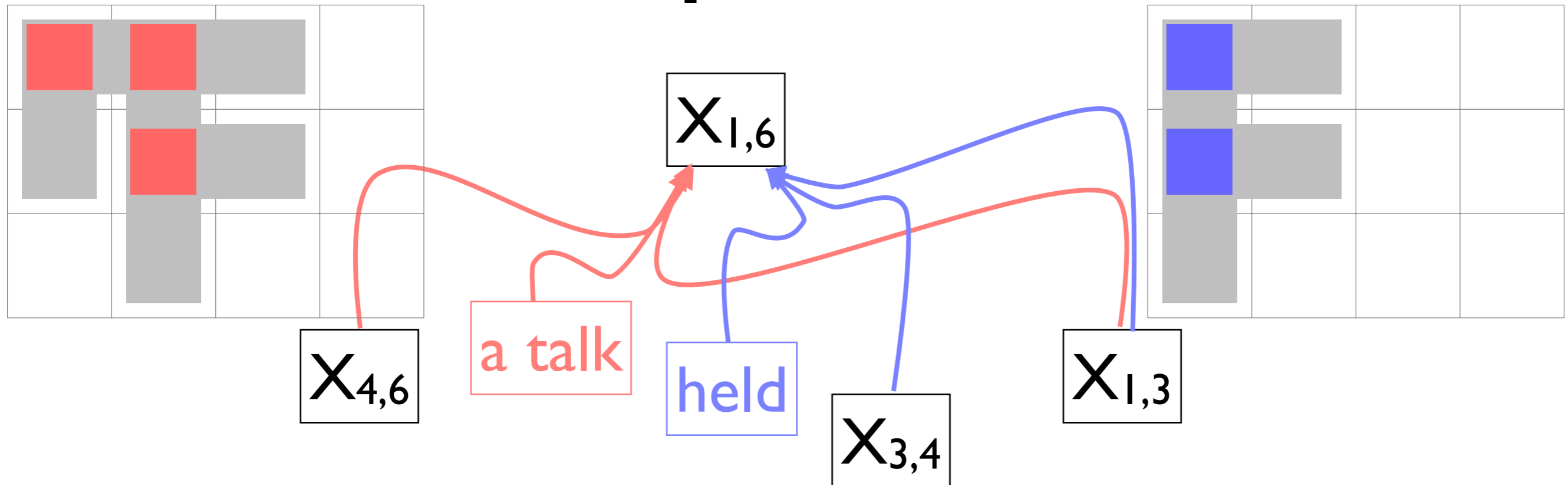|  |  | with * Sharon<br>1.5 | and * Sharon<br>1.7 | Sharon * with<br>2.6 | Sharon * and<br>3.2 |
|---|---|---|---|---|---|
| a * talk | 1.0 | 3.0 | 3.7 |  |  |
| talks | 1.3 | 3.1 | 4.5 |  |  |
| meeting | 2.2 | 4.2 |  |  |  |
| meetings | 2.6 |  |  |  |  |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (0,2) (1,1) (3,0)

k-best: (0,0) (0,1) (1,0)

| | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
| | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | 5.1 | |
| talks | 1.3 | 3.1 | 4.5 | | |
| meeting | 2.2 | 4.2 | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (1,1) (3,0)

k-best: (0,0) (0,1) (1,0) (0,2)

|  | | with Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
|  | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | 5.1 | |
| talks | 1.3 | 3.1 | 4.5 | | |
| meeting | 2.2 | 4.2 | | | |
| meetings | 2.6 | | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Cube Pruning

queue: (0,4) (1,1) (1,2) (3,0)

k-best: (0,0) (0,1) (1,0) (0,2)

| | | with * Sharon | and * Sharon | Sharon * with | Sharon * and |
|---|---|---|---|---|---|
| | | 1.5 | 1.7 | 2.6 | 3.2 |
| a * talk | 1.0 | 3.0 | 3.7 | 5.1 | |
| talks | 1.3 | 3.1 | 4.5 | | |
| meeting | 2.2 | 4.2 | 4.9 | | |
| meetings | 2.6 | 4.4 | | | |

- Starting from the upper-left corner, enumerate antecedent combinations

# Multiple Rules



- Multiple rules sharing the same span are queued

  - Each rule is associated with a cube

  - hypothesis = hyperedge + cube-position

# Further Faster Pruning

- Cube Growing (Huang and Chiang, 2007)

  - Top-down pruning combined with heuristic estimates

- Faster Cube Pruning (Gesmundo and Henderson, 2010)

  - Eliminate bookkeeping for inserted hypotheses by determining the ordering of cube enumerations

  - Push minimum hypotheses by looking up ancestors

- Incremental (Huang and Mi, 2010)

  - Top-down decoding as in (Watanabe et al., 2006)

# Conclusion

- Synchronous-CFG

  - paired CFG + shared non-terminal symbols

- Training is based on phrase-based MT by treating sub-phrase as a non-terminal

- Decoding: monolingual parsing

  - An efficient antecedent combination via cube-pruning

# Tree-based MT

- Backgrounds

  - CFG, parsing, hypergraph, deductive system semirings

- Tree-based SMT

  - Synchronous-CFG

  - **String-to-Tree, Tree-to-String**

# {Tree,String}-to-{Tree,String}



(Galley et al., 2004)

- Each synchronous rule has a subtree structure

- Flat structure + sharing the same non-terminal symbols = synchronous-CFG
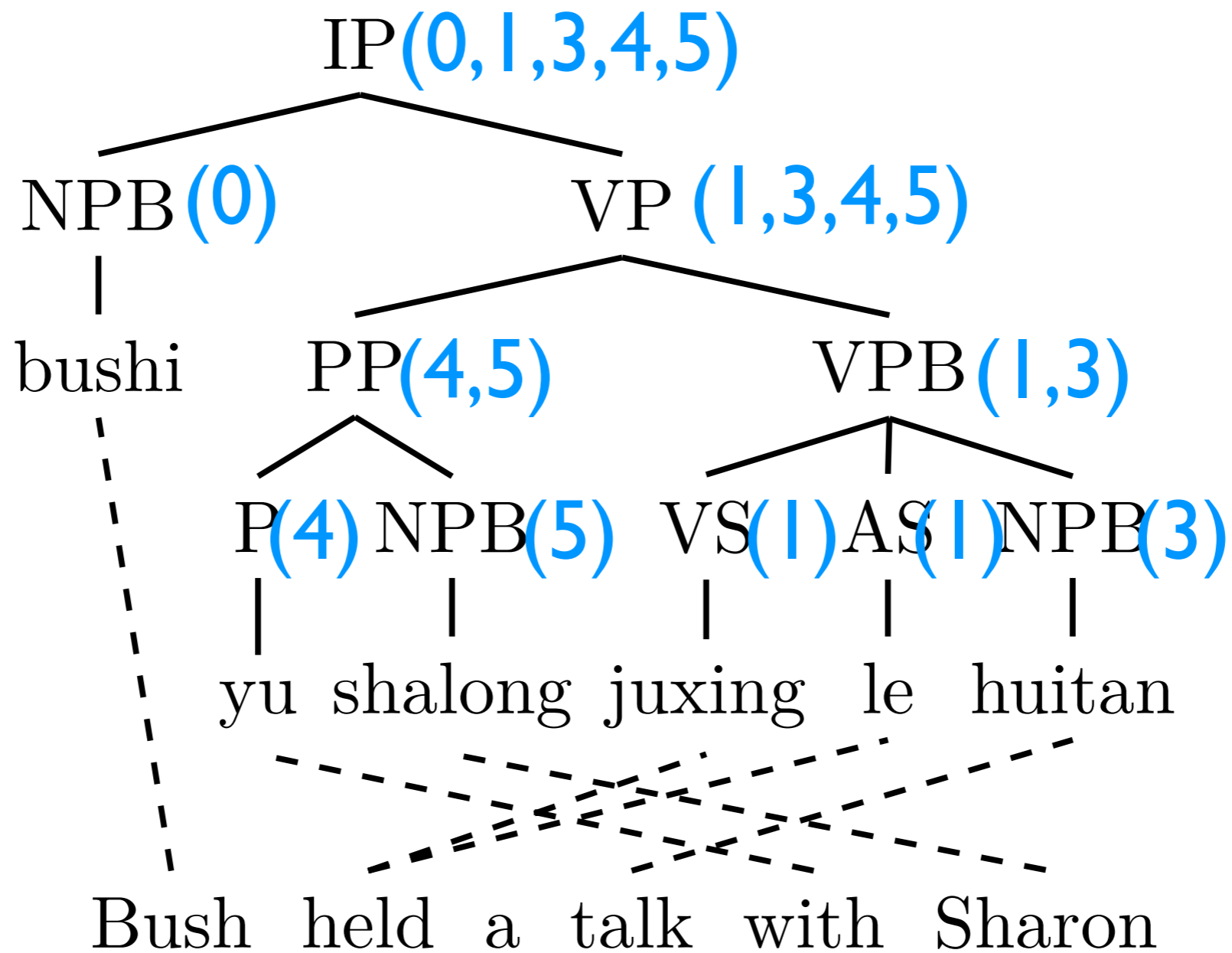
# Tree-to-String Rules

PP
├ P
│ └ yu → with
└ NPB
  └ shalong → Sharon

VPB
├ VS
│ └ juxing → held
├ AS
│ └ le → a
└ $x_1$:NPB → $x_1$

PP
├ P
│ └ dang → when
└ LCP
  ├ $x_1$:IP
  └ LC
    └ hou
→ when $x_1$

QP
├ $x_1$:CD
└ CLP
  └ ben
→ $x_1$

IP
├ $x_1$:NP
└ VP
  ├ $x_2$:IP
  └ $x_3$:VPB
→ $x_1$ $x_3$ $x_2$

NP
├ DNP
│ ├ $x_1$:NP
│ └ DEG
│   └ de
└ $x_2$:NP
→ $x_2$ of $x_1$

# ルールの抽出



(Galley et al., 2004)

- Compute "minimum rules" as in phrase-based MT

94

# Rule Extraction



IP(0,1,3,4,5)

NPB (0)    VP (1,3,4,5)

bushi    PP(4,5)    VPB(1,3)

P(4) NPB(5)  VS(1)AS(1)NPB(3)

yu shalong juxing le huitan

Bush  held  a  talk  with  Sharon

(Galley et al., 2004)

- Compute "spans" by propagating alignment in bottom-up

95

# Rule Extraction

IP(0,1,3,4,5) ()

NPB (0)(1,3,4,5)  VP (1,3,4,5)(0)

bushi

PP(4,5)(0,1,3)  VPB(1,3)(0,4,5)

(0,1,3,5)P(4) NPB(5)  VS(1) AS1 NPB(3)(0,1,4,5)
(0,1,3,4)(0,1,3,4,5)(0,1,3,4,5)

yu  shalong  juxing  le  huitan

Bush  held  a  talk  with  Sharon

- Compute "complements" in top-down

# Rule Extraction



- Compute "frontiers": The nodes in which the intersection of "spans" and "complements" is empty
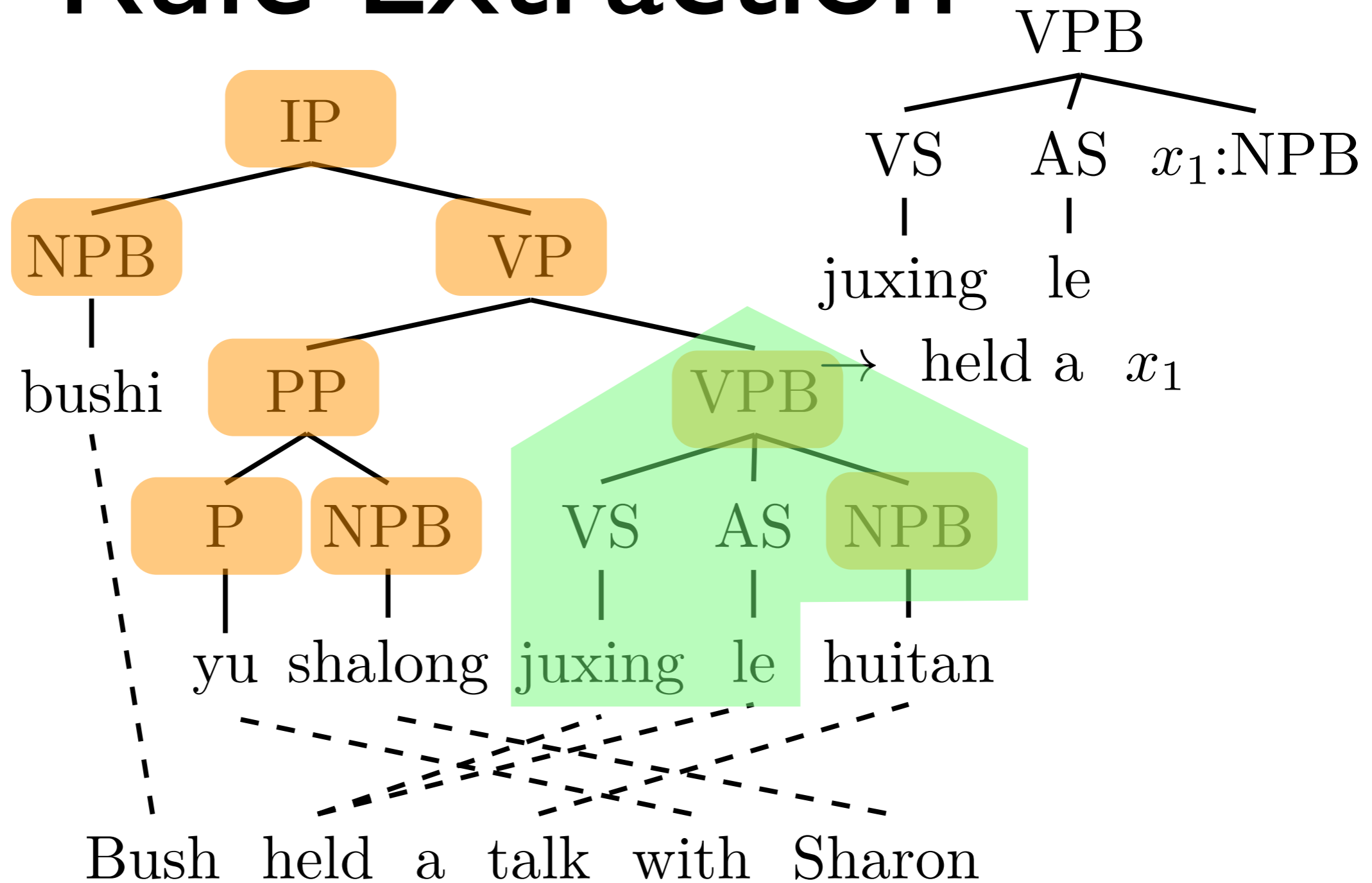
# Rule Extraction



$$\text{IP} \to x_2\ x_1$$

with tree:
- IP
  - $x_1$:NPB  $x_1$:VP

- Extract minimum rules using frontiers

# Rule Extraction



- Extract minimum rules using frontiers

# Rule Extraction



- Extract minimum rules using frontiers

# Rule Extraction
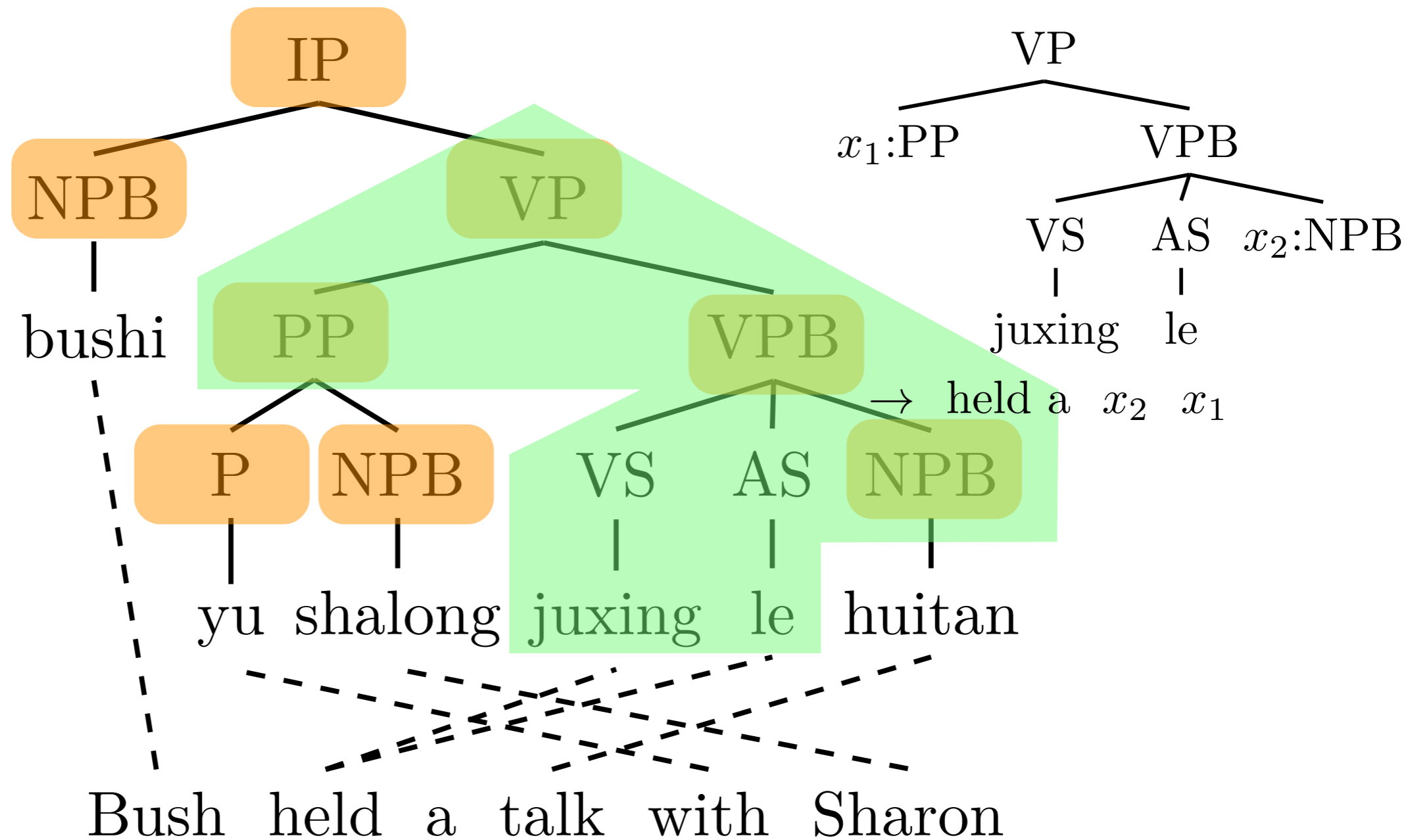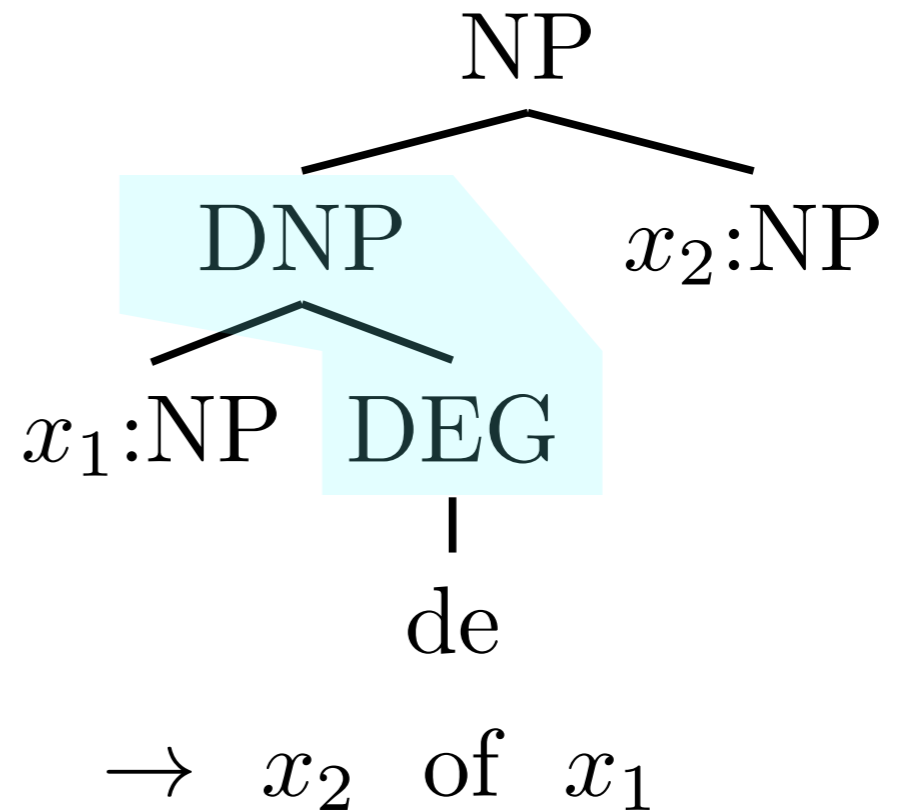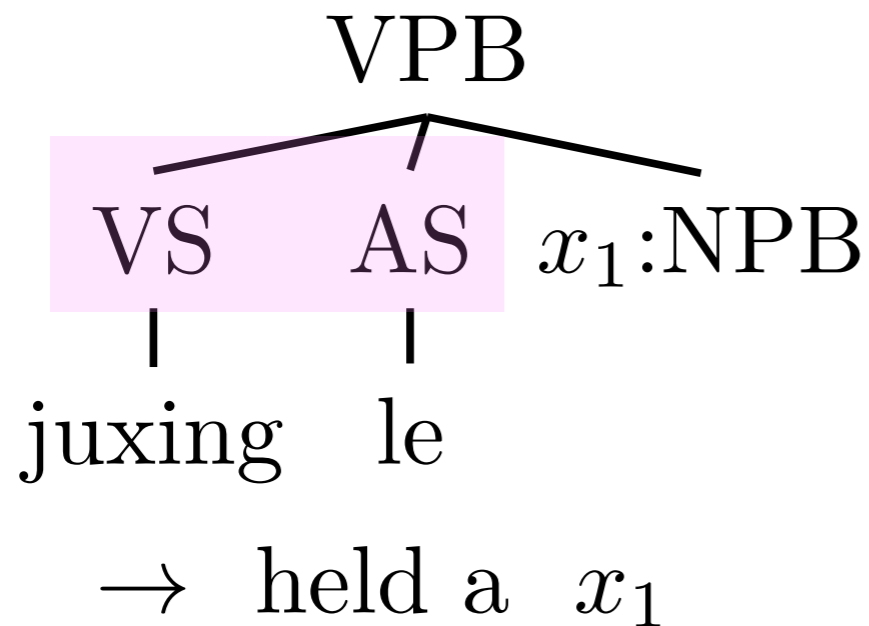


VP
  x₁:PP    VPB
         VS   AS   x₂:NPB
         juxing  le
→ held a  x₂  x₁

(Galley et al., 2006)

- Extract "compound rules" by combining minimum rules (i.e. longer phrases)

# Decoding: String-{String, Tree}



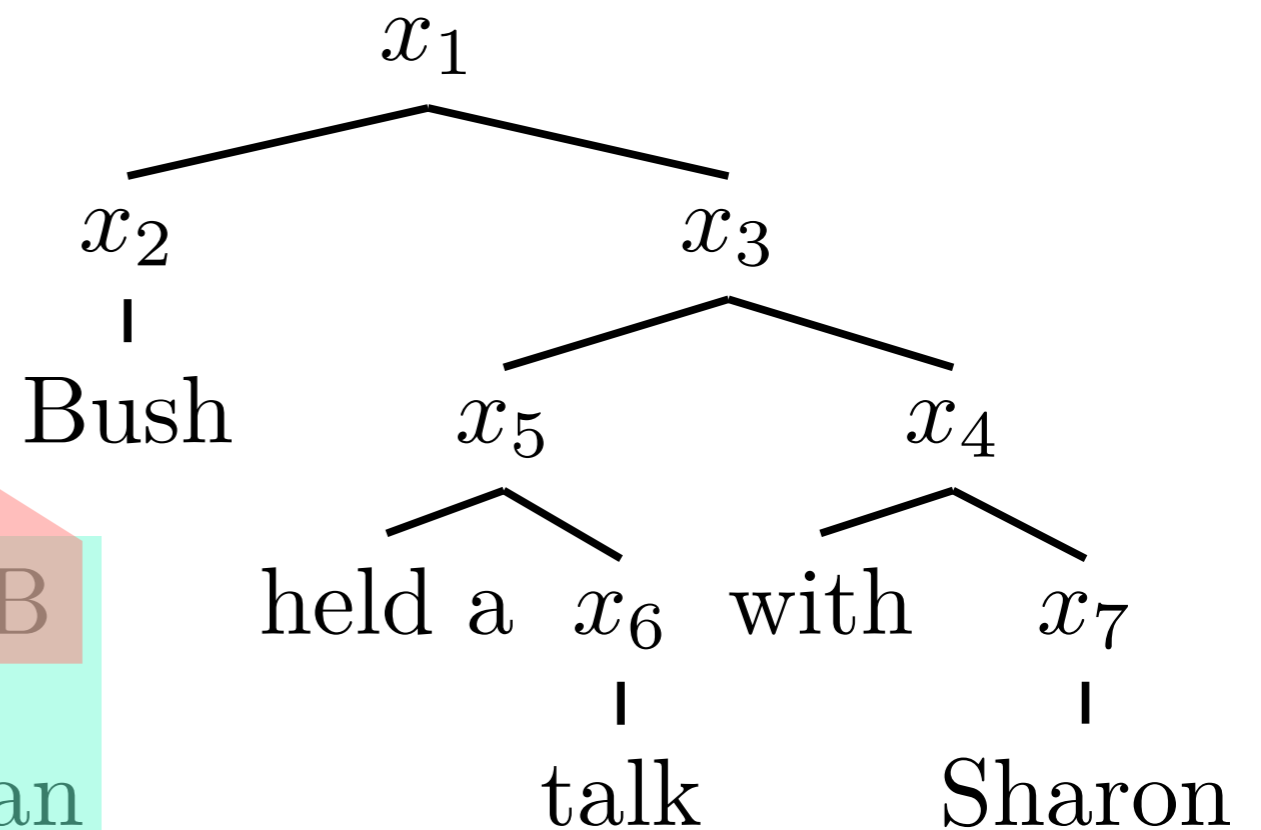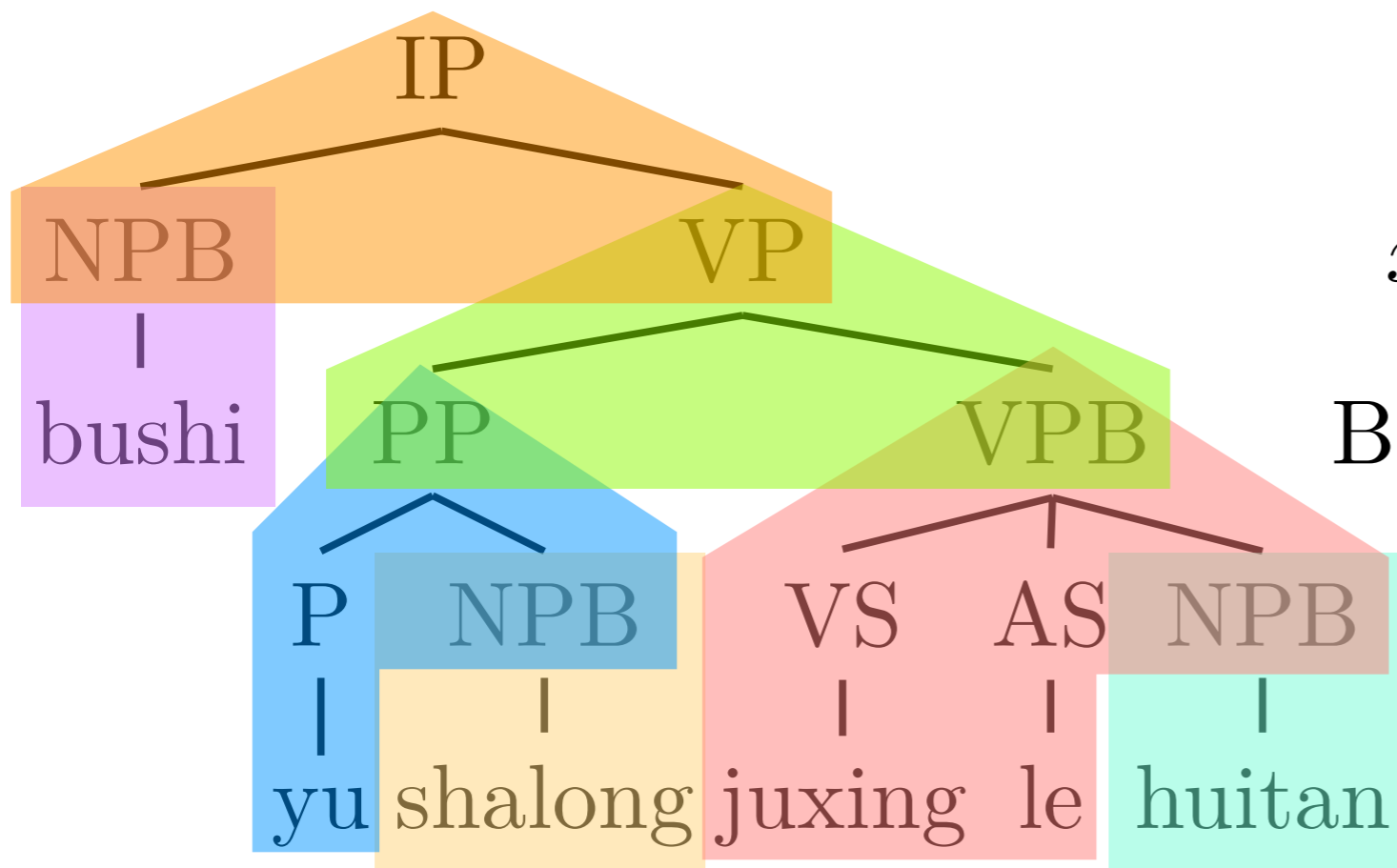$$\langle \text{VPB} \rightarrow \text{juxing le NPB}_1,$$
$$x \rightarrow \text{hold a } x_1 \rangle$$

$$\langle \text{NP} \rightarrow \text{NP}_1 \text{ de NP}_2,$$
$$x \rightarrow x_2 \text{ of } x_1 \rangle$$

(Galley et al., 2004)

- Similar to SCFG decoding: Use the "collapsed" source side rule to perform CKY parsing

- Construct a translation forest using the target side

# Decoding: Tree-{String, Tree}



(Huang et al., 2006)

- First, an input sentence is parsed

- Input tree is transformed into a translation forest by tree rewriting

# Forest Rescoring

- Translation by {tree,string}-to-{tree,string}

  - string-to-{tree,sting}: parsing using the source-side grammar

  - tree-to-{tree,string}: parse input sentences + tree-match-rewrite

- Construct forest by the projected target side

- From forests, compute the best derivation (Huang and Chiang, 2005)

# Conclusion

- {String, Tree}-to-{String, Tree} translation models

- Rules extraction by GHKM (Galley et al., 2004)

  - Galley M, Hopkins M, Knight K, Marcu D, 2004

- Decoding:

  - String-to-{String, Tree} by CKY

  - Tree-to-{String, Tree} by tree-rewrite

# More on Tree-based Models

- Forest-based approach: instead of 1-best parse, use forest encoding k-bests (Mi and Huang, 2008; Mi et al., 2008)

- "Binarized forest" as an alternative to represent multiple parses (Zhang et al., 2011)

- Fuzzy tree-to-tree as a way to overcome "stricktness" of tree-based models (Chiang, 2010)

- Use of dependency (Mi and Liu, 2010; Xie et al., 2011)

# Tree-based MT

- Backgrounds

  - CFG, parsing, hypergraph, deductive system semirings

- Tree-based SMT

  - Synchronous-CFG

  - String-to-Tree, Tree-to-String

# Structures in SMT

- Tutorial

  - Phrase-based MT

  - Tree-based MT

- Syntactic Structures in System Combination

# MT System Combination by Confusion Forest

Taro Watanabe and Eiichiro Sumita
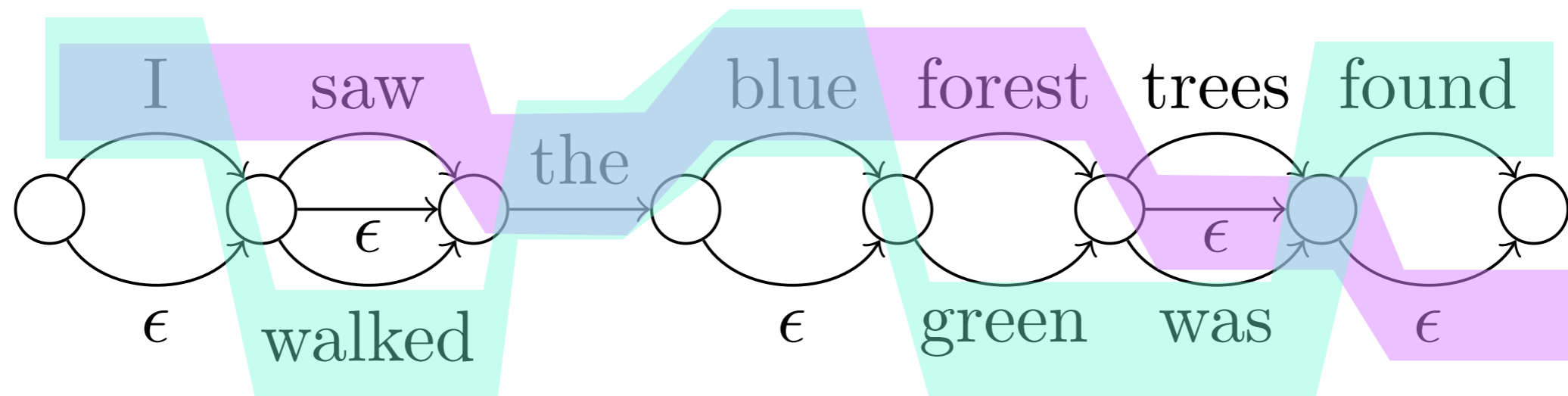@ NICT

# MT System Combination

- Better translation by combining multiple system outputs:

  - Sentence selection(Nomoto, 2004; etc.)

  - Phrasal combination (Frederking and Nirenburg, 1994; etc.)

  - Word level combination (Bangalore et al., 2001; Matusov et al., 2006; etc.)

- This Work: Syntactic combination, not word-wise combination

# Confusion Network

★
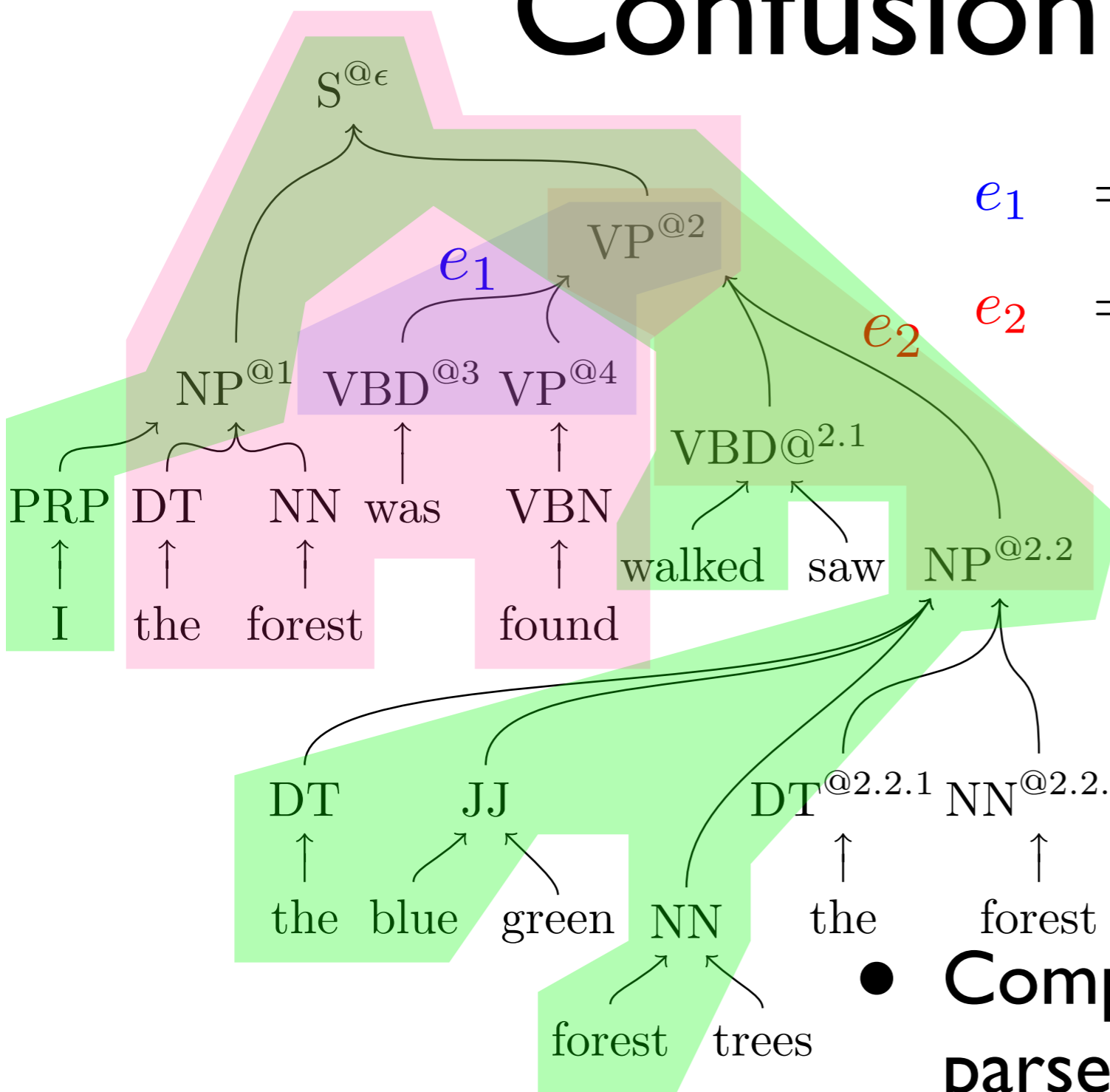|  | I | saw | the | forest |
|---|---|---|---|---|
| I | walked | the blue | forest |
| I | saw | the | green trees |
| the | forest | was | found |

- State-of-the-art: Confusion Network

- Choose a skeleton, compute word alignment against the skeleton

  - Edit-distance-based alignment (TER etc.) (Sim et al., 2007)

  - Model-based alignment(GIZA++ etc.) (Matsov et al., 2006)

# Confusion Network



- Construct a network with each arc representing alternative translation

  - Best path = Best translation

  - Syntactically different language pairs: i.e. active/passive voices

  - Spurious insertion/repetition due to alignment error

  - Incremental alignment/construction + merge multiple networks into one (Rosti et al., 2008)
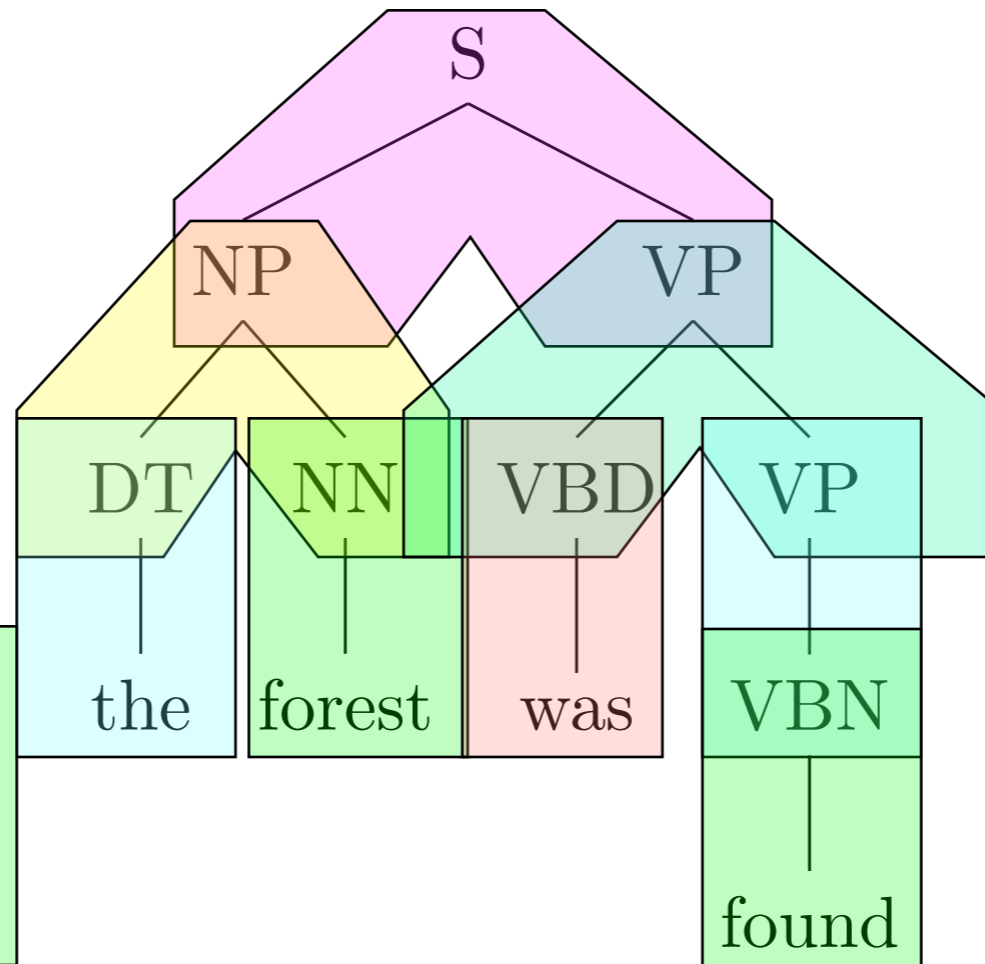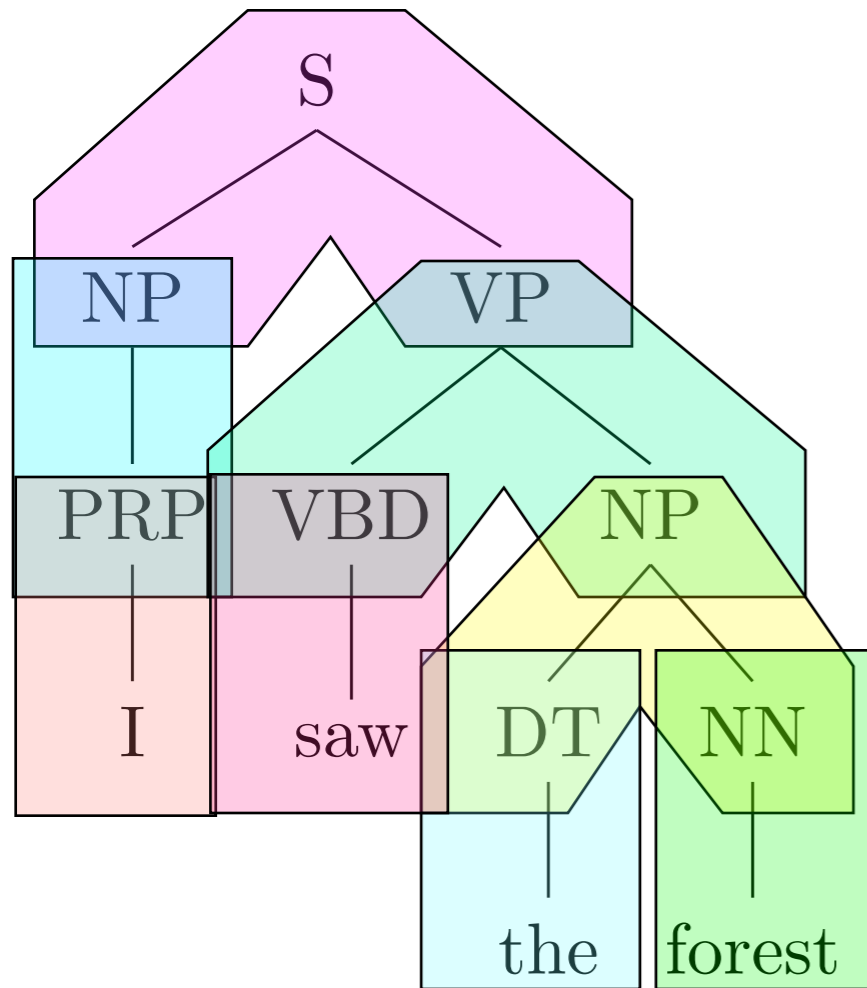
# Confusion Forest



$$e_1 = \langle \mathrm{VP}^{@2}, \{\mathrm{VBD}^{@3}, \mathrm{VP}^{@4}\} \rangle$$

$$e_2 = \langle \mathrm{VP}^{@2}, \{\mathrm{VBD}^{@2.1}, \mathrm{NP}^{@2.2}\} \rangle$$

- Compactly represent multiple parses by sharing nodes

- Represented by "hypergraph"

# Rule Extraction



S → NP VP
NP → PRP
PRP → I
VP → VBD NP
VBD → saw
NP → DT NN
DT → the
NN → forest
VP → VBD VP
VBD → was
VP → VBN
VBN → found

- Parse each system output by a parser

- Extract rules from parsed trees: local grammar

# Generation by Earley

Scan:

$$\frac{[\mathrm{X} \to \alpha \bullet x\beta, h] : u}{[\mathrm{X} \to \alpha x \bullet \beta, h] : u}$$
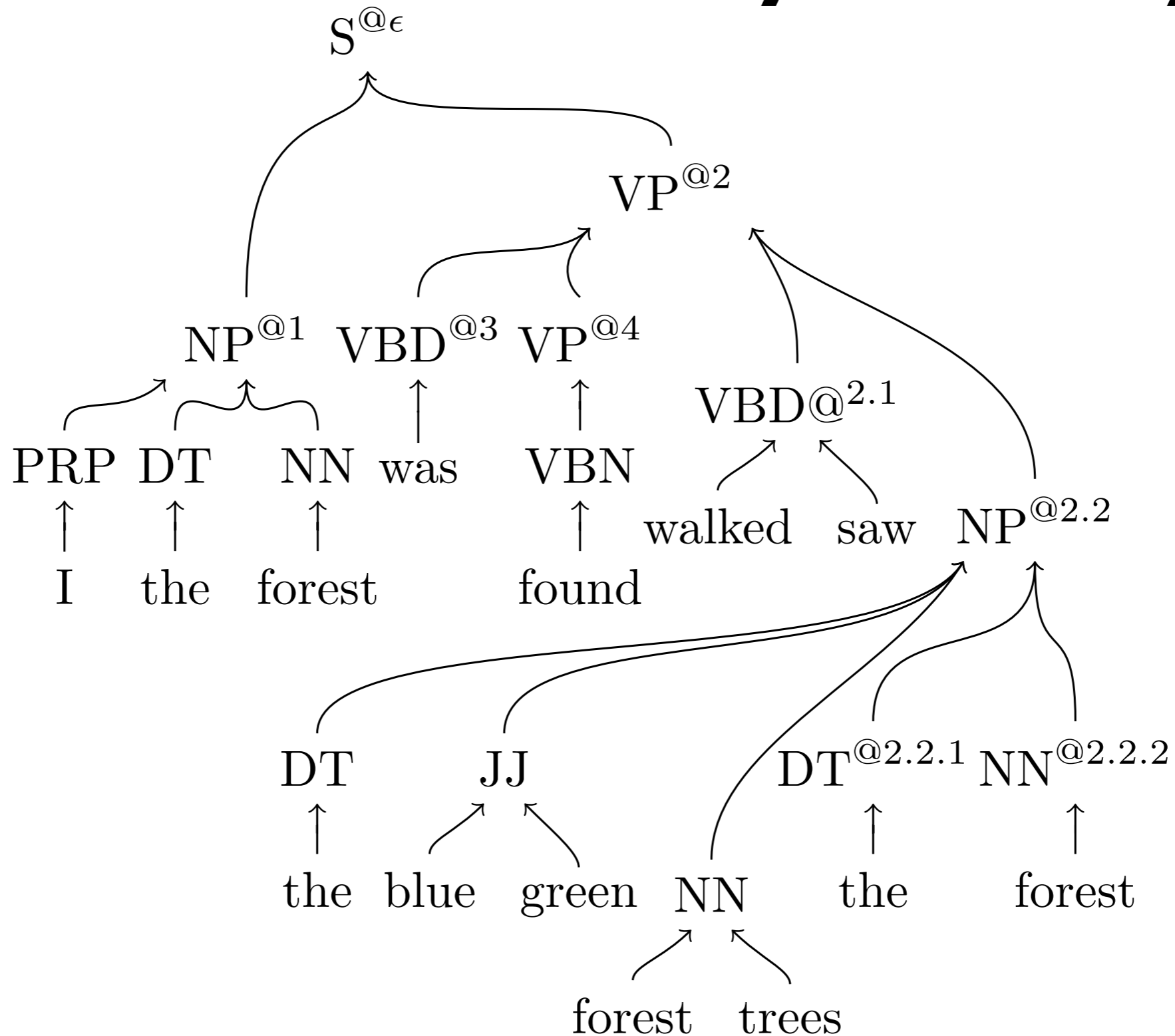
Predict:

$$\frac{[\mathrm{X} \to \alpha \bullet \mathrm{Y}\beta, h]}{[\mathrm{Y} \to \bullet\gamma, h+1] : u} \quad \mathrm{Y} \xrightarrow{u} \gamma \in \mathcal{G}, h < H$$

Complete:

$$\frac{[\mathrm{X} \to \alpha \bullet \mathrm{Y}\beta, h] : u \quad [\mathrm{Y} \to \gamma\bullet, h+1] : v}{[\mathrm{X} \to \alpha\mathrm{Y} \bullet \beta, h] : u \otimes v}$$
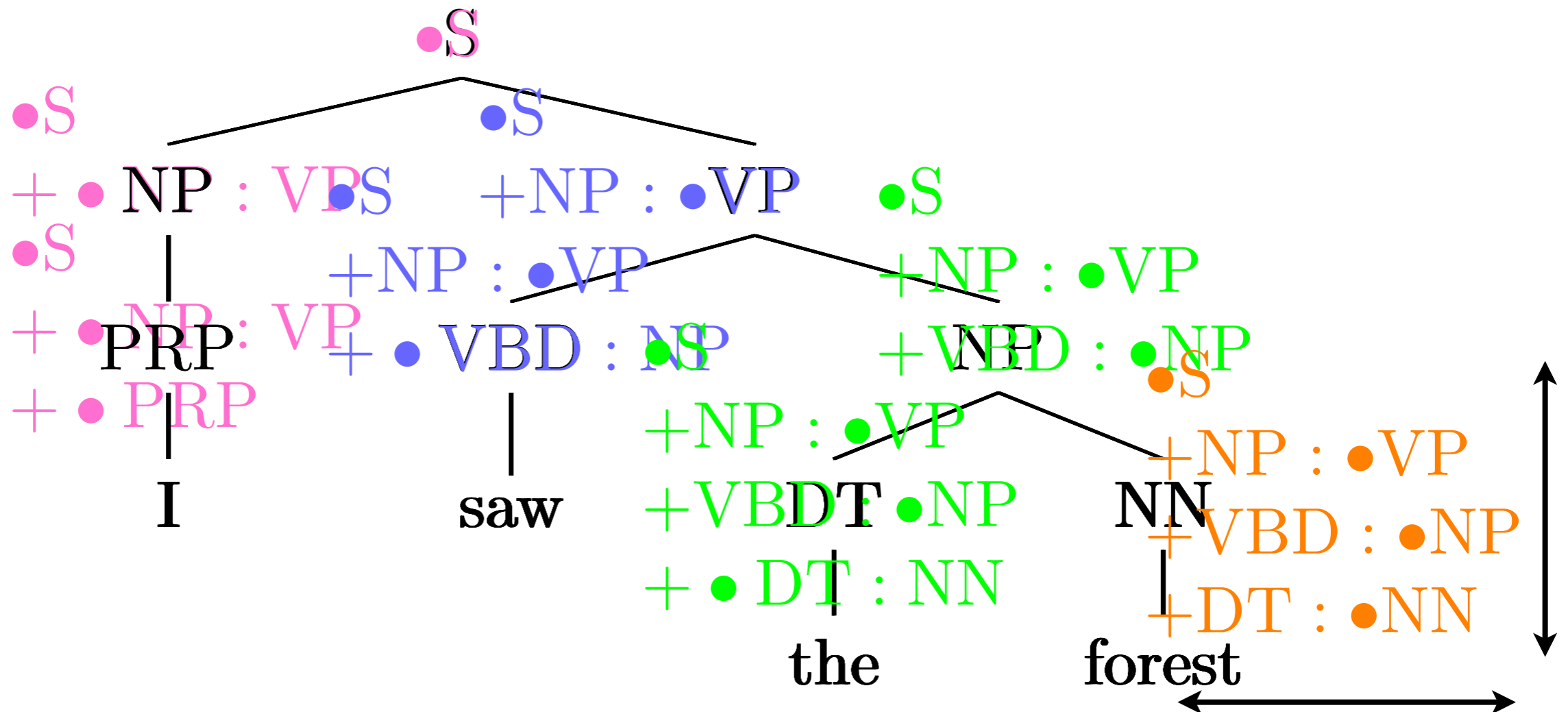
- Generation from the extracted grammar

- Scanning always succeed: constraint by height

# Generation by Earley

# Spurious Ambiguity



- Memorize the (partial) tree structures in each node

- Employ the sequence of Ealrye state as a node

  - Horizontal/Vertical Markovization (Klein and Manning, 2003)

# Forest Reranking

$$\hat{d} = \arg\max_{d \in D} \mathbf{w}^\top \cdot \mathbf{h}(d, F)$$

- Choose the best derivation d among all possible derivations D in a forest F

  - Terminal yield of the best derivation = the best translation

  - Approximately apply non-local features (ngram language models) by Cube Pruning (Huang and Chiang, 2007)

  - Efficient k-best by Algorithm 3 (Huang and Chiang, 2005)

# Experiments

- WMT10 System Combination Task
  - Czech, German, Spanish, French → English
  - tune/test: 455/2,034 sentences

|         | cz-en  | de-en  | es-en  | fr-en  |
|---------|--------|--------|--------|--------|
| systems | 6      | 16     | 8      | 14     |
| tune    | 10.6K  | 10.9K  | 10.9K  | 11.0K  |
| test    | 50.5K  | 52.1K  | 52.1K  | 52.4K  |

# Systems

- CF: Stanford parser + "cicada" (a hypergraph-based toolkit based on SEMIring parsing framework)

- CN: Single network by merging multiple networks + conversion into hypergraph by lattice parsing

- features: tuned by hypergraph-MERT(Kumar et al. 2009)

  - Language Models, # of terminals, # of hyperedges

  - # of rules in a derivation originally in $n_{th}$ system output

  - BLEUs by treating each system output as a reference translation

  - Network distance (only used for CN)

# BLEU

| | cz-en | de-en | es-en | fr-en |
|---|---|---|---|---|
| system min | 14.09 | 15.62 | 21.79 | 16.79 |
| max | 23.44 | 24.10 | 29.97 | **29.17** |
| CN | 23.70 | 24.09 | **30.45** | **29.15** |
| CF,$v=\infty$,$h=\infty$ | **24.13** | 24.18 | **30.41** | **29.57** |
| CF,$v=\infty$,$h=2$ | **24.14** | **24.58** | **30.52** | 28.84 |
| CF,$v=\infty$,$h=1$ | **24.01** | 23.91 | **30.46** | **29.32** |

# Oracle BLEU

|  | cz-en | de-en | es-en | fr-en |
|---|---|---|---|---|
| rerank | 29.40 | 32.32 | 36.83 | 36.59 |
| CN | 38.52 | 34.97 | 47.65 | 46.37 |
| CF,v=∞,h=∞ | 30.51 | 34.07 | 38.69 | 38.94 |
| CF,v=∞,h=2 | 30.61 | 34.25 | 38.87 | 39.10 |
| CF,v=∞,h=1 | 31.09 | 34.65 | 39.27 | 39.51 |

# Hypergraph size

|  | cz-en | de-en | es-en | fr-en |
|---|---|---|---|---|
| CN | 2,222.68 | 47,231.20 | 2,932.24 | 11,969.40 |
| CF,v=∞,h=1 | 230.08 | 540.03 | 262.30 | 386.79 |
| CF,v=5,h=1 | 254.45 | 651.10 | 302.01 | 477.51 |
| CF,v=4,h=1 | 286.01 | 802.79 | 349.21 | 575.17 |

- Average # of hyperedges
- (rough) estimates for speed

# Conclusion

- System combination by Confusion Forest which employs syntactic distance, not word-level distance

- Forest construction by the grammar extracted from system outputs

  - Parser: assign tree structure to the similar expressions

- Compact date structure + comparable performance against Confusion Network

# Structures in SMT

- Tutorial
  - Phrase-based MT
  - Tree-based MT
- Syntactic Structures in System Combination

# Research on MT

- Reading: at least 50 papers are related to MT "every year"

- Specialist: solve a sub-problem

- Language Neutral: a solution which works only for a particular language pair is boring