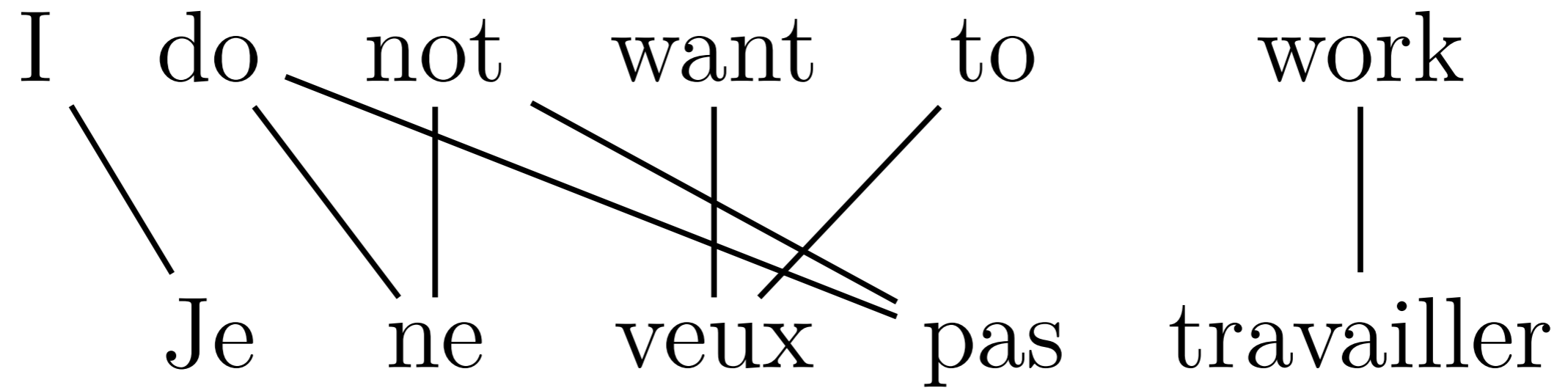


# Tree-based Models for SMT

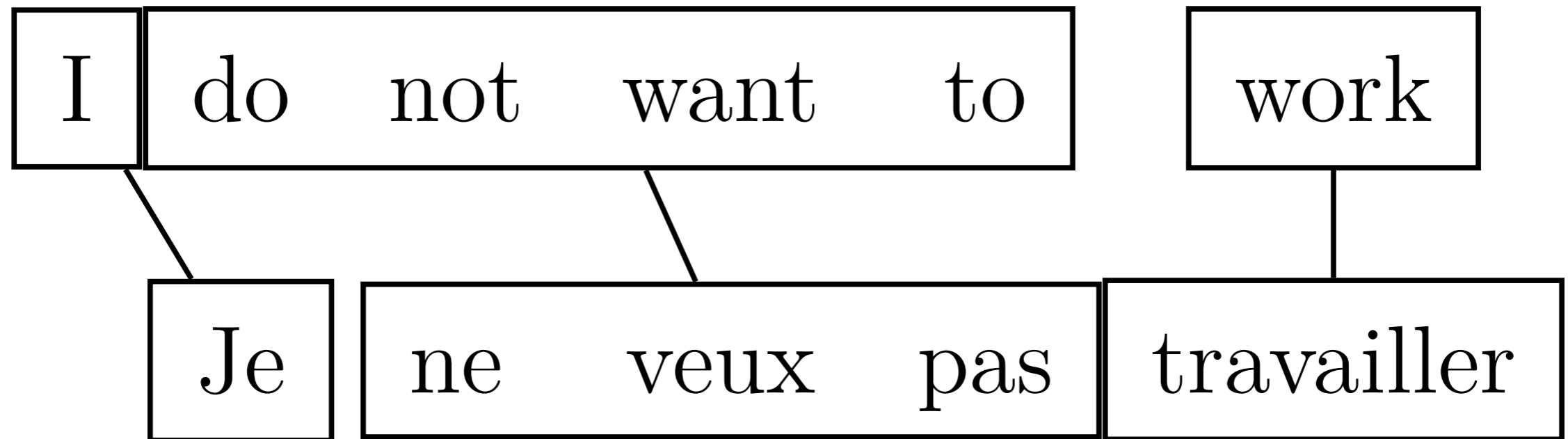
Taro Watanabe

# Word-based MT



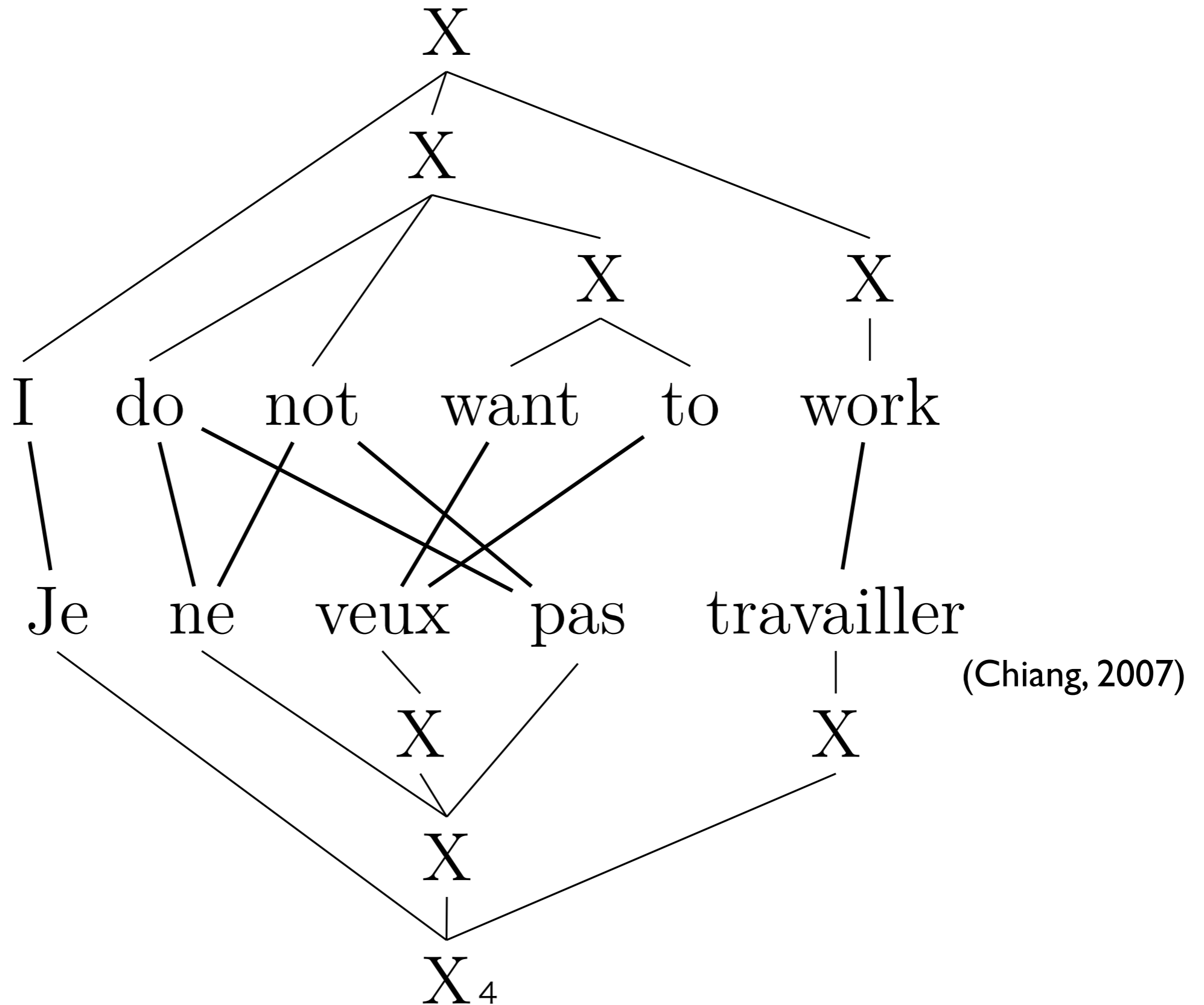
(Brown et al., 1993)

# Phrase-based MT

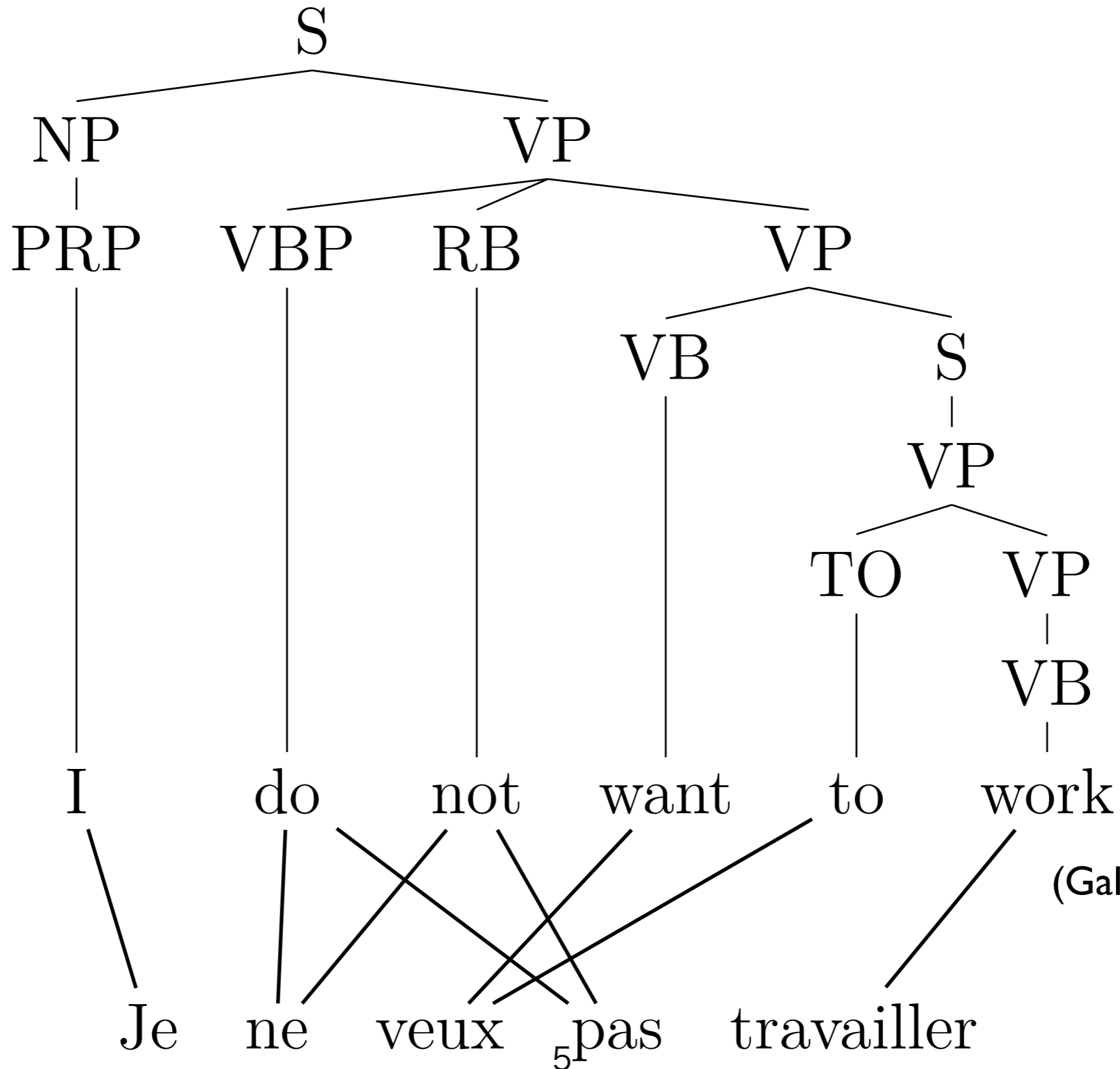


(Koehn et al., 2003)

# Hierarchical PBMT



# Syntax-based MT



(Galley et al., 2004)

# Variations

tree	(partial) examples
none	Chiang (2007), Zollman and Venugopal (2006)
source	Huang et al. (2006), Liu et al. (2006), Quirk et al. (2005)
target	Galley et al. (2004), Shen et al. (2008)
both	Ding and Palmer (2005), Liu et al. (2009)

- formally syntactical, linguistically syntactical
- dependency structure and constituency structure
- {tree,string}-to-{tree,string}
- In this talk, we will summarize them as “tree-based MT”

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing



# Backgrounds: CFG

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

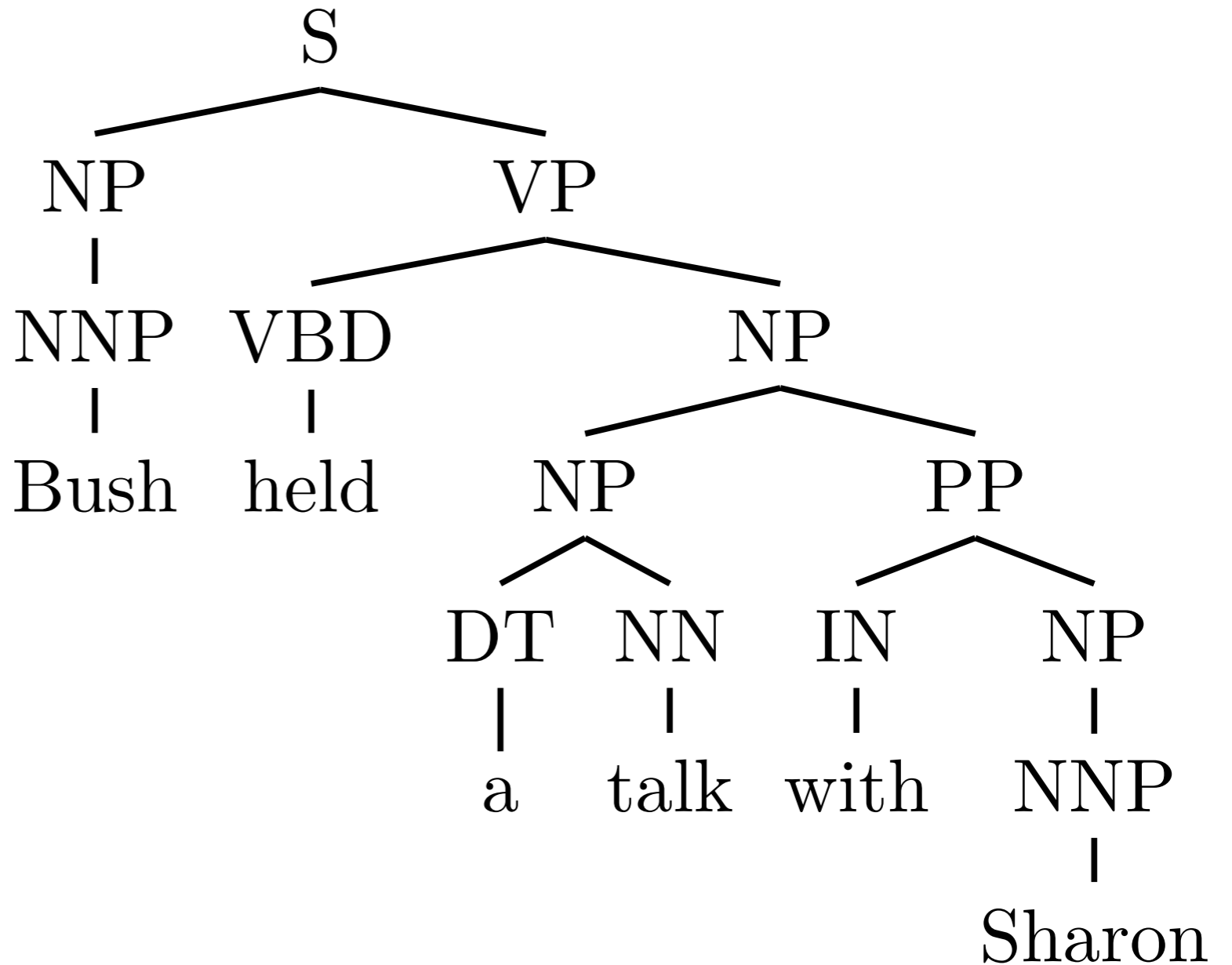
# Backgrounds: CFG

S → NP VP  
NP → NNP  
NP → NP PP  
NP → DP NN  
NP → DT NN  
VP → VBD NP  
NNP → Bush  
VBD → held  
⋮

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

# Backgrounds: CFG

S → NP VP  
 NP → NNP  
 NP → NP PP  
 NP → DP NN  
 NP → DT NN  
 VP → VBD NP  
 NNP → Bush  
 VBD → held  
 ⋮

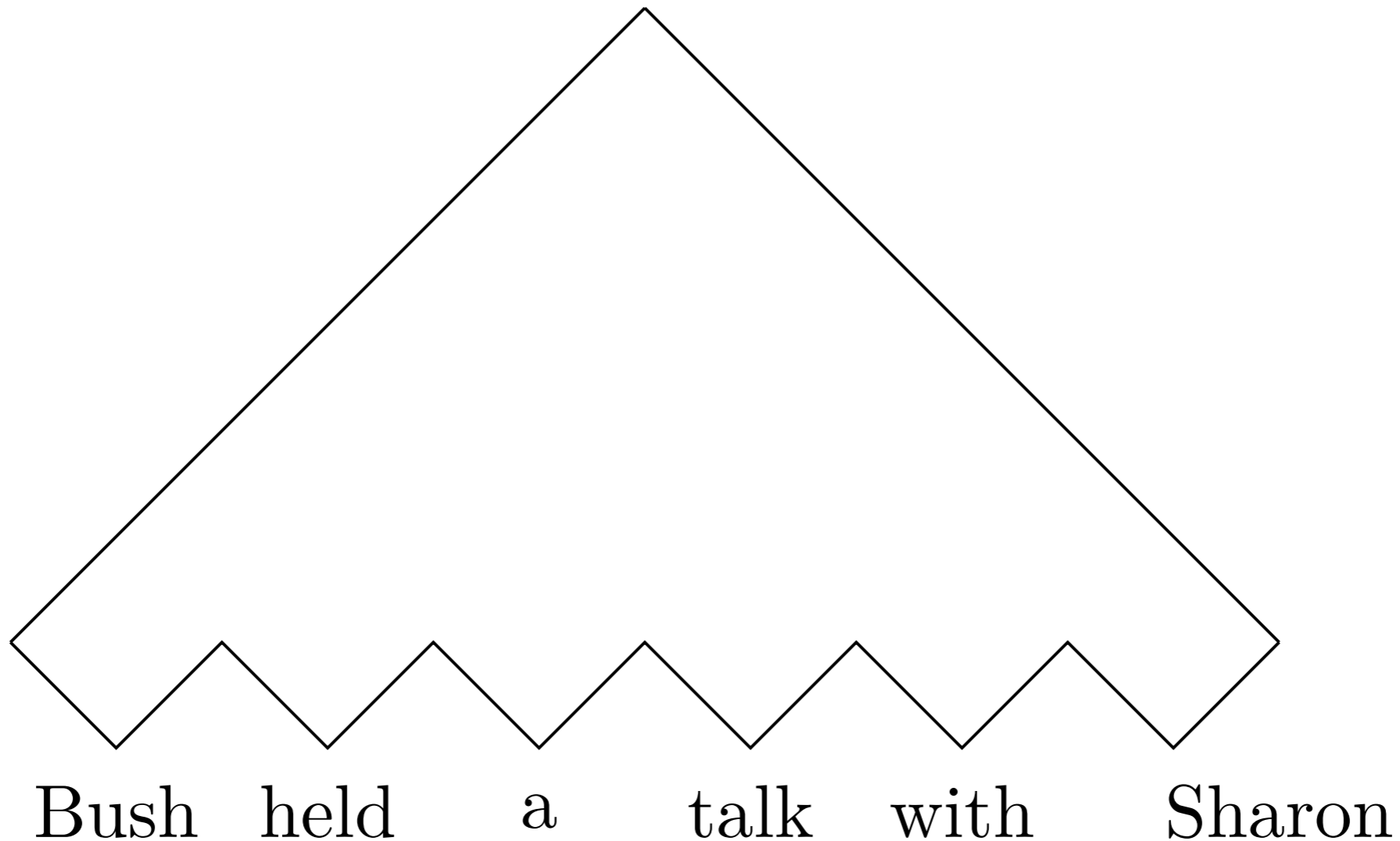


- parsing = intersection of CFG with a string (regular grammar)

# Parsing: CKY

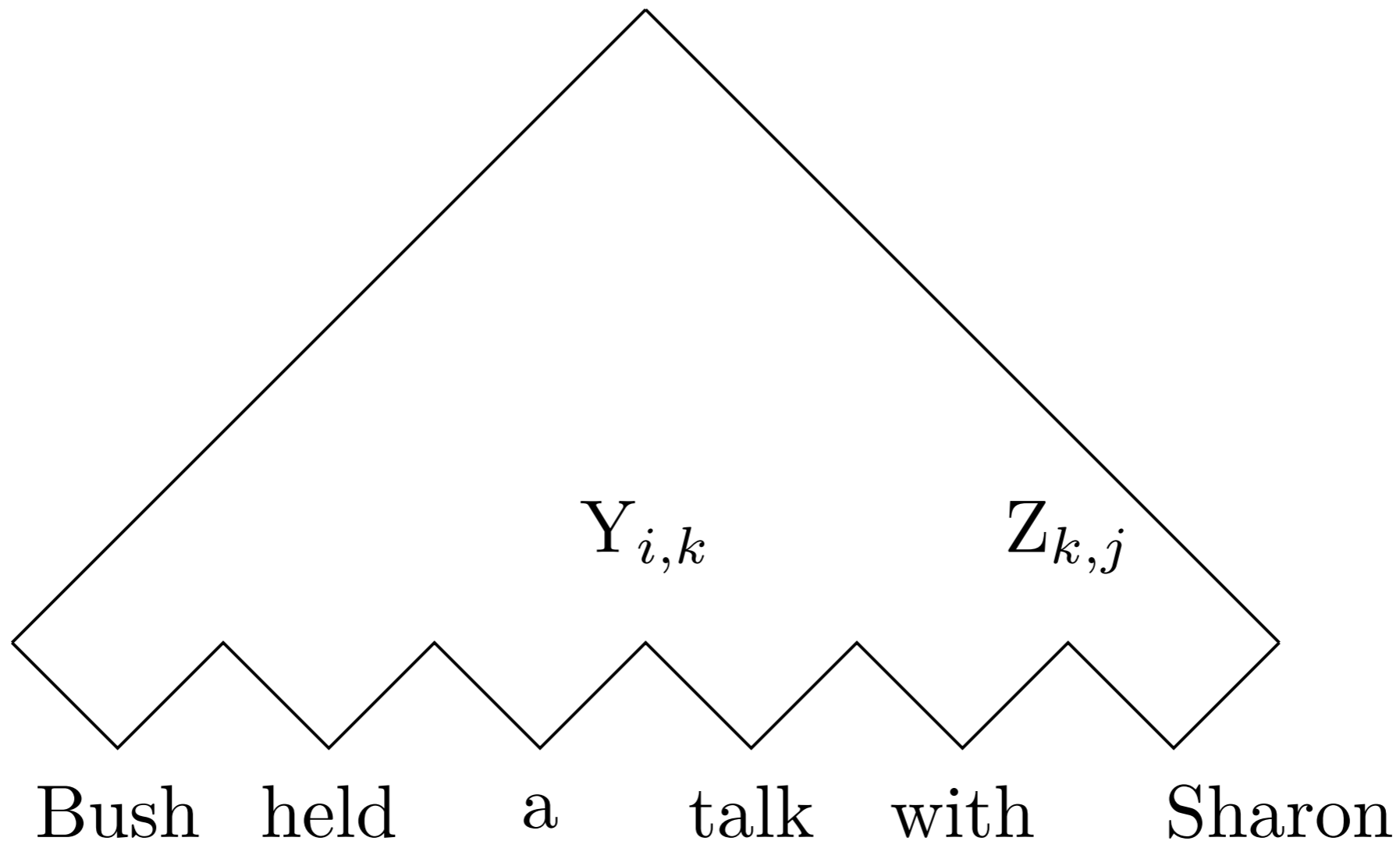
- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

# Parsing: CKY



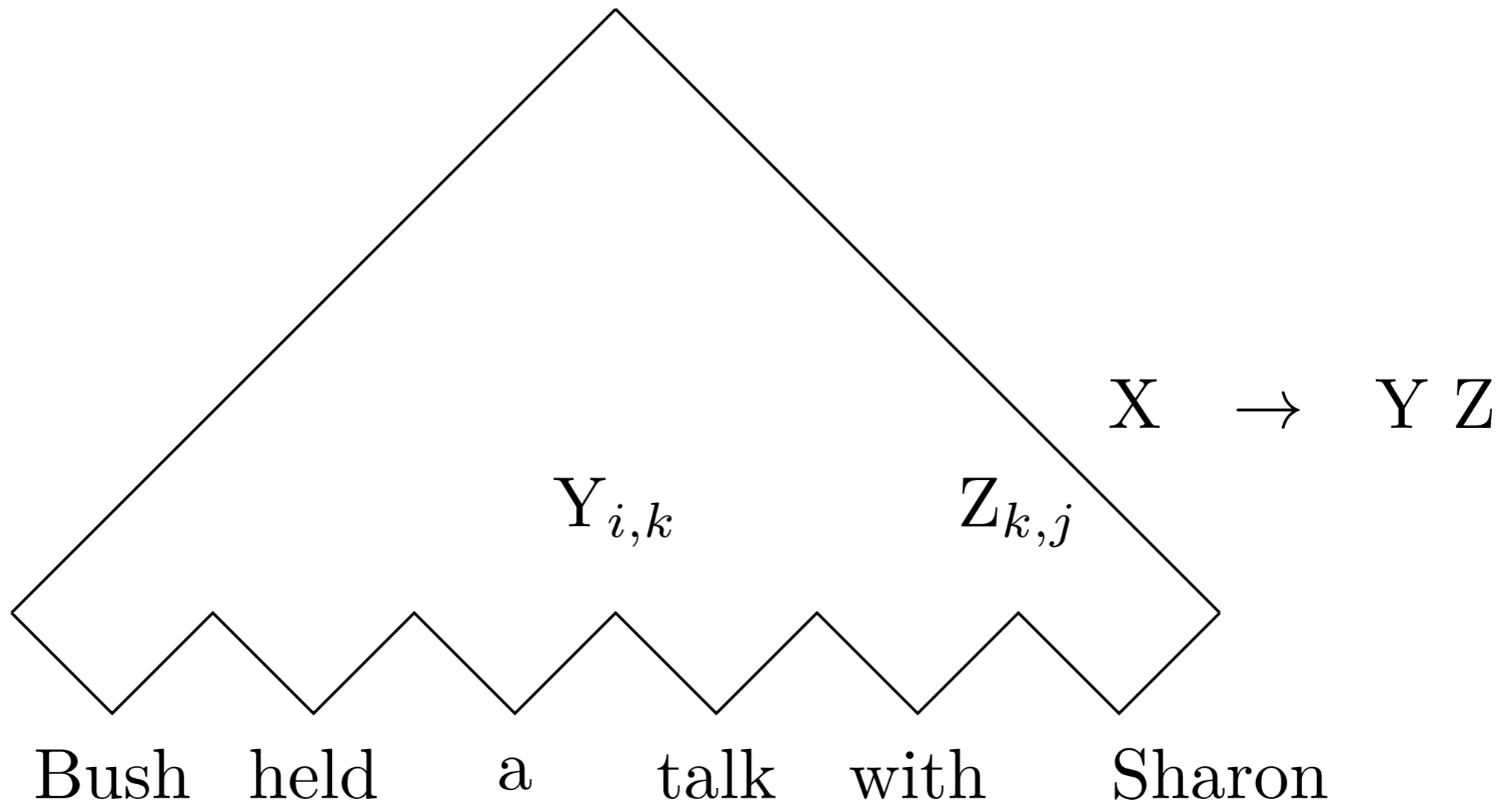
- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

# Parsing: CKY



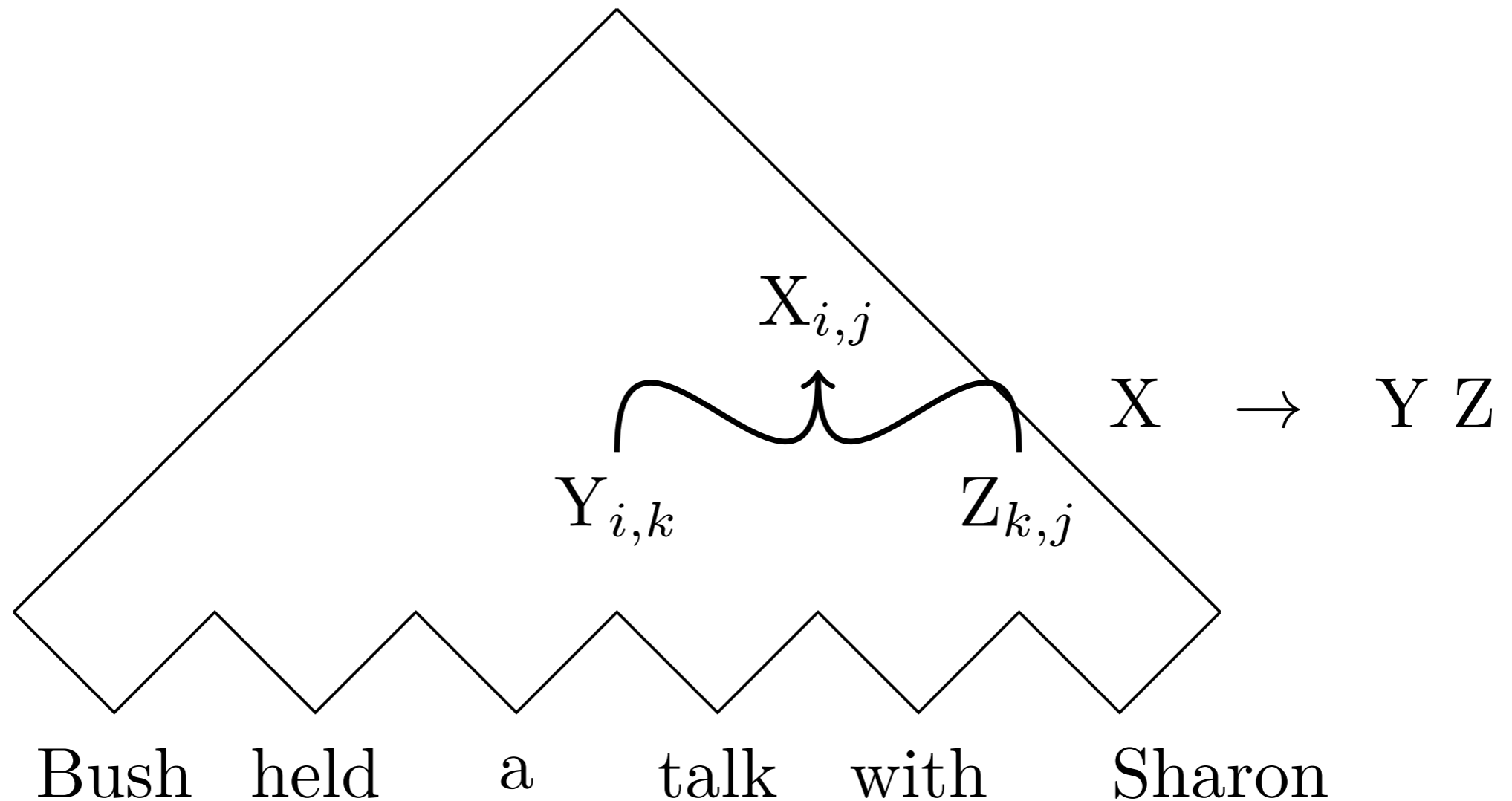
- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

# Parsing: CKY



- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

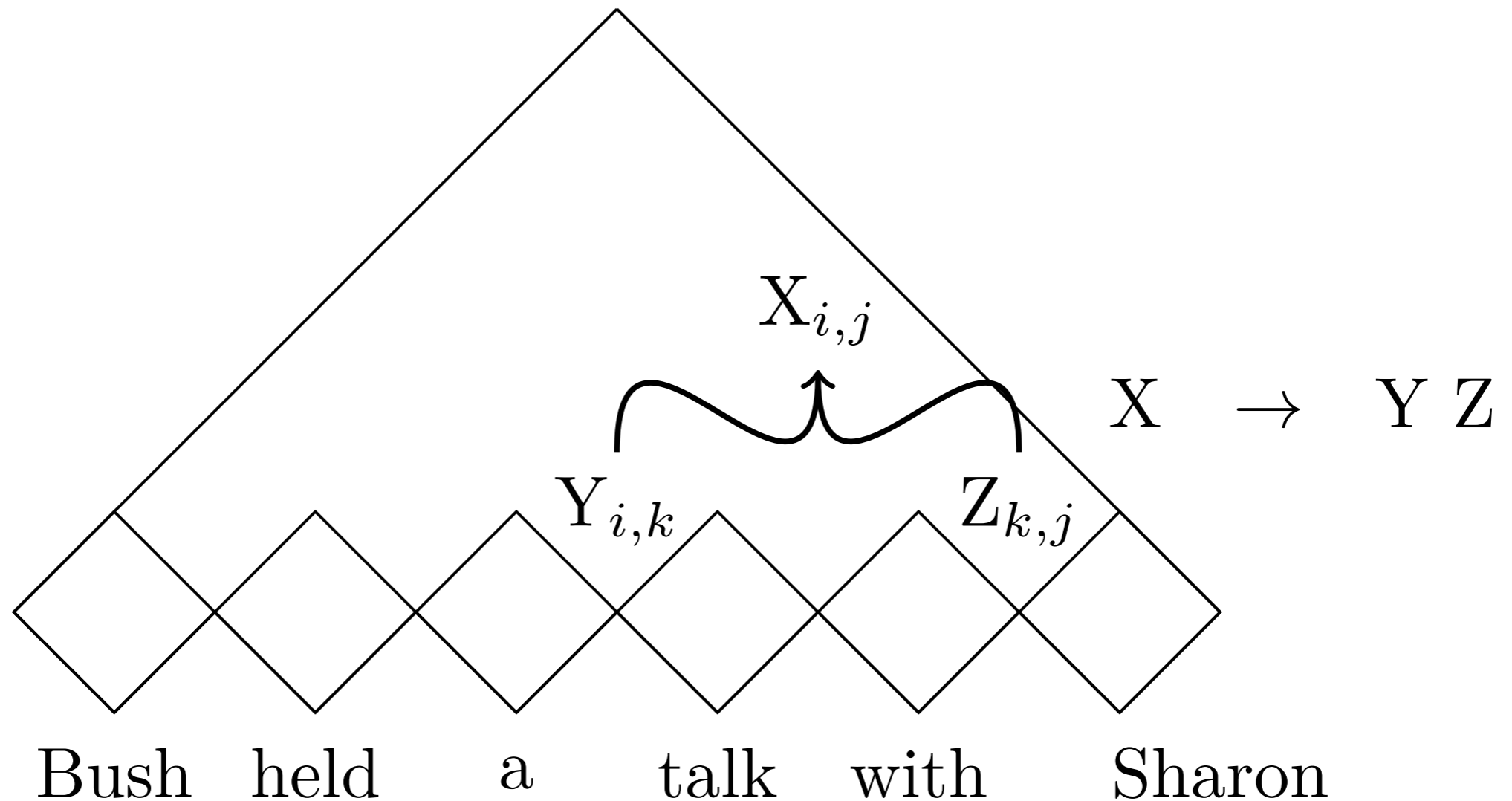
# Parsing: CKY



- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

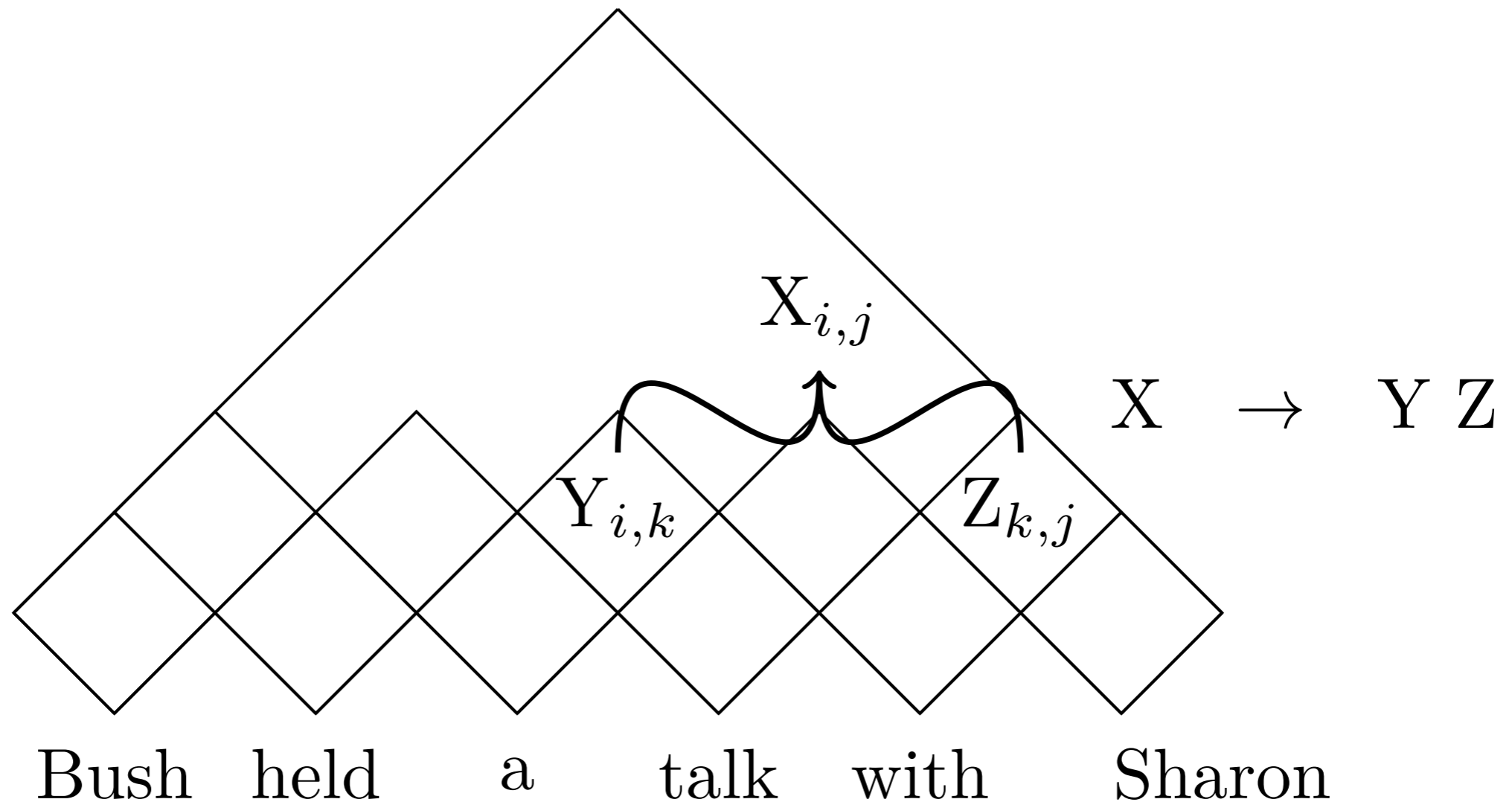


# Parsing: CKY



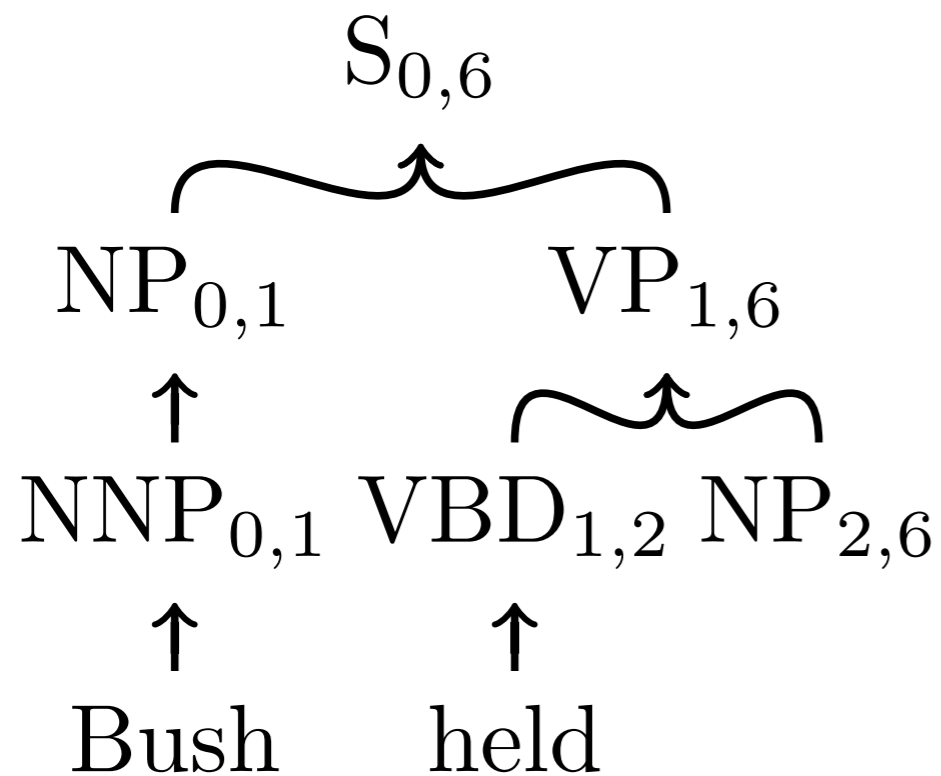
- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

# Parsing: CKY



- $O(n^3)$  : For each length  $n$ , for each position  $i$ , for each rule  $X \rightarrow Y Z$ , for each split point  $k$
- (Bottom-up) topological order

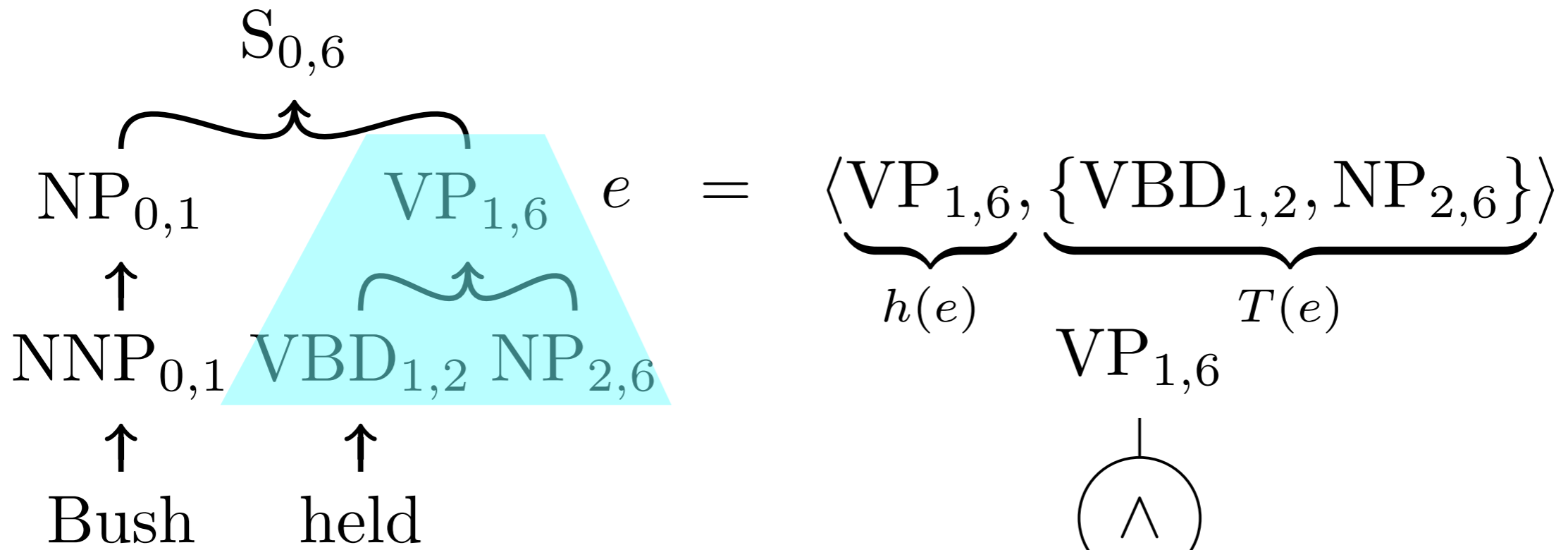
# Hypergraph



(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

# Hypergraph



(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

(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

# Deductive System

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

(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

# Deductive System

$$\begin{array}{c}
 \text{VBD}_{1,2} \quad \text{NP}_{2,6} \\
 \underbrace{\hspace{10em}} \\
 \text{VP}_{1,6}
 \end{array}
 \quad
 \frac{\overbrace{\text{VBD}_{1,2} \quad \text{NP}_{2,6}}^{\text{antecedents}}}{\underbrace{\text{VP}_{1,6}}_{\text{consequent}}} \text{VP}_{[i,j]} \rightarrow \text{VBZ}_{[j,k]} \text{NP}_{[i,k]}$$

(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

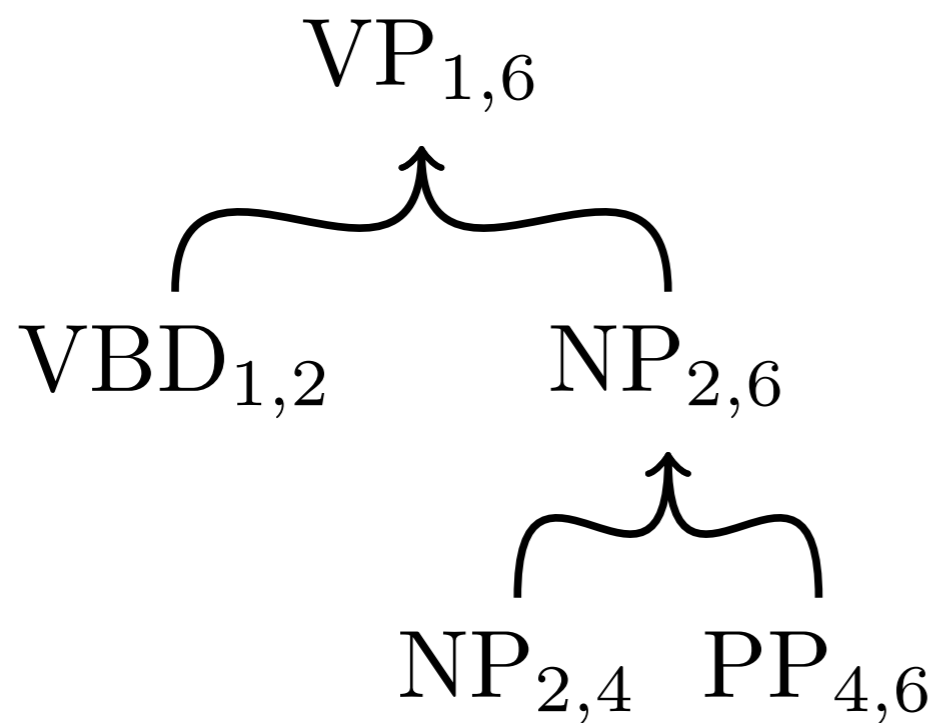
# Packed Forest

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

- A polynomial space encoding of exponentially many parses by sharing common sub-derivations
- Single derivation = tree



# Packed Forest

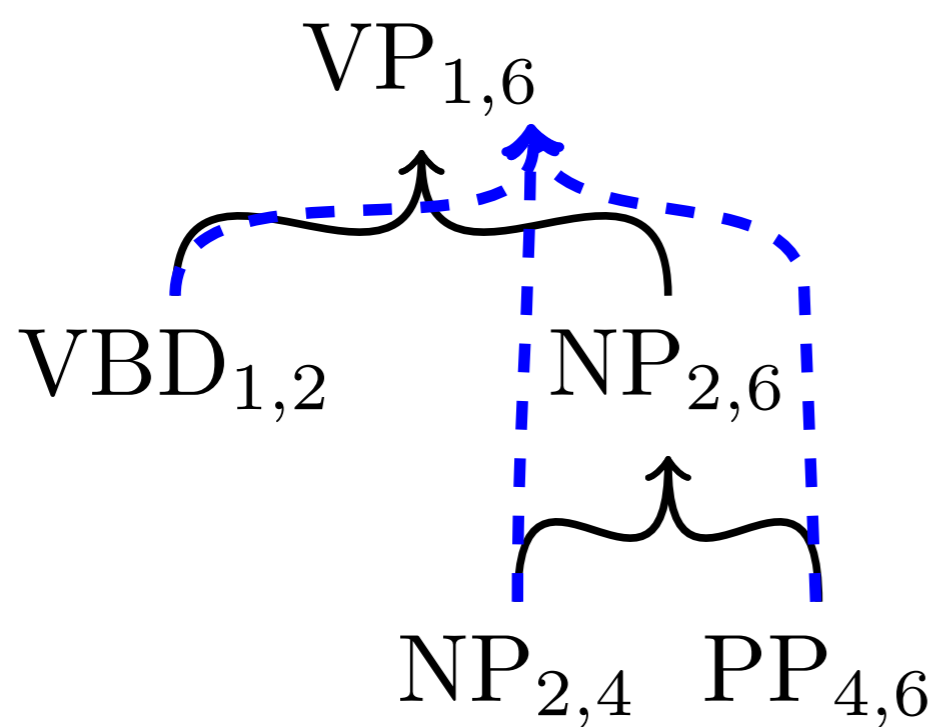


$$\frac{\text{VBD}_{1,2} \frac{\text{NP}_{2,4} \text{PP}_{4,6}}{\text{NP}_{2,6}}}{\text{VP}_{1,6}}$$

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

- A polynomial space encoding of exponentially many parses by sharing common sub-derivations
- Single derivation = tree

# Packed Forest



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

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

(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

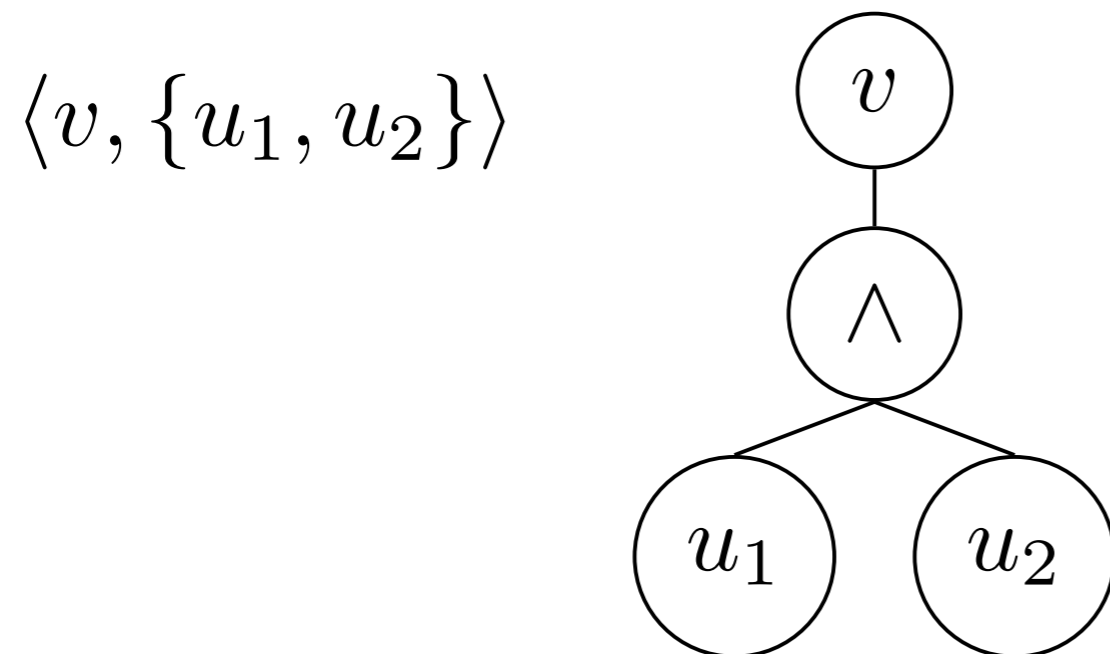
# Summary of Formalisms

hypergraph
vertex source-vertex target-vertex
hyperedge

$\langle v, \{u_1, u_2\} \rangle$

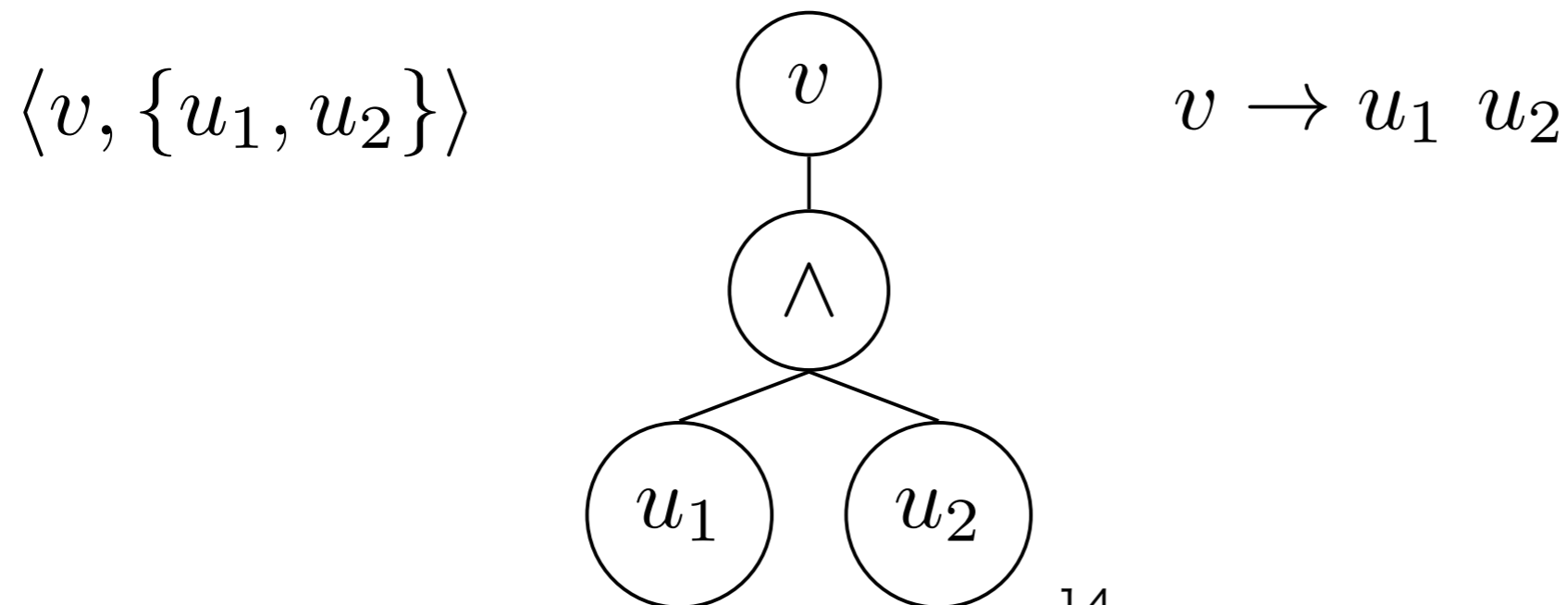
# Summary of Formalisms

hypergraph	AND/OR graph
vertex	OR-node
source-vertex	leaf OR-node
target-vertex	root OR-node
hyperedge	AND-node



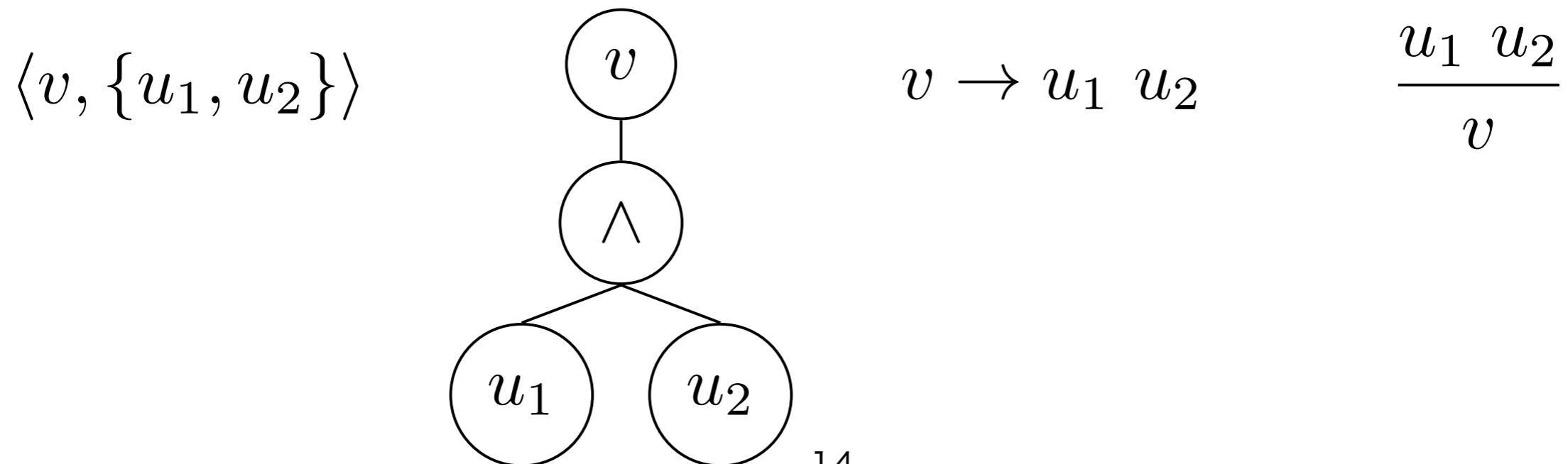
# Summary of Formalisms

hypergraph	AND/OR graph	CFG
vertex	OR-node	symbol
source-vertex	leaf OR-node	terminal
target-vertex	root OR-node	start symbol
hyperedge	AND-node	production



# 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



# Weights and Semirings

- Associate weights as in WFST
- $\otimes$  : extension (multiplicative),  $\oplus$  : summary (additive)



# Weights and Semirings

VP  $\xrightarrow{w_1}$  VBD NP  
NP  $\xrightarrow{w_2}$  NP PP

- Associate weights as in WFST
- $\otimes$  : extension (multiplicative),  $\oplus$  : summary (additive)

# Weights and Semirings

VP  $\xrightarrow{w_1}$  VBD NP

NP  $\xrightarrow{w_2}$  NP PP

$VP_{1,6} : w_1 \otimes c \otimes d$

$\underbrace{\hspace{10em}}_{\substack{\text{VBD}_{1,2} : c \quad \text{NP}_{2,6} : d}}$

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

- Associate weights as in WFST
- $\otimes$  : extension (multiplicative),  $\oplus$  : summary (additive)

# Weights and Semirings

VP  $\xrightarrow{w_1}$  VBD NP

NP  $\xrightarrow{w_2}$  NP PP

VP<sub>1,6</sub> :  $w_1 \otimes c \otimes d$

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

VBD<sub>1,2</sub> :  $c$  NP<sub>2,6</sub> :  $d$

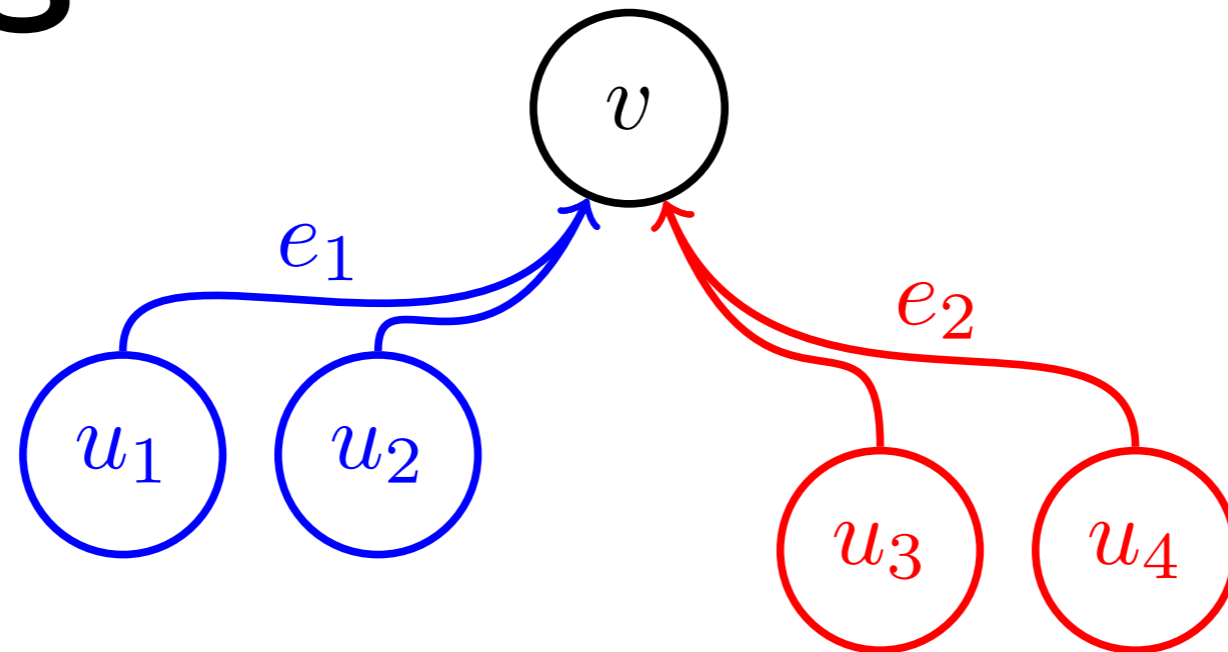
NP<sub>2,6</sub> :  $w_2 \otimes a \otimes b$

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

NP<sub>2,4</sub> :  $a$  PP<sub>4,6</sub> :  $b$

- Associate weights as in WFST
- $\otimes$  : extension (multiplicative),  $\oplus$  : summary (additive)

# Weights and Semirings



$$d(v) = (w(e_1, u_1, u_2) \otimes d(u_1) \otimes d(u_2)) \oplus (w(e_2, u_3, u_4) \otimes d(u_3) \otimes d(u_4))$$

- 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	$\times$	0	1
Real	R	+	$\times$	0	1
Log	R	logsumexp	+	$+\infty$	0
Tropical	R	min	+	$+\infty$	0

# Conclusion

- Review important concepts from “parsing”
- CFG, parsing, hypergraph, deductive system, weights, semirings

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- **Tree-based SMT**
  - **Synchronous-CFG**
  - String-to-Tree, Tree-to-String
  - Bitext parsing

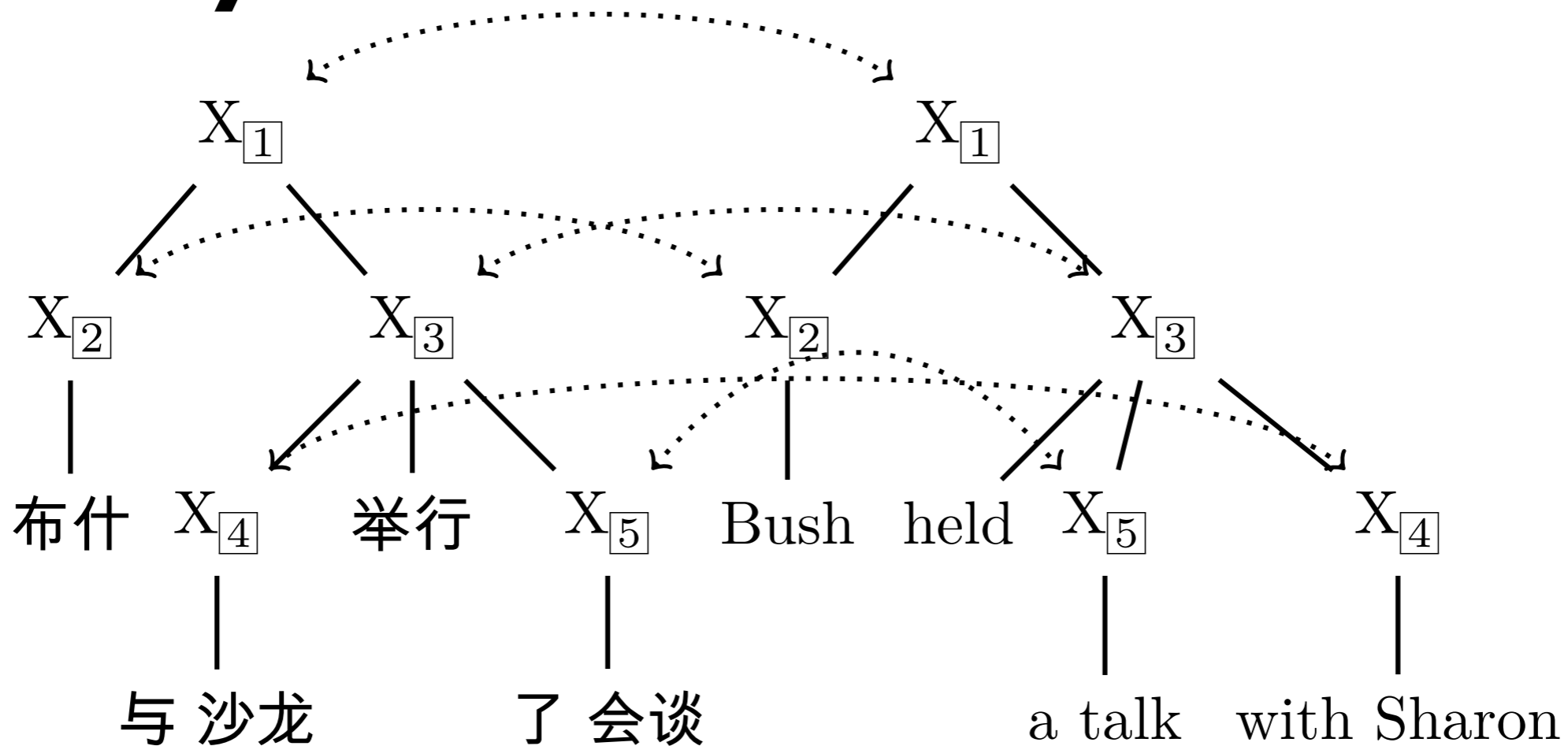


# Synchronous-CFG

$$\begin{aligned}\hat{\mathbf{e}} &= \operatorname{argmax}_{\mathbf{e}} \frac{\exp(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, D, \mathbf{f}))}{\sum_{\mathbf{e}', D'} \exp(\mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}', D', \mathbf{f}))} && \text{(Chiang, 2007)} \\ &= \operatorname{argmax}_{\mathbf{e}} \mathbf{w}^\top \cdot \mathbf{h}(\mathbf{e}, D, \mathbf{f})\end{aligned}$$

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

# Synchronous-CFG



$$\hat{e} = \operatorname{argmax}_{e} \frac{\exp(\mathbf{w}^{\top} \cdot \mathbf{h}(e, D, \mathbf{f}))}{\sum_{e', D'} \exp(\mathbf{w}^{\top} \cdot \mathbf{h}(e', D', \mathbf{f}))} \quad (\text{Chiang, 2007})$$

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

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

# Synchronous-CFG: Model

- We use two categories, S and X (Chiang, 2007)
- Or, borrow linguistic categories from syntactic parse (Zollman and Venugopal, 2006)

# Synchronous-CFG: Model

$$S \rightarrow \langle S_{[1]} X_{[2]}, S_{[1]} X_{[2]} \rangle$$
$$S \rightarrow \langle X_{[1]}, X_{[1]} \rangle$$
$$X \rightarrow \langle X_{[1]} \text{ 举行 } X_{[2]}, \text{hold } X_{[2]} X_{[1]} \rangle$$
$$X \rightarrow \langle \text{与 沙龙}, \text{with Sharon} \rangle$$

- We use two categories, S and X (Chiang, 2007)
- Or, borrow linguistic categories from syntactic parse (Zollman and Venugopal, 2006)

# Synchronous-CFG: Model

$$S \rightarrow \langle S_{[1]} X_{[2]}, S_{[1]} X_{[2]} \rangle$$
$$S \rightarrow \langle X_{[1]}, X_{[1]} \rangle$$
$$X \rightarrow \langle X_{[1]} \text{举行} X_{[2]}, \text{hold } X_{[2]} X_{[1]} \rangle$$
$$X \rightarrow \langle \text{与沙龙}, \text{with Sharon} \rangle$$
$$VP \rightarrow \langle VBD_{[1]} NP_{[2]}, NP_{[2]} VBD_{[1]} \rangle$$
$$NP \rightarrow \langle NP_{[1]} PP_{[2]}, NP_{[1]} PP_{[2]} \rangle$$
$$VP \rightarrow \langle VBD_{[1]} NP_{[2]} PP_{[3]}, NP_{[2]} PP_{[3]} VBD_{[1]} \rangle$$

- 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

Bush	■				
held			■		
a					
talk					■
with		■			
Sharon			■		

(Example from Huang and Chiang, 2007)

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

# Rule Extraction

布什 与 沙龙举行了 会谈

Bush	■				
held			■		
a					
talk					■
with	■				
Sharon		■			

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

(Example from Huang and Chiang, 2007)

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

# Rule Extraction

布什 与 沙龙举行了 会谈

Bush	■				
held			■		
a					
talk					■
with	■				
Sharon		■			

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

⟨with Sharon, 与 沙龙⟩

(Example from Huang and Chiang, 2007)

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



# Rule Extraction

布什 与 沙龙举行了 会谈

Bush	■				
held			■		
a					
talk					■
with	■				
Sharon		■			

⟨held a talk with Sharon,  
与沙龙举行了会谈⟩

⟨with Sharon, 与沙龙⟩

⟨held, 举行⟩

(Example from Huang and Chiang, 2007)

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

# Rule Extraction

布什 与 沙龙举行了 会谈

Bush	■				
held			■		
a					
talk					■
with	■	■			
Sharon	■	■	X	→	⟨X <sub>1</sub> X <sub>2</sub> 了 会谈, X <sub>2</sub> a talk X <sub>1</sub> ⟩

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

⟨with Sharon, 与 沙龙⟩

⟨held, 举行⟩

X<sub>2</sub> 了 会谈, X<sub>2</sub> a talk X<sub>1</sub>

(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

Bush	■				
held			■		
a					
talk					■
with		■			
Sharon			■		

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

# 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⟩

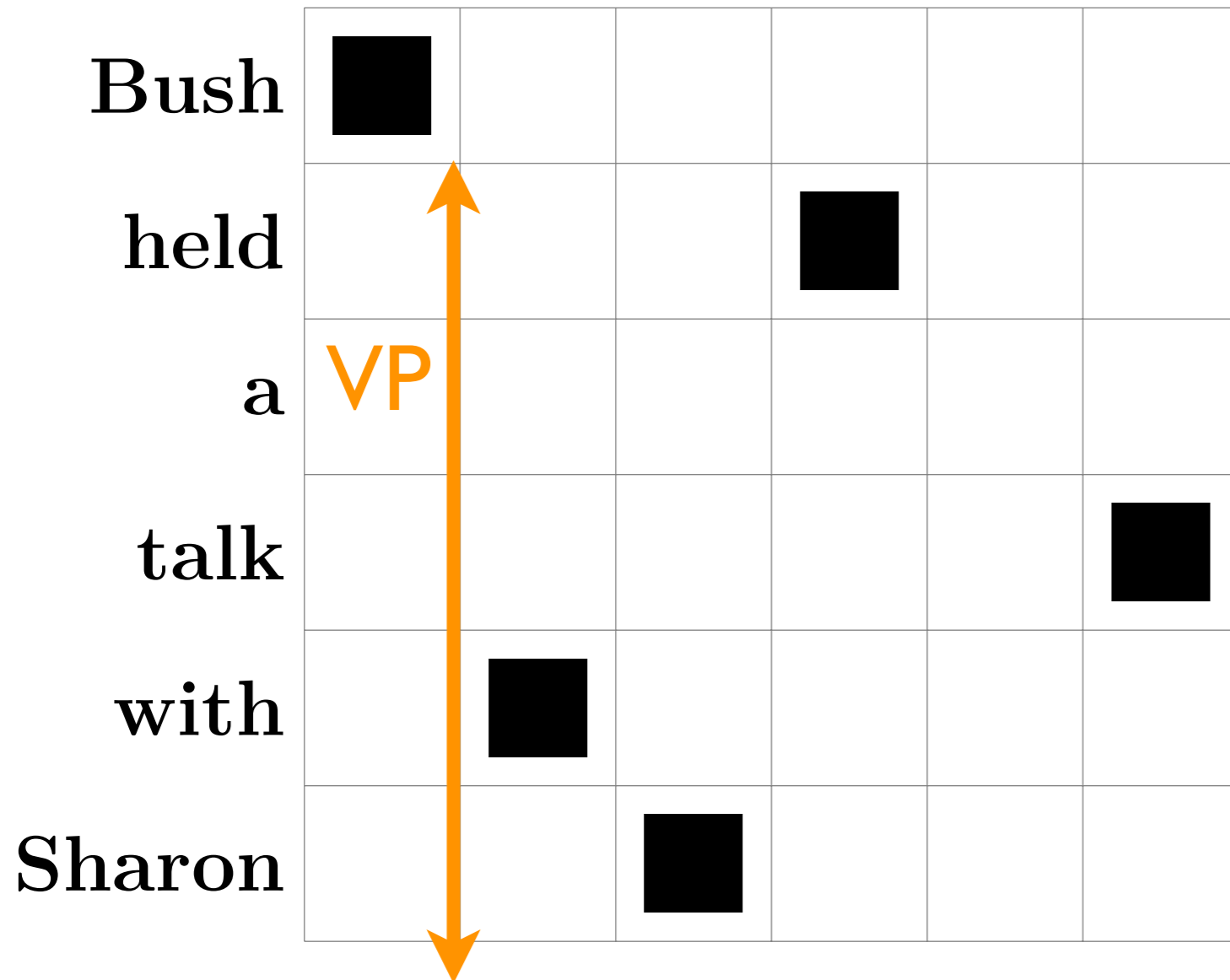
Bush	■				
held			■		
a					
talk					■
with		■			
Sharon			■		

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

# Syntactic Categories

布什 与 沙龙举行了会谈

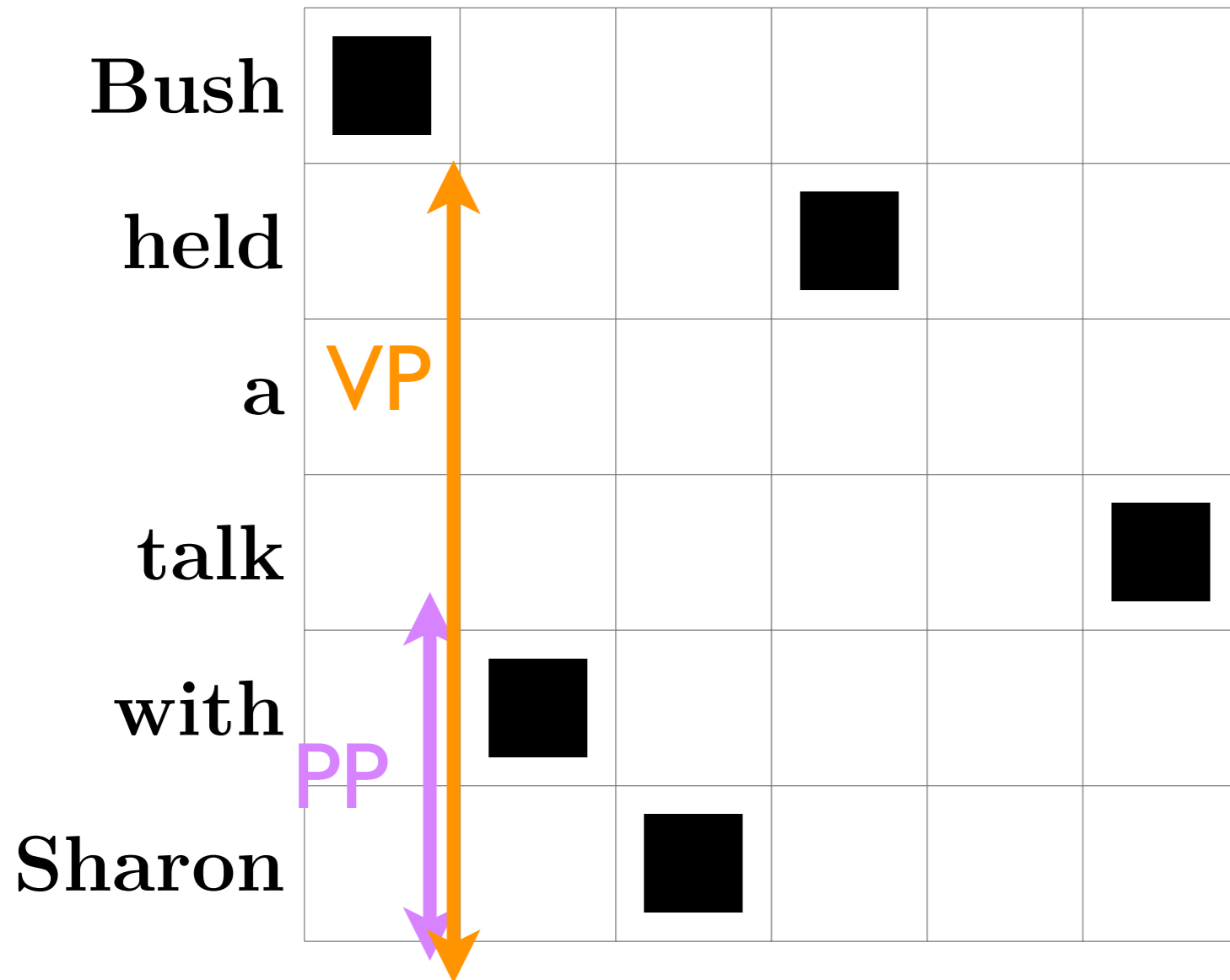
⟨held a talk with Sharon,  
与沙龙举行了会谈⟩



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

# Syntactic Categories

布什 与 沙龙举行了 会谈



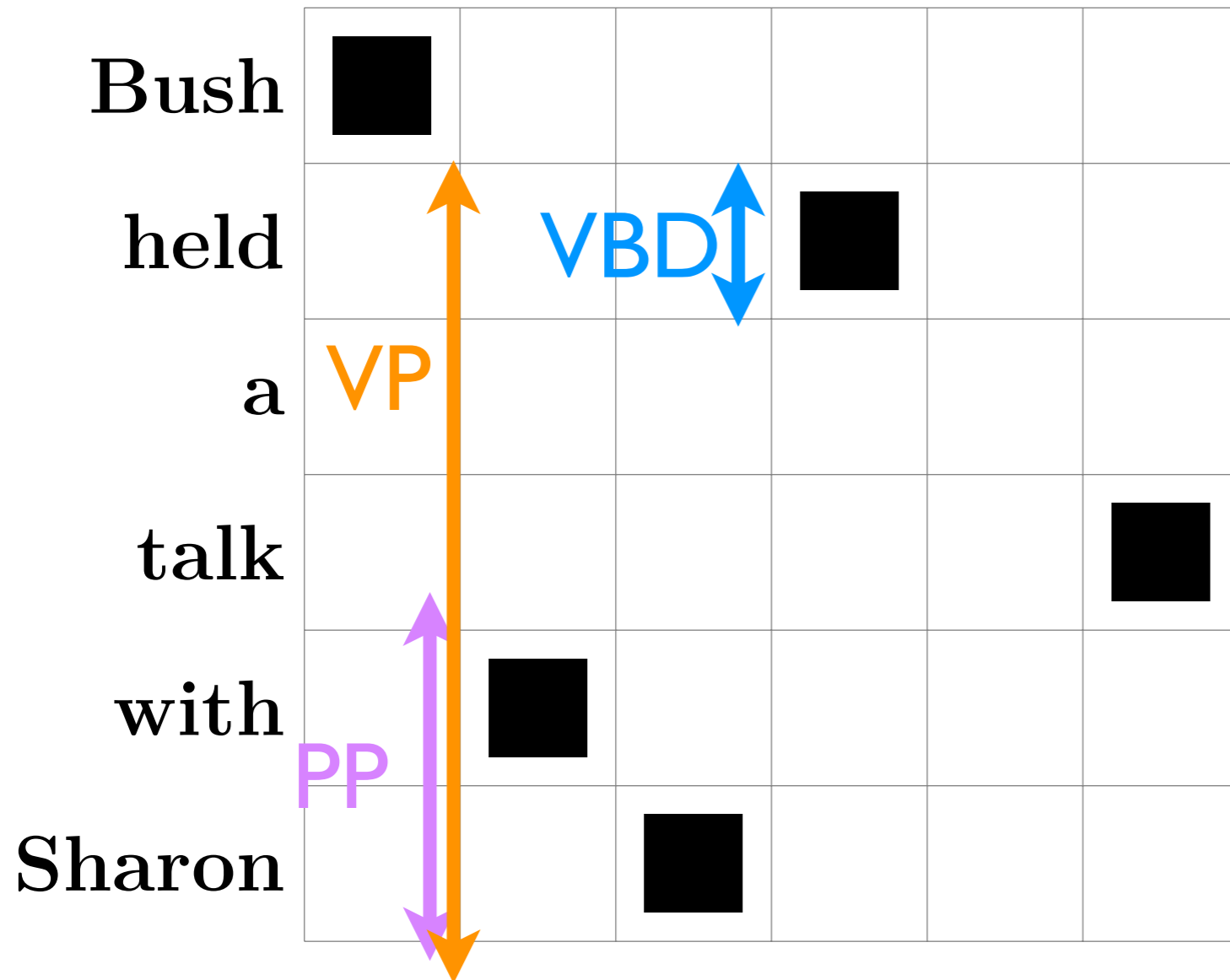
⟨held a talk with Sharon,  
与 沙龙 举行了 会谈⟩

⟨with Sharon, 与 沙龙⟩

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

# 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⟩

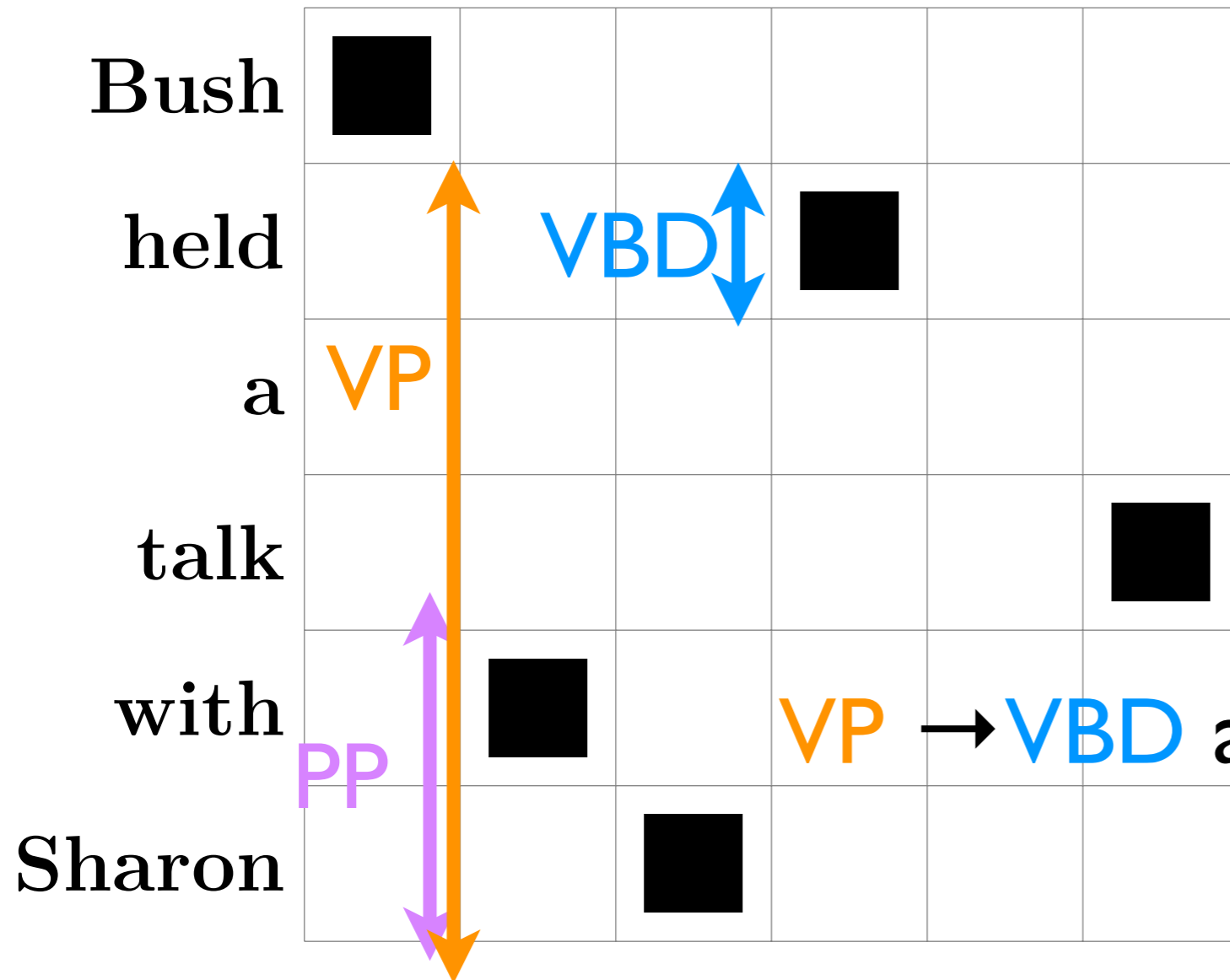
⟨with Sharon, 与沙龙⟩

⟨held, 举行⟩

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

# 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〉

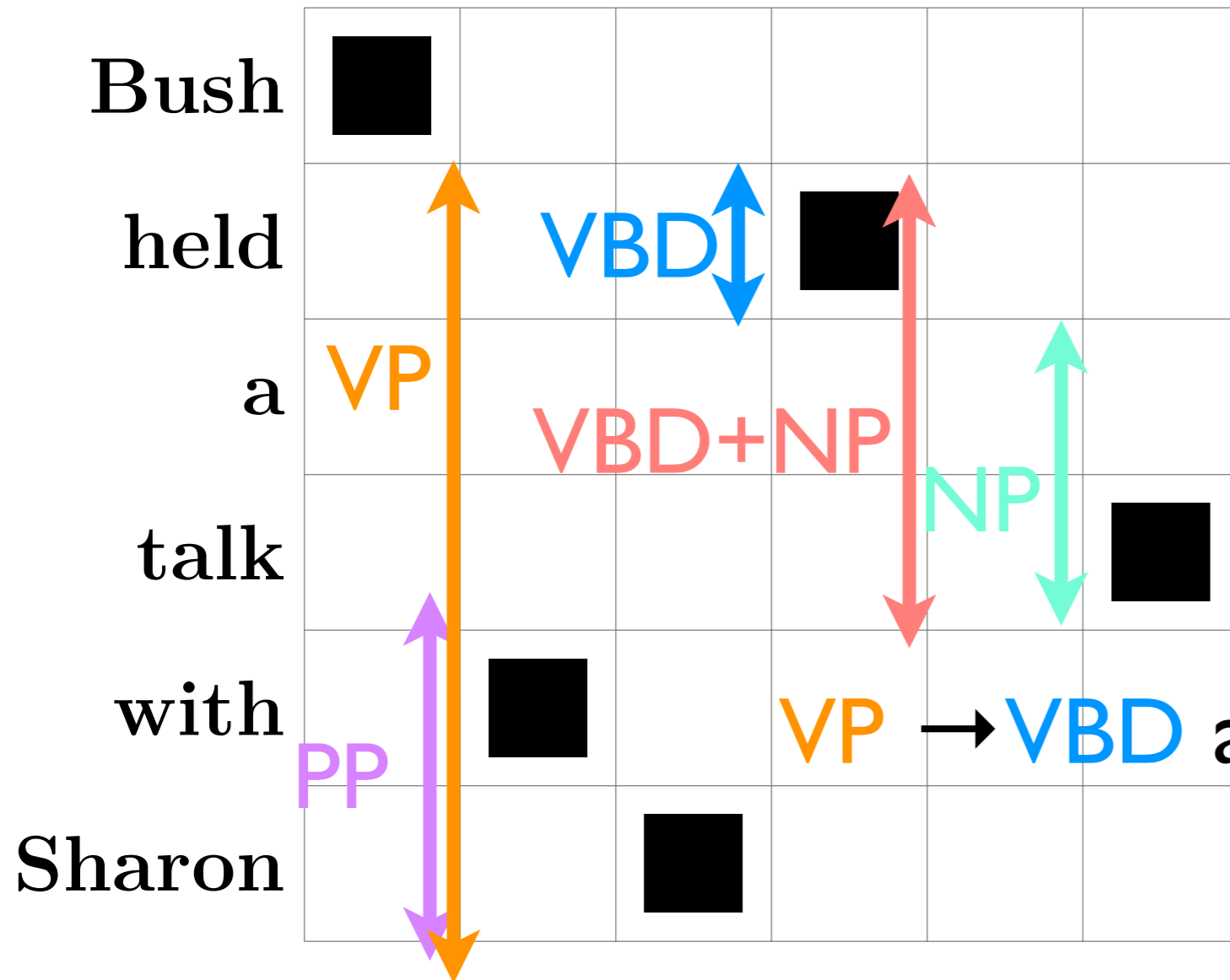
a talk PP, PP VBD 了 会谈

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



# 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〉

with a talk PP, PP VBD 了 会谈

- 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

Bush	■				
held			■		
a					
talk					■
with	■				
Sharon		■			

# Exhaustive Extraction

布什 与 沙龙举行了会谈

Bush	■				
held			■		
a					
talk					■
with	■				
Sharon		■			

- Exhaustively extract rules as in phrase-based MT

# Exhaustive Extraction

布什 与 沙龙举行了会谈

					$X_1$	$X_2$	了 会谈	$X_2$	a talk	$X_1$
Bush	■				$X_1$	$X_2$	会谈	$X_2$	a talk	$X_1$
held			■		$X_1$	$X_2$	会谈	$X_2$	talk	$X_1$
a					$X_1$		举行 $X_2$	held	$X_2$	$X_1$
talk					$X_1$		举行了 $X_2$	held a	$X_2$	$X_1$
with		■					与 沙龙 $X_1$	$X_1$	with Sharon	
Sharon			■				与 $X_1$ $X_2$	$X_2$	with	$X_1$

- Exhaustively extract rules as in phrase-based MT

# Exhaustive Extraction

布什 与 沙龙举行了会谈

					$X_1$	$X_2$	了 会谈	$X_2$	a talk	$X_1$
Bush	■				$X_1$	$X_2$	会谈	$X_2$	a talk	$X_1$
held			■		$X_1$	$X_2$	会谈	$X_2$	talk	$X_1$
a					$X_1$		举行 $X_2$	held	$X_2$	$X_1$
talk					$X_1$		举行了 $X_2$	held a	$X_2$	$X_1$
with		■					与 沙龙 $X_1$	$X_1$	with Sharon	
Sharon			■				与 $X_1$ $X_2$	$X_2$	with $X_1$	
					S		$\rightarrow$	$\langle S_1$	$X_2, S_1$	$X_2 \rangle$
					S		$\rightarrow$	$\langle X_1,$	$X_1 \rangle$	

- 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  $(\alpha, \beta)$  from the data:
- $\alpha$  = source side string,  $\beta$  = 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

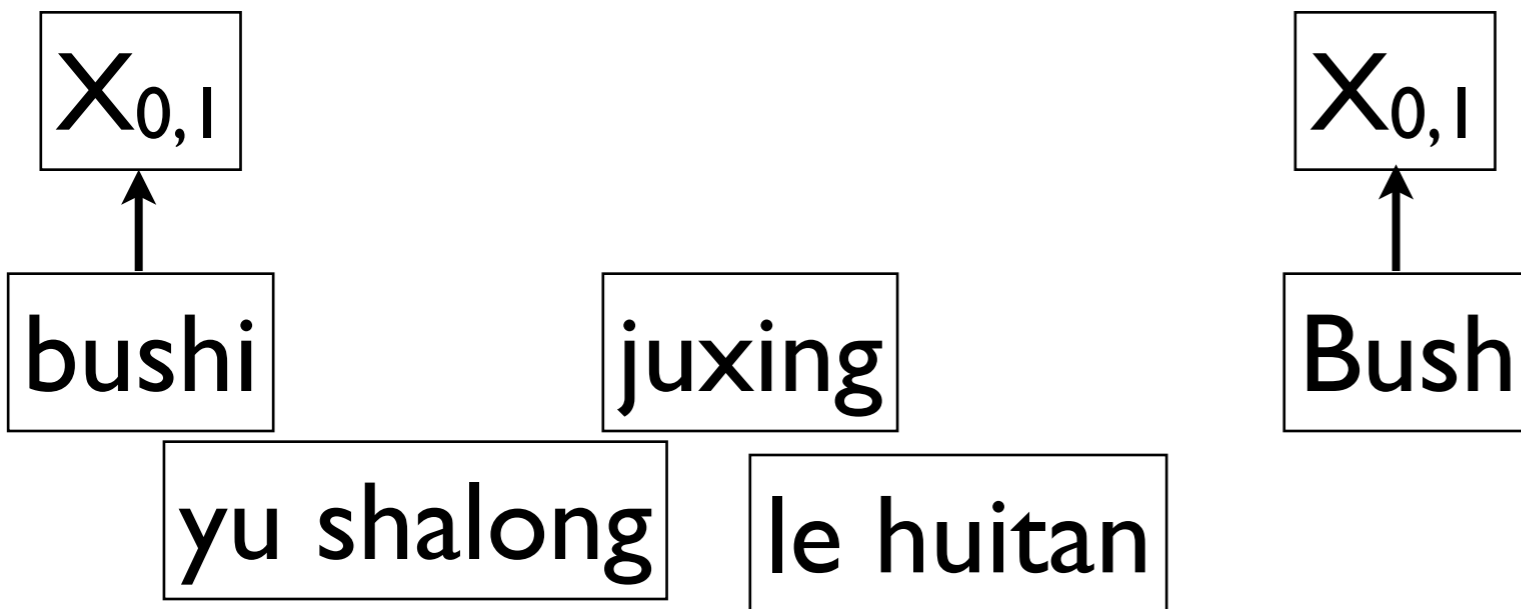


# Synchronous-CFG: Parsing

bushi                      juxing  
yu shalong              le huitan

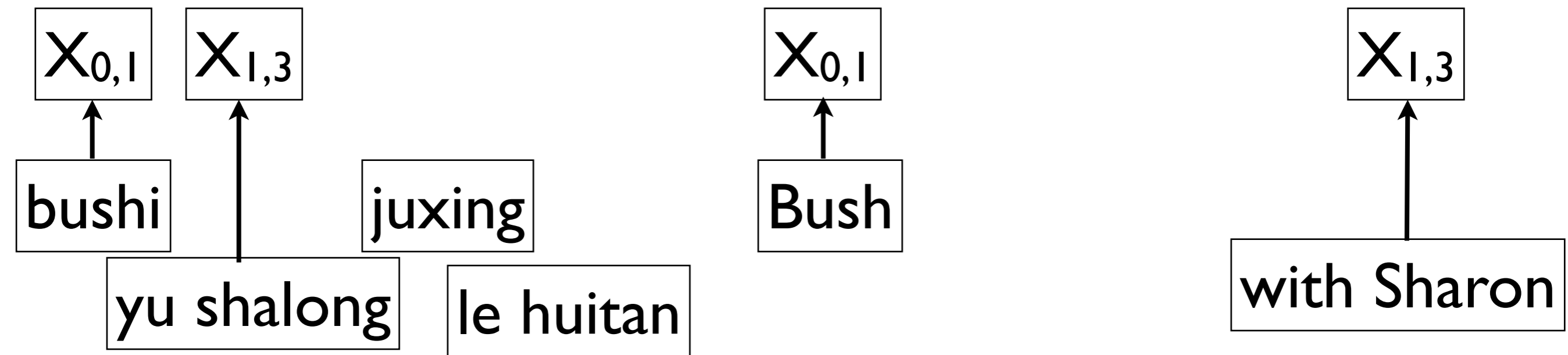
- Parse input sentence using the source side, and construct a translation forest by target side

# Synchronous-CFG: Parsing



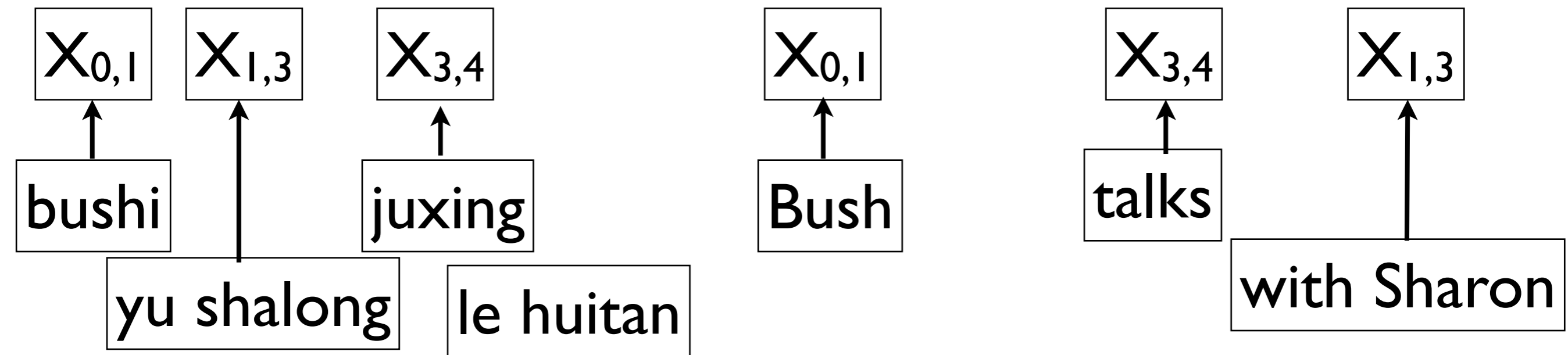
- Parse input sentence using the source side, and construct a translation forest by target side

# Synchronous-CFG: Parsing



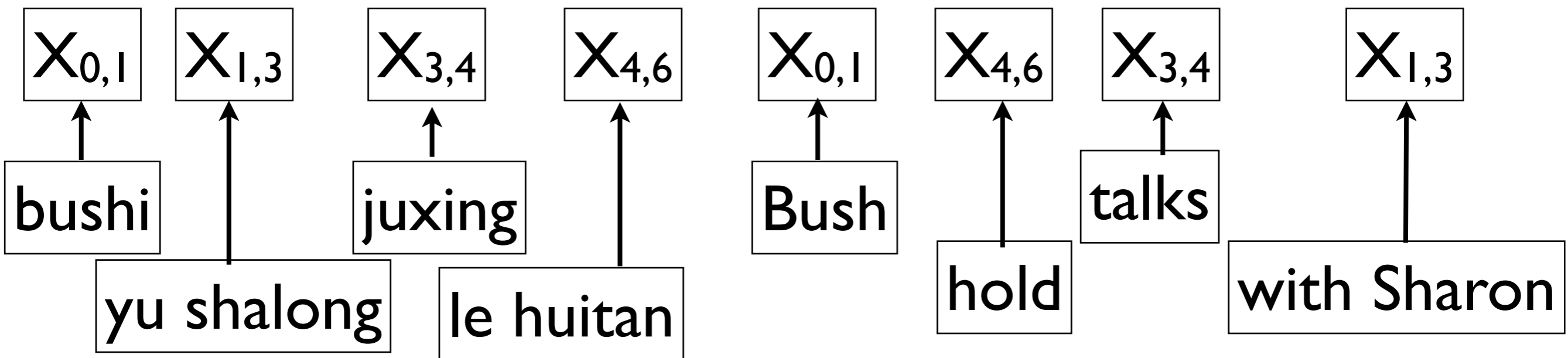
- Parse input sentence using the source side, and construct a translation forest by target side

# Synchronous-CFG: Parsing



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

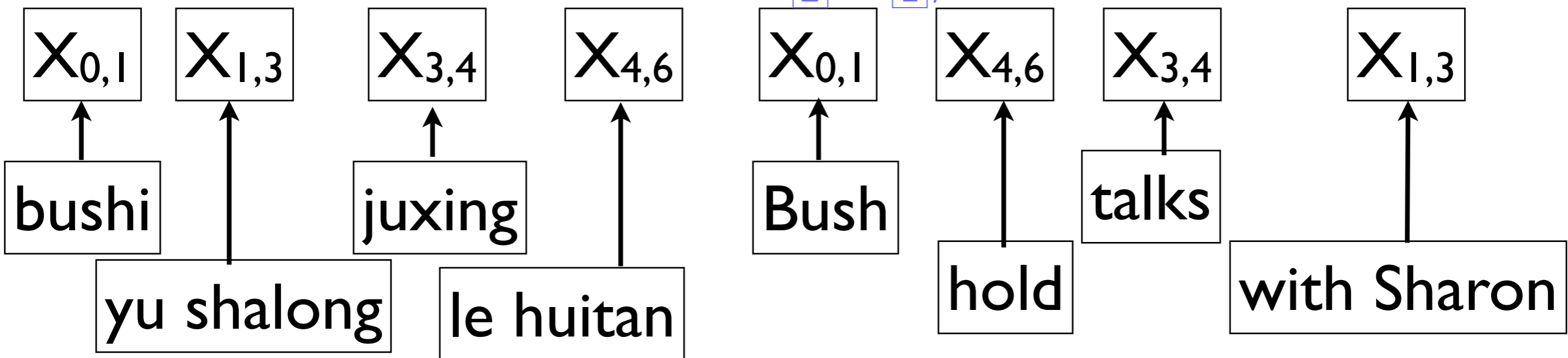
# Synchronous-CFG: Parsing



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

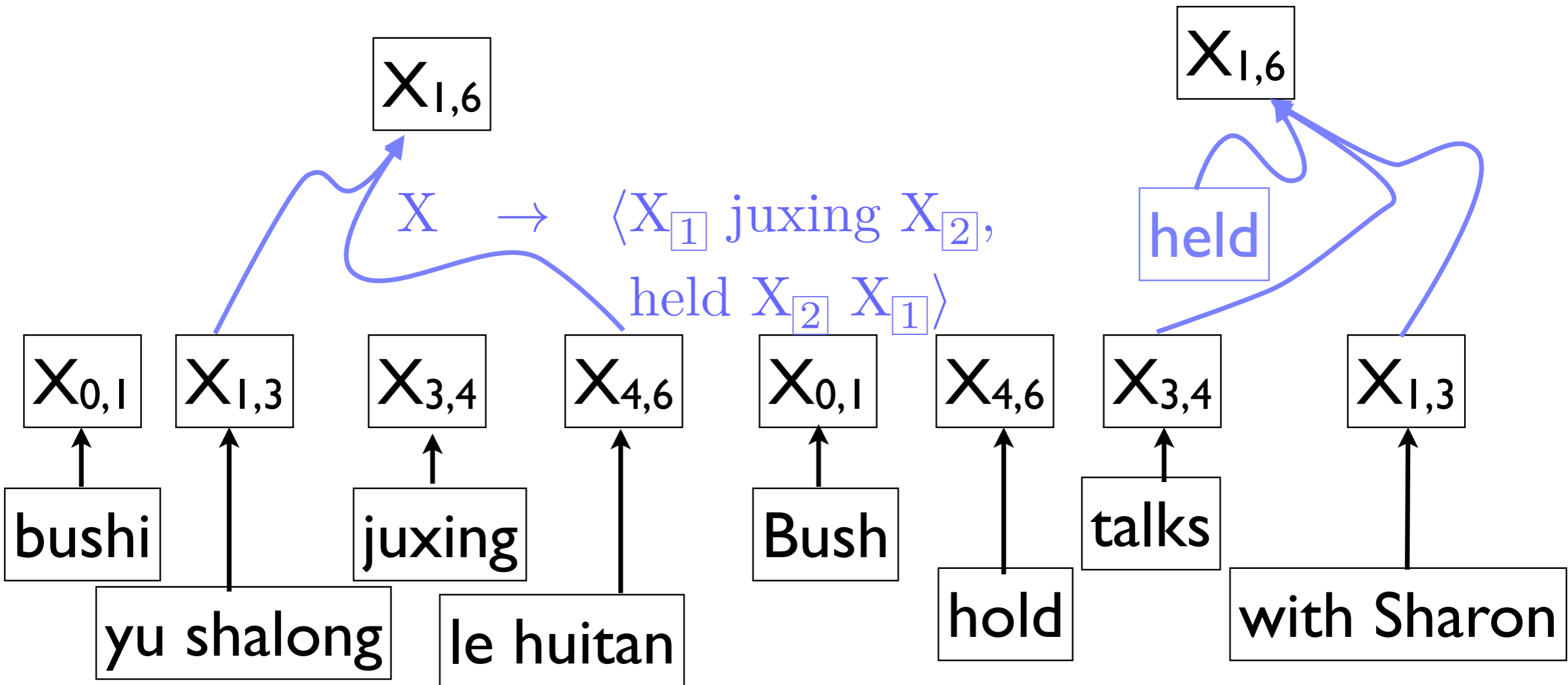
# Synchronous-CFG: Parsing

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



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

# Synchronous-CFG: Parsing



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

# Synchronous-CFG: Parsing

$X \rightarrow \langle X_1 X_2 \text{ le huitan},$   
 $X_2 \text{ a talk } X_1 \rangle$

$X_{1,6}$

$X_{1,6}$

a talk

held

$X \rightarrow \langle X_1 \text{ juxing } X_2,$   
 $\text{held } X_2 X_1 \rangle$

$X_{0,1}$

$X_{1,3}$

$X_{3,4}$

$X_{4,6}$

$X_{0,1}$

$X_{4,6}$

$X_{3,4}$

$X_{1,3}$

bushi

juxing

Bush

talks

yu shalong

le huitan

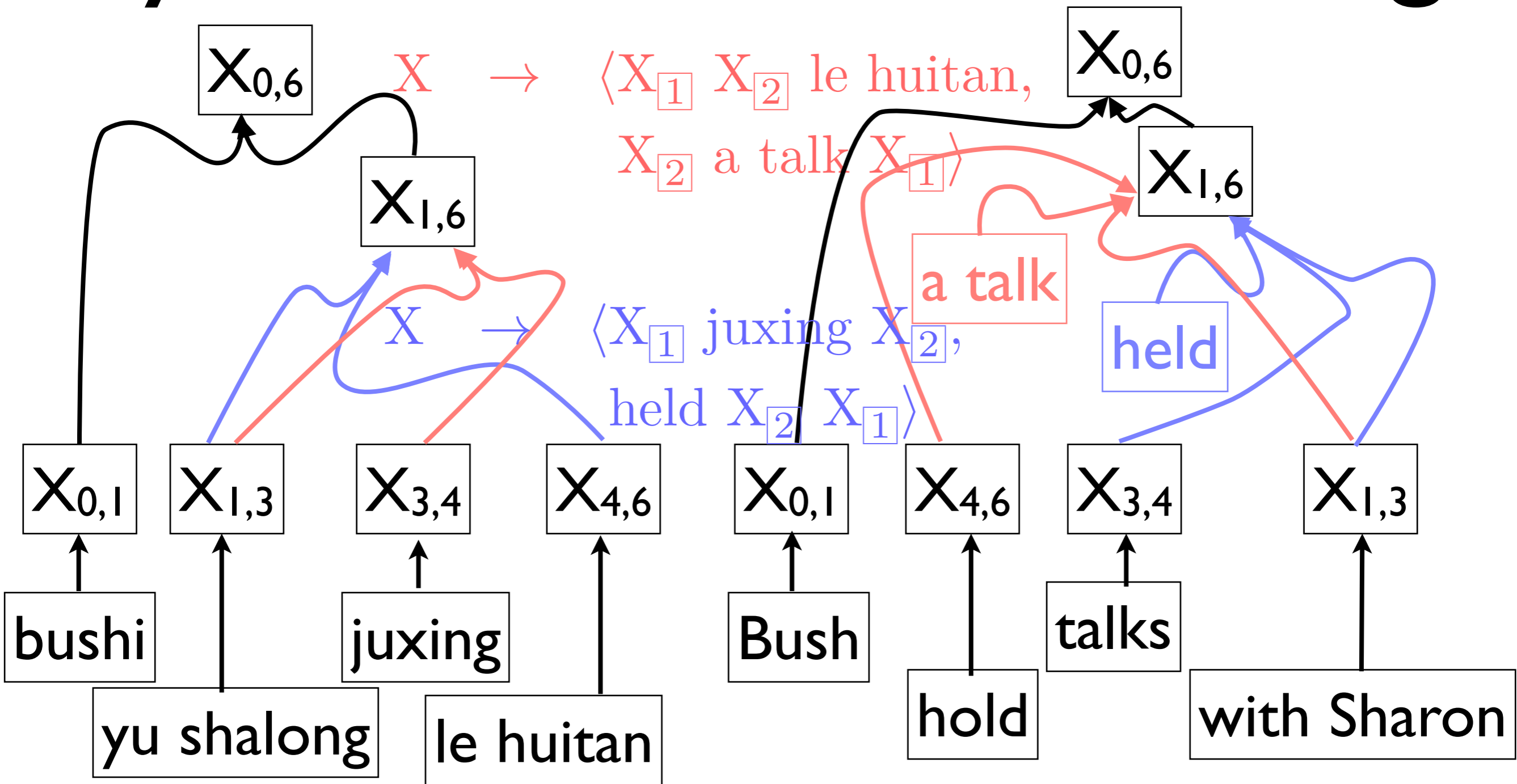
hold

with Sharon

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



# Synchronous-CFG: Parsing



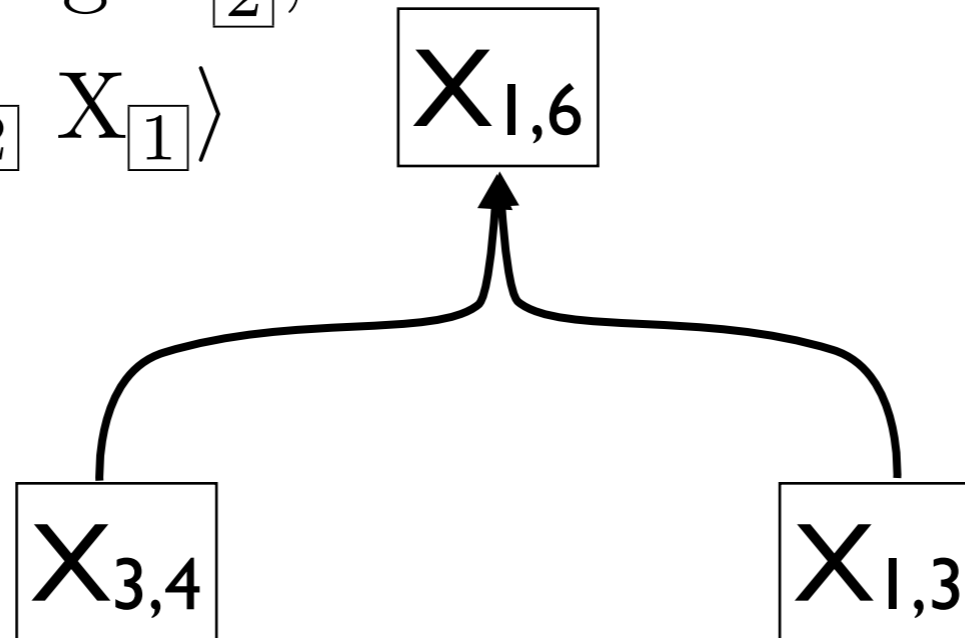
- Parse input sentence using the source side, and construct a translation forest by target side

# Synchronous-CFG: Parsing

- Translation by SCFG = monolingual parsing using the source side grammar
- Complexity:  $O(n^3)$  as in monolingual CKY

# Non-Local Features

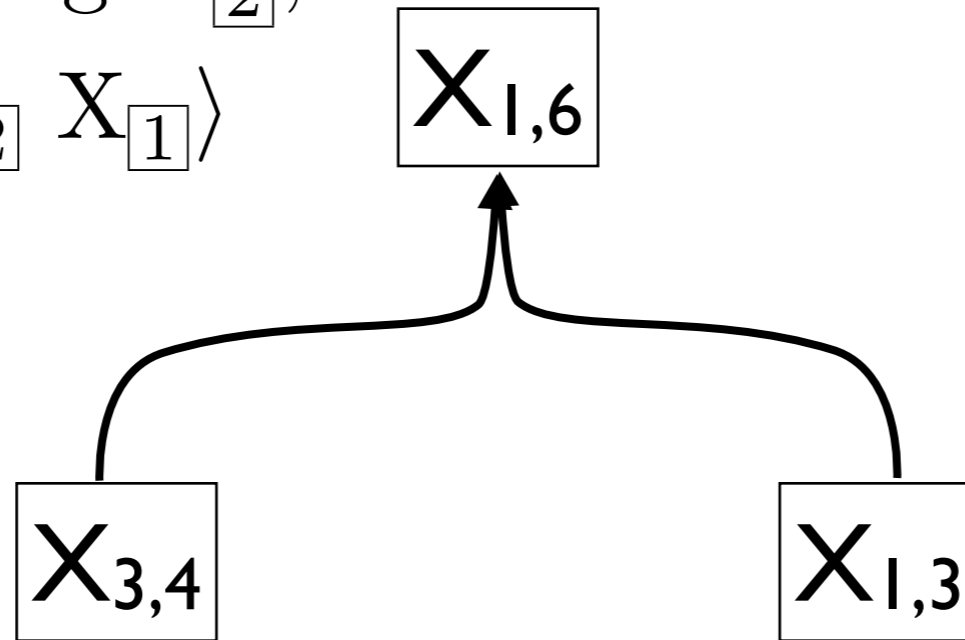
$X \rightarrow \langle X_{[1]} \text{ juxing } X_{[2]},$   
 $\text{held } X_{[2]} X_{[1]} \rangle$



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

# Non-Local Features

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



a talk

talks

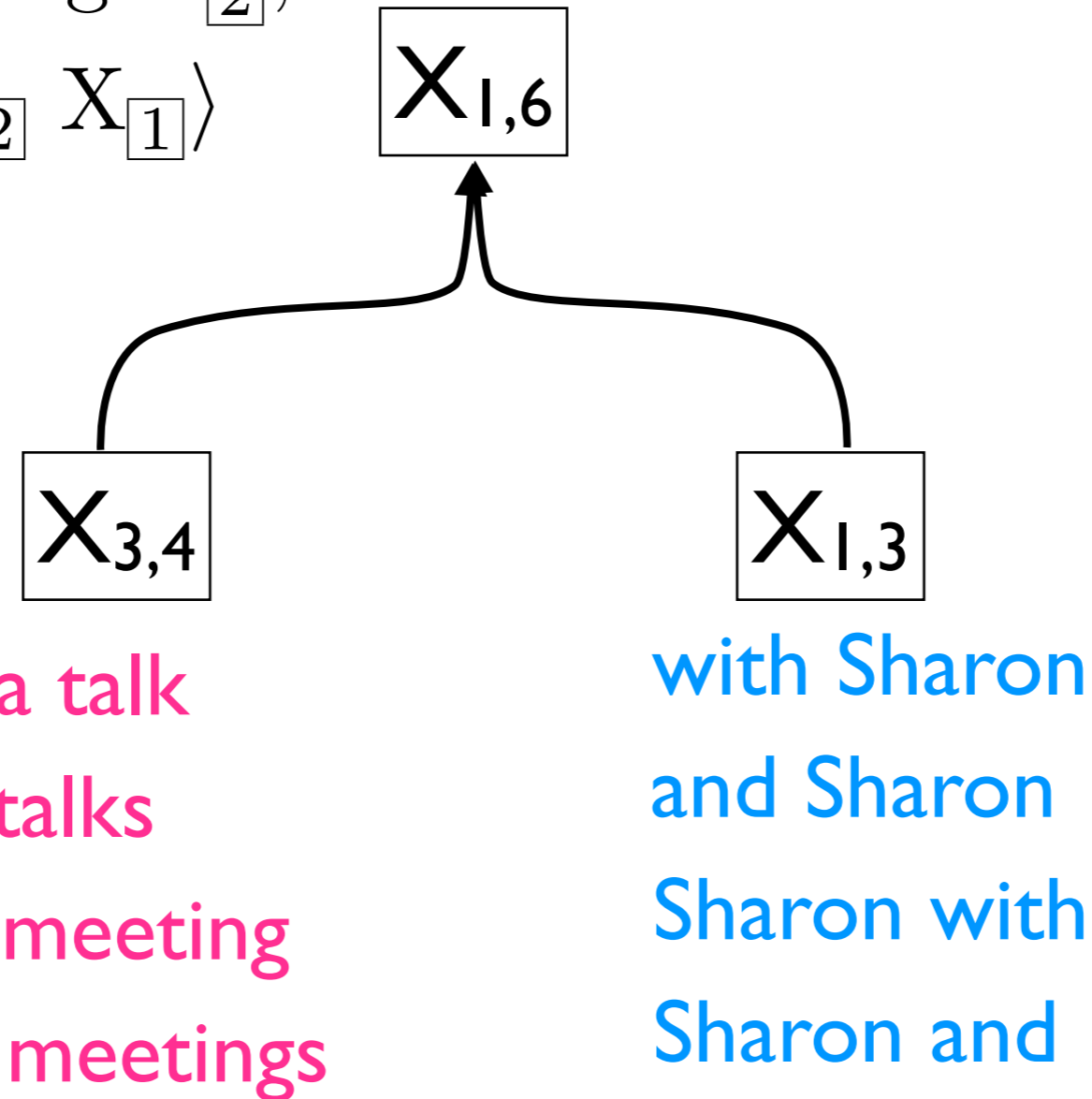
meeting

meetings

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

# Non-Local Features

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

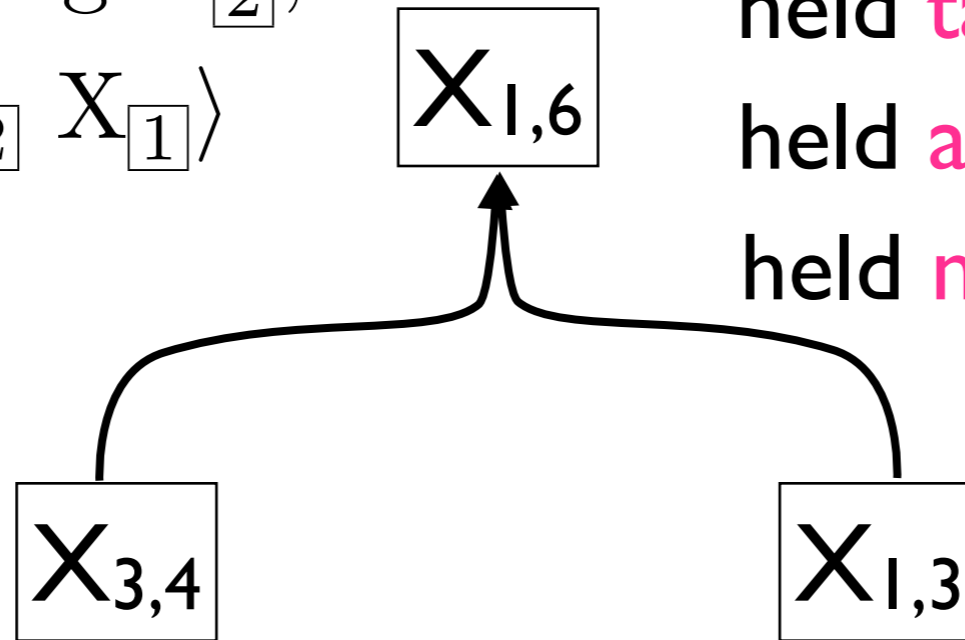


- non-local features which requires out-of-span context, i.e. bigram LM

# Non-Local Features

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

held **a talk** with Sharon  
 held **talks** with Sharon  
 held **a talk** and Sharon  
 held **meeting** Sharon with



**a talk**  
**talks**  
**meeting**  
**meetings**

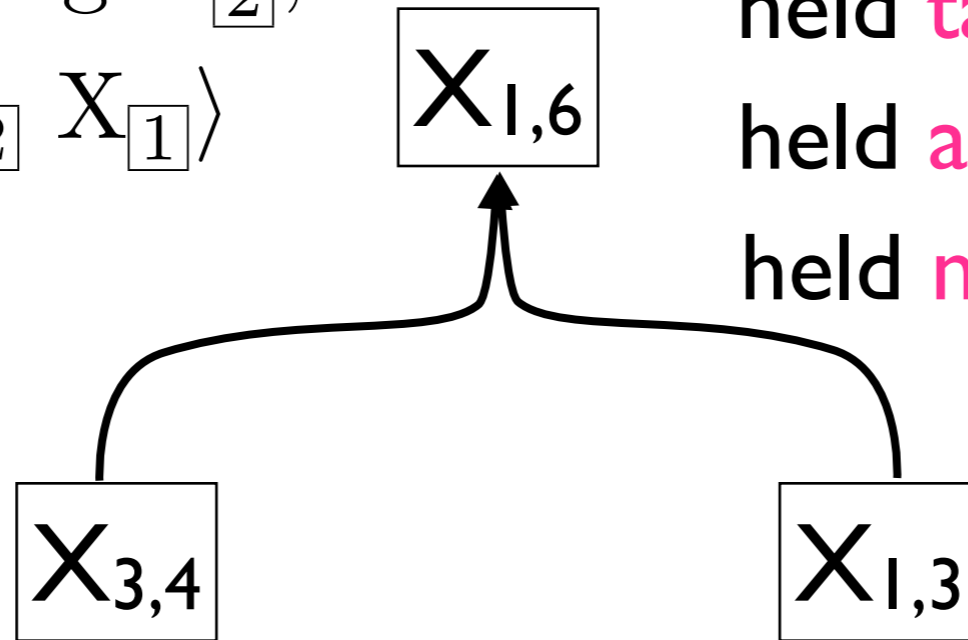
**with Sharon**  
**and Sharon**  
**Sharon with**  
**Sharon and**

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

# Non-Local Features

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

held **a talk** with Sharon  
 held **talks** with Sharon  
 held **a talk** and Sharon  
 held **meeting** Sharon with



$p(\text{talk} \mid a)$  **a talk**  
**talks**  
**meeting**  
**meetings**

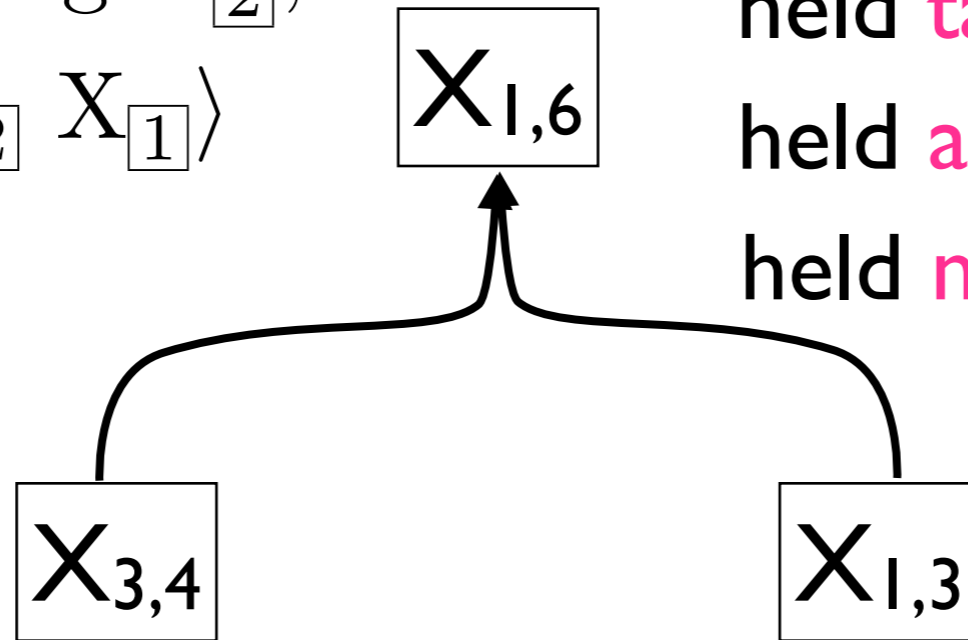
**with Sharon**  
**and Sharon**  
**Sharon with**  
**Sharon and**

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

# Non-Local Features

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

held **a talk** with Sharon  
 held **talks** with Sharon  
 held **a talk** and Sharon  
 held **meeting** Sharon with



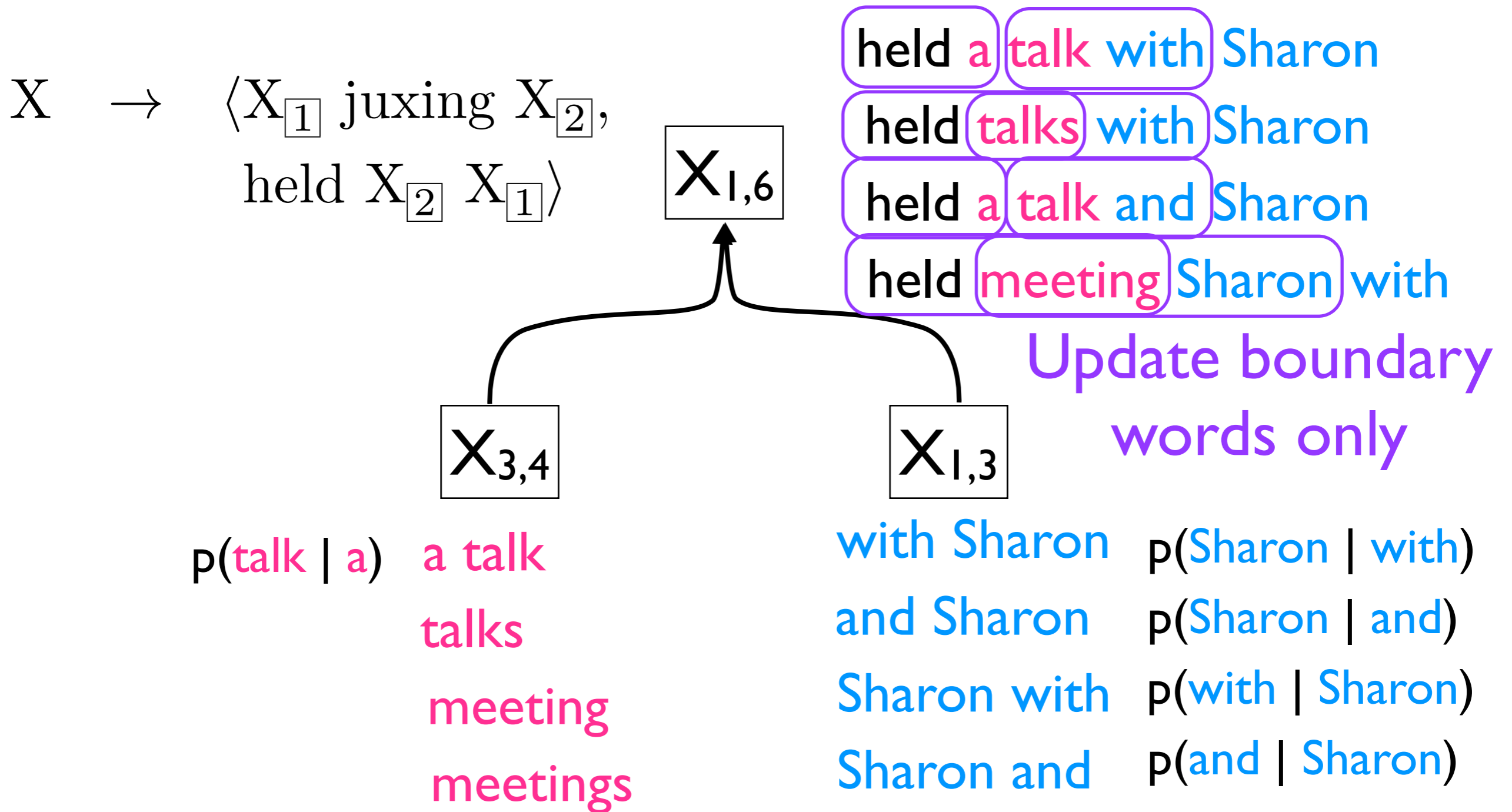
$p(\text{talk} \mid \text{a})$  **a talk**  
**talks**  
**meeting**  
**meetings**

**with Sharon**  $p(\text{Sharon} \mid \text{with})$   
**and Sharon**  $p(\text{Sharon} \mid \text{and})$   
**Sharon with**  $p(\text{with} \mid \text{Sharon})$   
**Sharon and**  $p(\text{and} \mid \text{Sharon})$

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

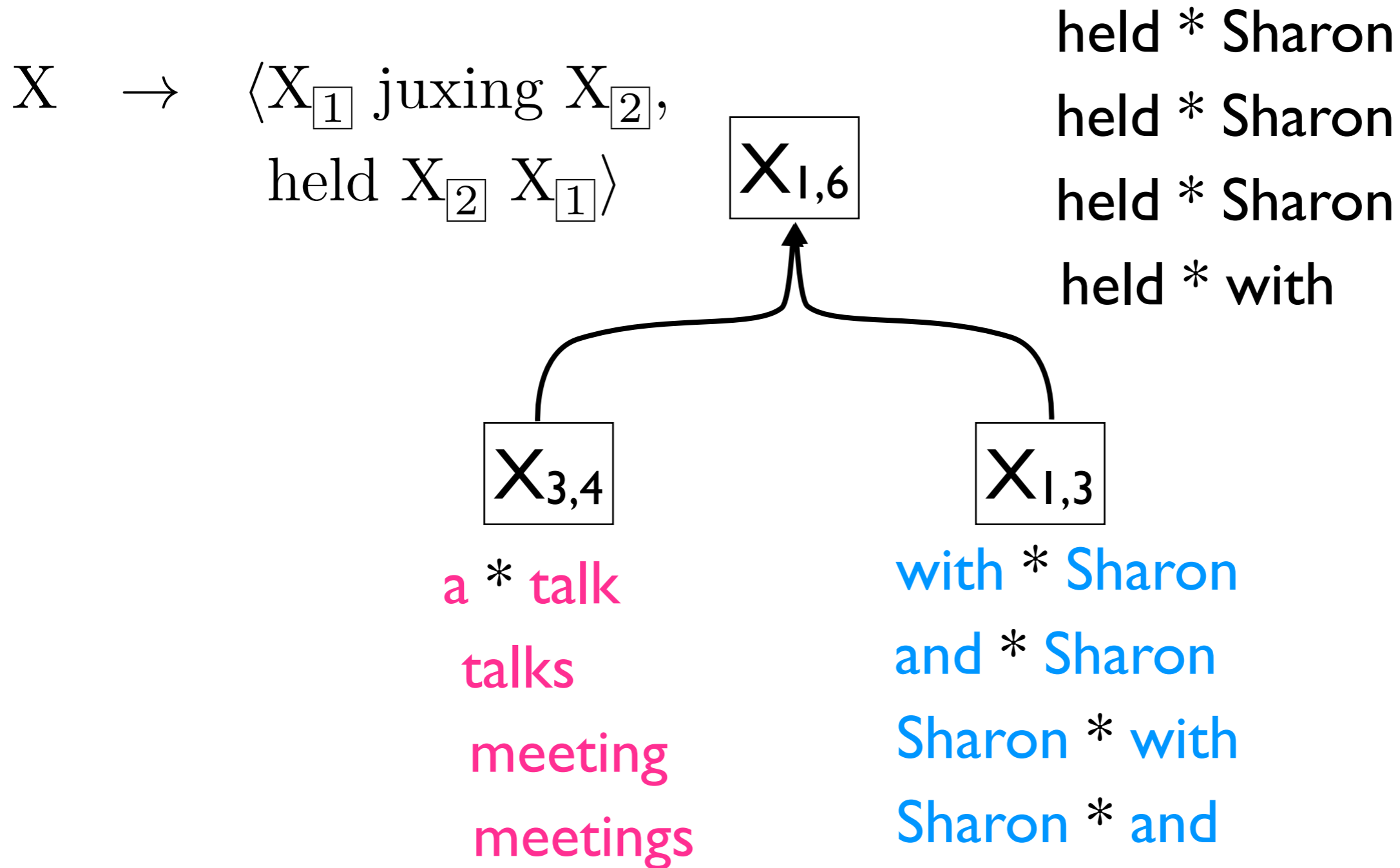


# Non-Local Features



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

# Bigram Features



- We keep only bigram states

# 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:  $O(n^3V^{2(m-1)})$
- Very inefficient: we need to explicitly enumerate all the hypotheses in antecedents

# Cube Pruning

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

*with \* Sharon* 1.5      *and \* Sharon* 1.7      *Sharon \* with* 2.6      *Sharon \* and* 3.2

<i>a * talk</i>	1.0	2.5	2.7	3.6	4.2
<i>talks</i>	1.3	2.8	3.0	3.9	4.5
<i>meeting</i>	2.2	3.7	3.9	4.8	5.4
<i>meetings</i>	2.6	4.1	4.3	5.2	5.8

- For each rule, create a “cube” representing combinations of antecedents (Huang and Chiang, 2007)

# Cube Pruning

$X \rightarrow \langle X_{[1]} \text{ juxting } X_{[2]}, \text{ held } X_{[2]} X_{[1]} \rangle$

*with \* Sharon* 1.5      *and \* Sharon* 1.7      *Sharon \* with* 2.6      *Sharon \* and* 3.2

<i>a * talk</i>	1.0	2.5 <i>+0.5</i>	2.7 <i>+1.0</i>	3.6 <i>+1.5</i>	4.2 <i>+2.5</i>
<i>talks</i>	1.3	2.8 <i>+0.3</i>	3.0 <i>+1.5</i>	3.9 <i>+2.0</i>	4.5 <i>+2.0</i>
<i>meeting</i>	2.2	3.7 <i>+0.5</i>	3.9 <i>+1.0</i>	4.8 <i>+1.5</i>	5.4 <i>+2.5</i>
<i>meetings</i>	2.6	4.1 <i>+0.3</i>	4.3 <i>+1.5</i>	5.2 <i>+2.0</i>	5.8 <i>+2.0</i>

- Bigrams require contexts from antecedents:  
non-monotonic scoring

# Cube Pruning

queue: (0,0)

k-best:

		<i>with * Sharon</i> 1.5	<i>and * Sharon</i> 1.7	<i>Sharon * with</i> 2.6	<i>Sharon * and</i> 3.2
<i>a * talk</i>	1.0	3.0			
<i>talks</i>	1.3				
<i>meeting</i>	2.2				
<i>meetings</i>	2.6				

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

# Cube Pruning

queue:

k-best: (0,0)

		<i>with * Sharon</i> 1.5	<i>and * Sharon</i> 1.7	<i>Sharon * with</i> 2.6	<i>Sharon * and</i> 3.2
<i>a * talk</i>	1.0	3.0			
<i>talks</i>	1.3				
<i>meeting</i>	2.2				
<i>meetings</i>	2.6				

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

# Cube Pruning

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

k-best: (0,0)

		<i>with * Sharon</i> 1.5	<i>and * Sharon</i> 1.7	<i>Sharon * with</i> 2.6	<i>Sharon * and</i> 3.2
<i>a * talk</i>	1.0	3.0	3.7		
<i>talks</i>	1.3	3.1			
<i>meeting</i>	2.2				
<i>meetings</i>	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

# 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

<i>a * talk</i>	1.0	3.0	3.7		
<i>talks</i>	1.3	3.1	4.5		
<i>meeting</i>	2.2	4.2			
<i>meetings</i>	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*

*and \* Sharon*

*Sharon \* with*

*Sharon \* and*

1.5

1.7

2.6

3.2

<i>a * talk</i>	1.0	3.0	3.7		
<i>talks</i>	1.3	3.1	4.5		
<i>meeting</i>	2.2	4.2			
<i>meetings</i>	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

<i>a * talk</i>	1.0	3.0	3.7	5.1	
<i>talks</i>	1.3	3.1	4.5		
<i>meeting</i>	2.2	4.2			
<i>meetings</i>	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 = rule + cube-position

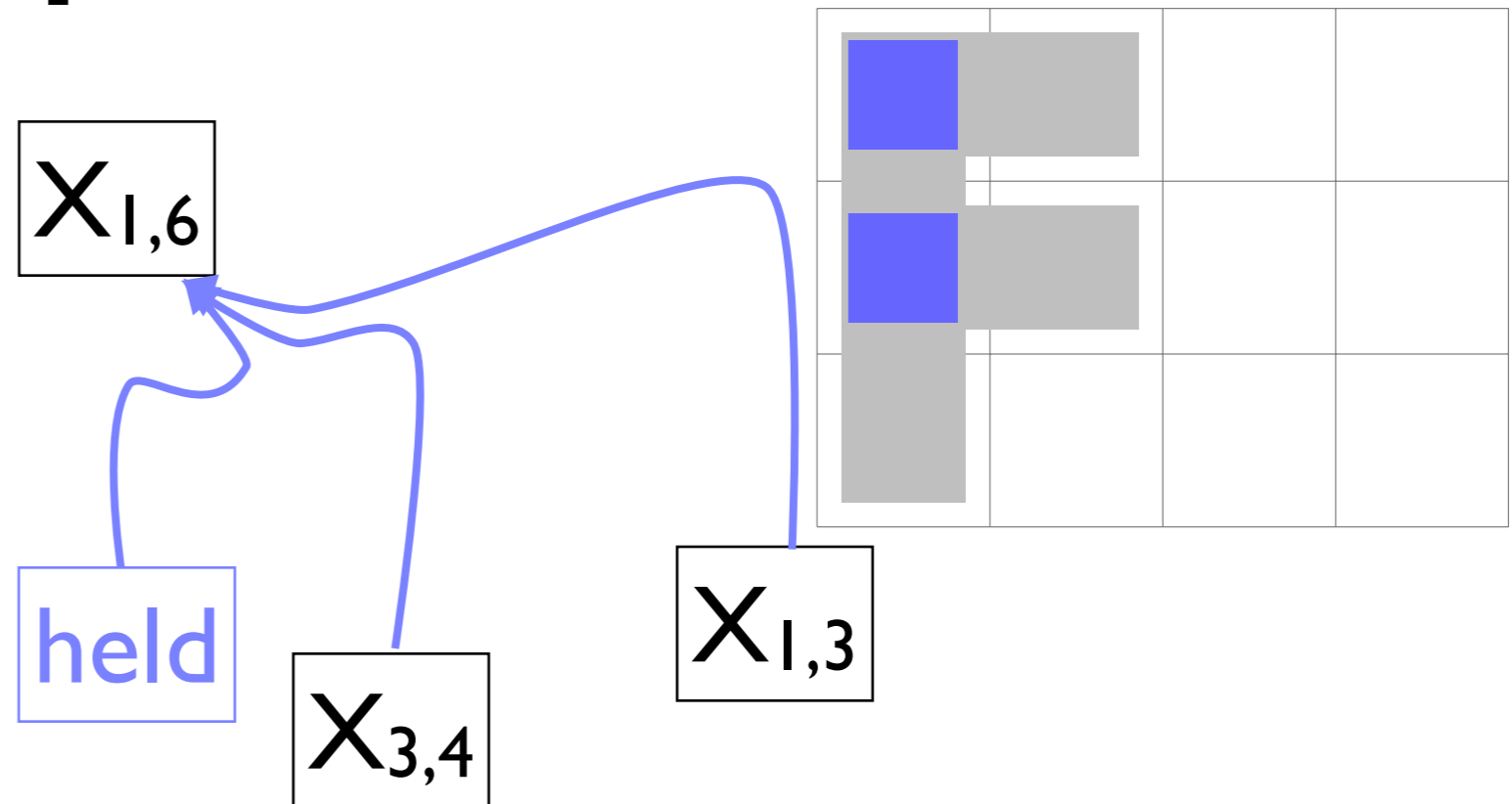
# Multiple Rules

$X_{1,6}$

- Multiple rules sharing the same span are queued
- Each rule is associated with a cube
- hypothesis = rule + cube-position

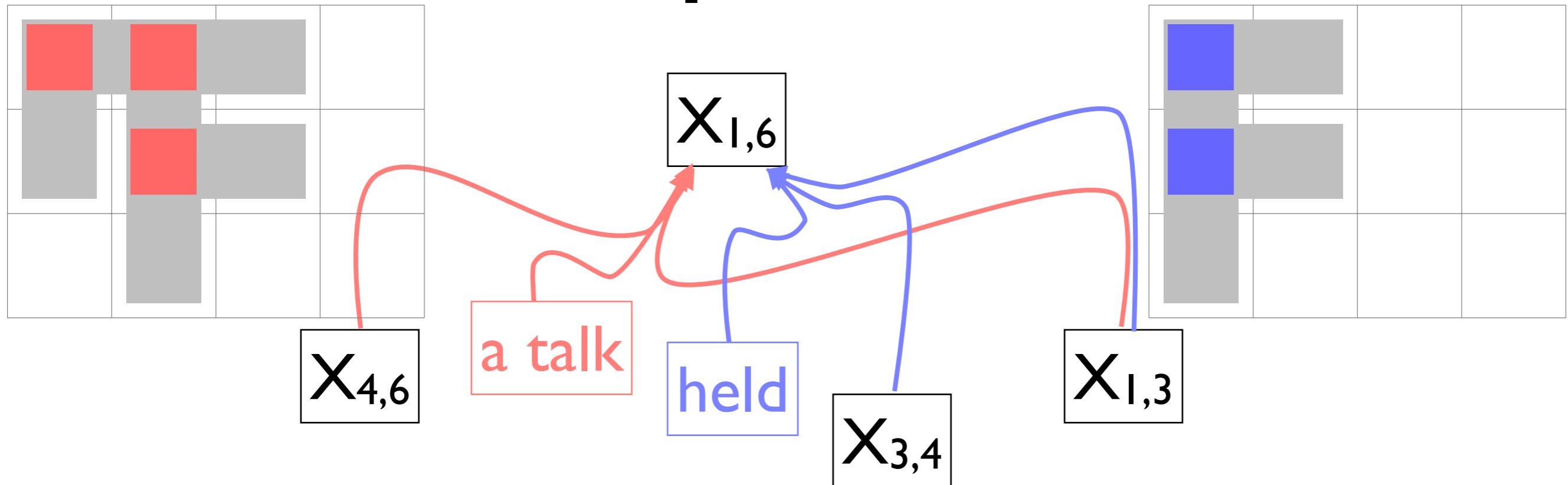


# Multiple Rules



- Multiple rules sharing the same span are queued
- Each rule is associated with a cube
- hypothesis = rule + cube-position

# Multiple Rules



- Multiple rules sharing the same span are queued
- Each rule is associated with a cube
- hypothesis = rule + 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

# 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

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing

# Overview

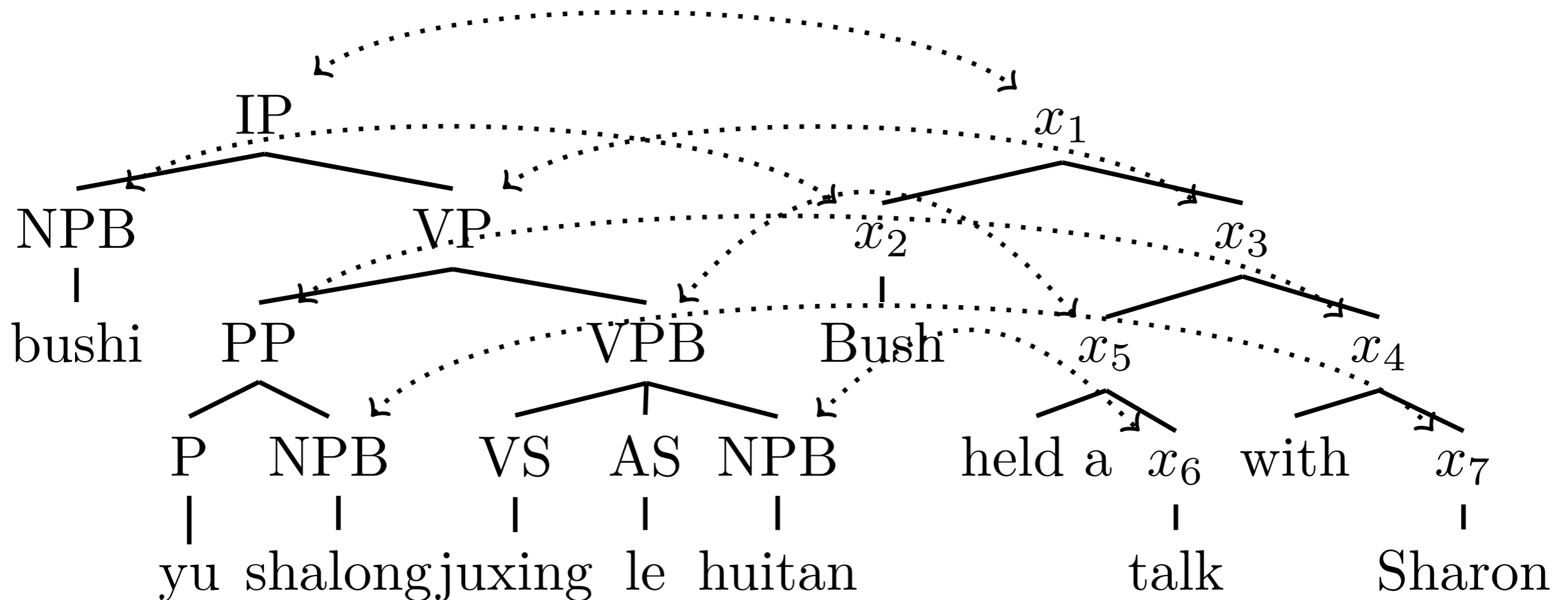
- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- **Tree-based SMT**
  - Synchronous-CFG
  - **String-to-Tree, Tree-to-String**
  - Bitext parsing

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

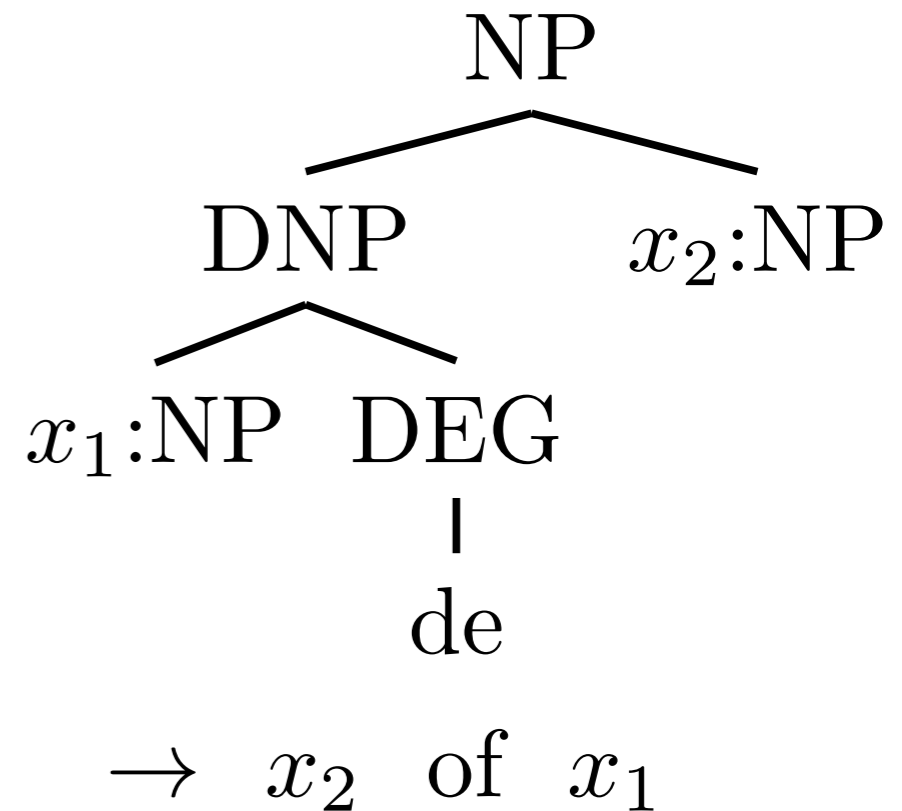
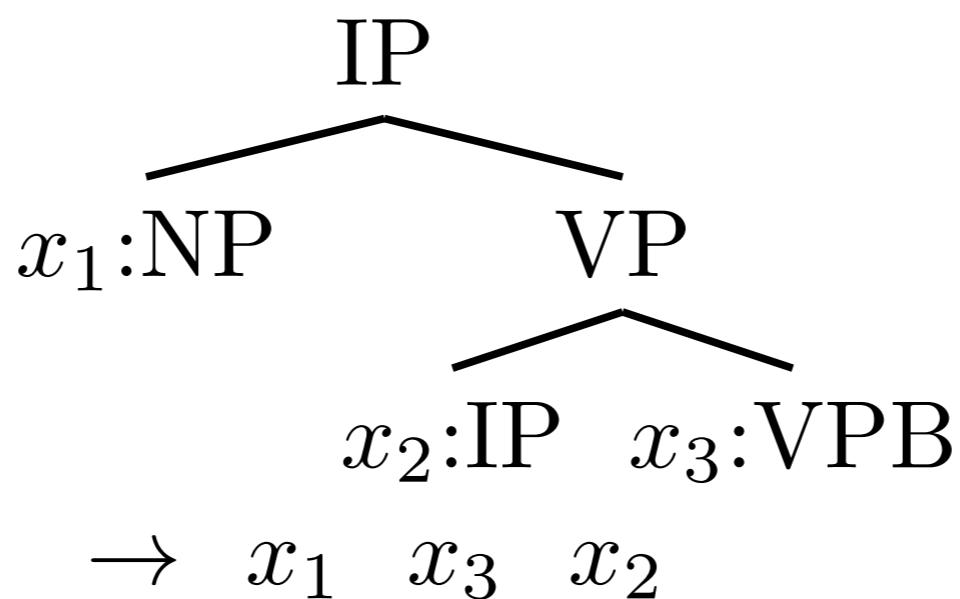
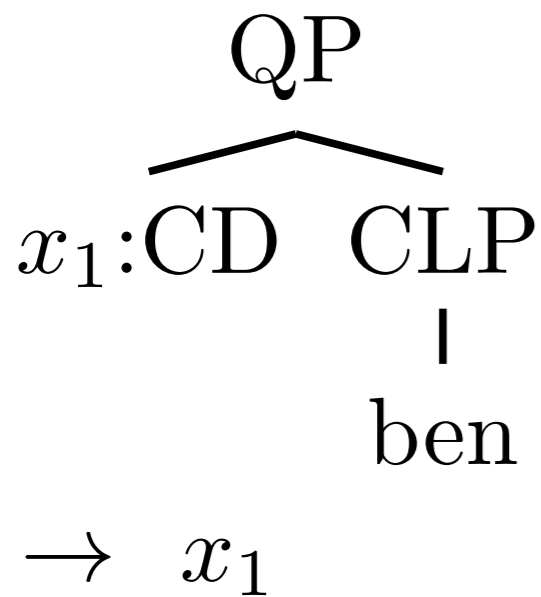
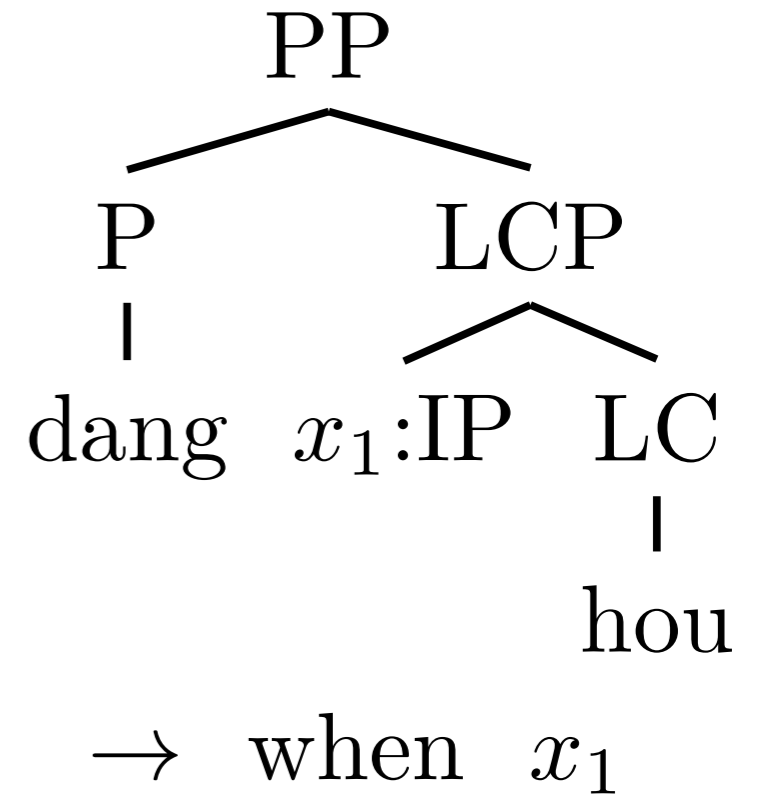
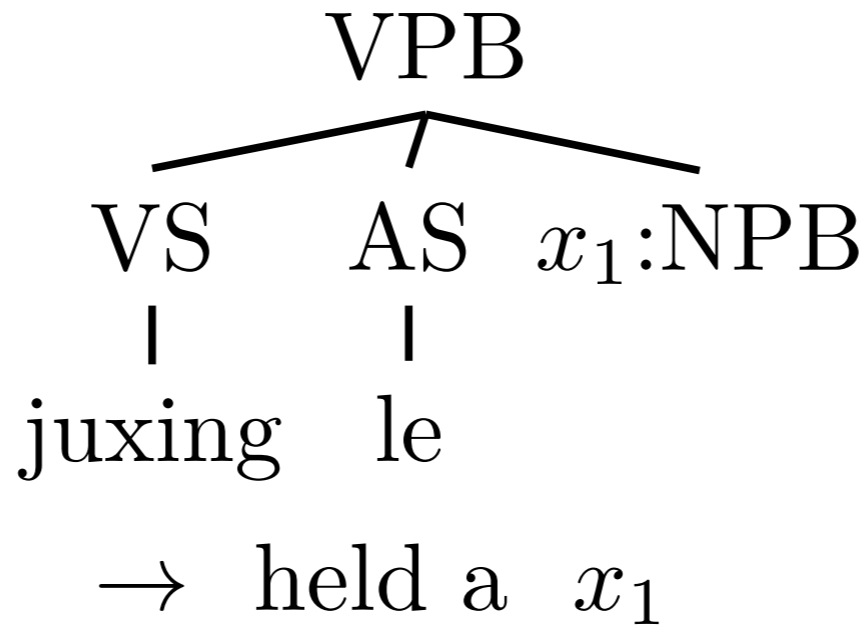
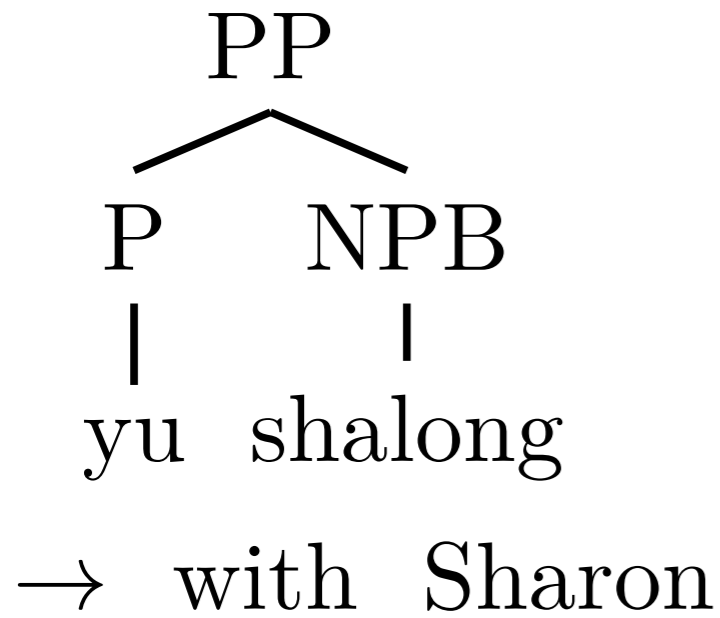


(Galley et al., 2004)

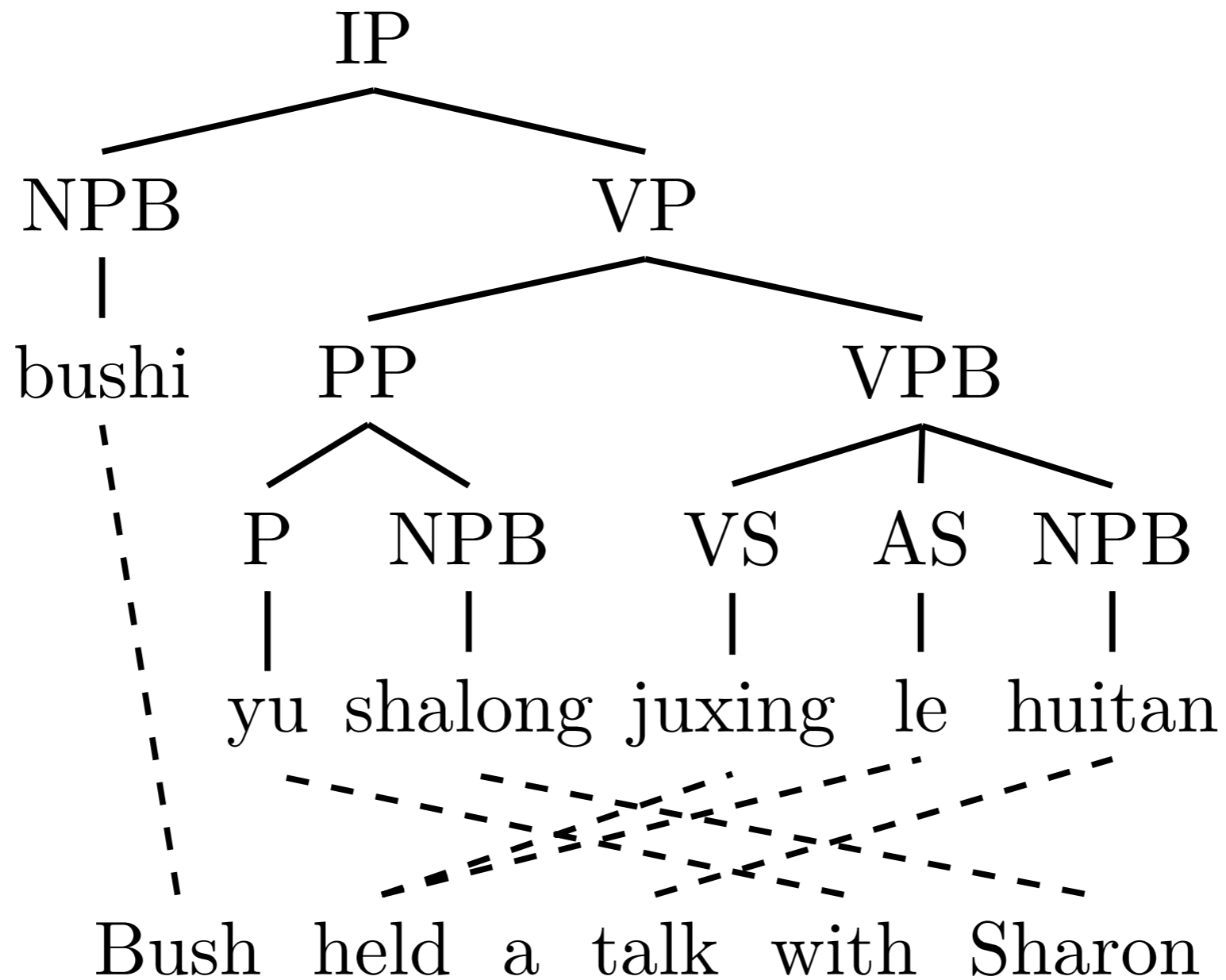
- Each synchronous rule has a subtree structure
- Flat structure + sharing the same non-terminal symbols = synchronous-CFG



# Tree-to-String Rules



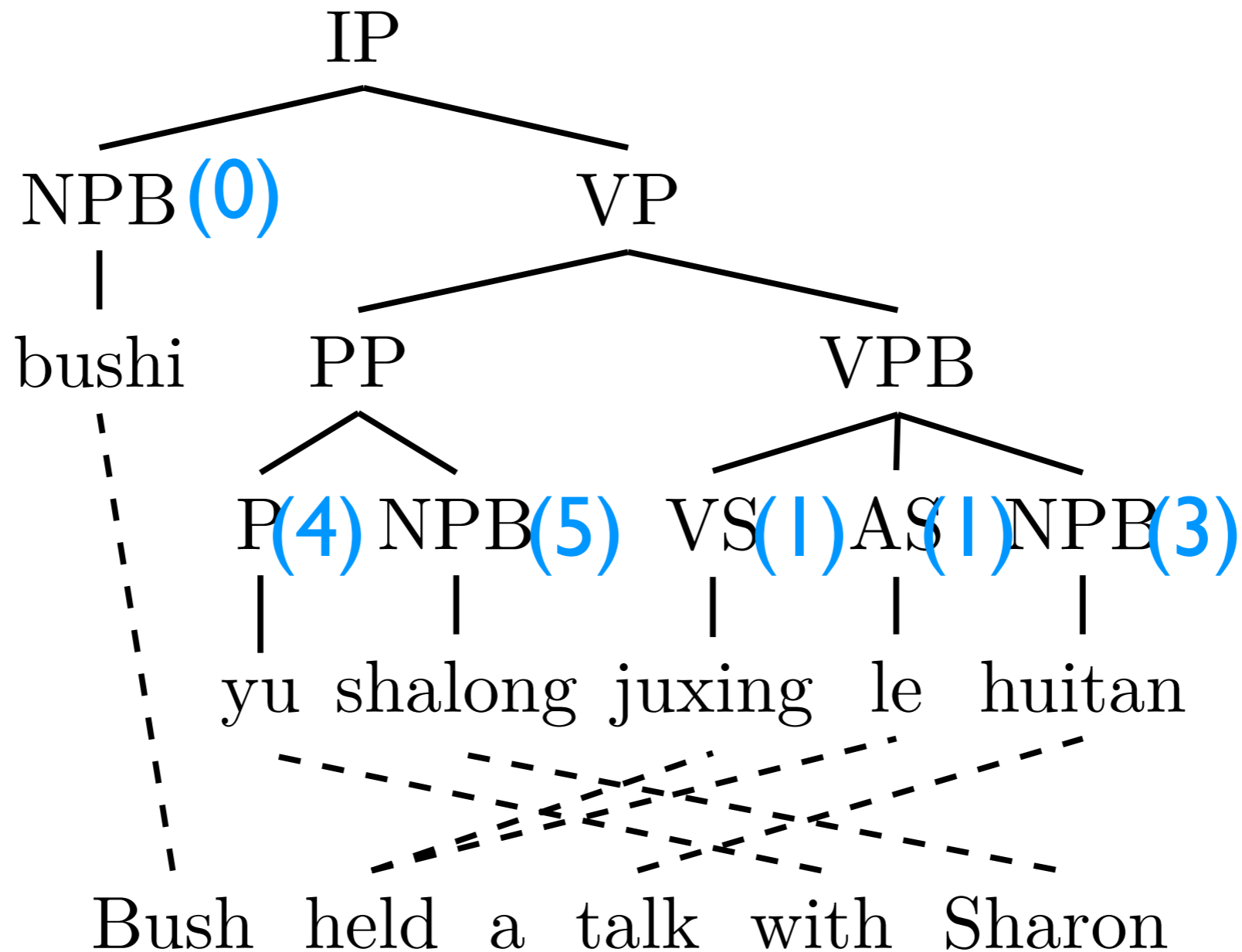
# Rule Extraction



(Galley et al., 2004)

- Compute “spans” by propagating alignment in bottom-up

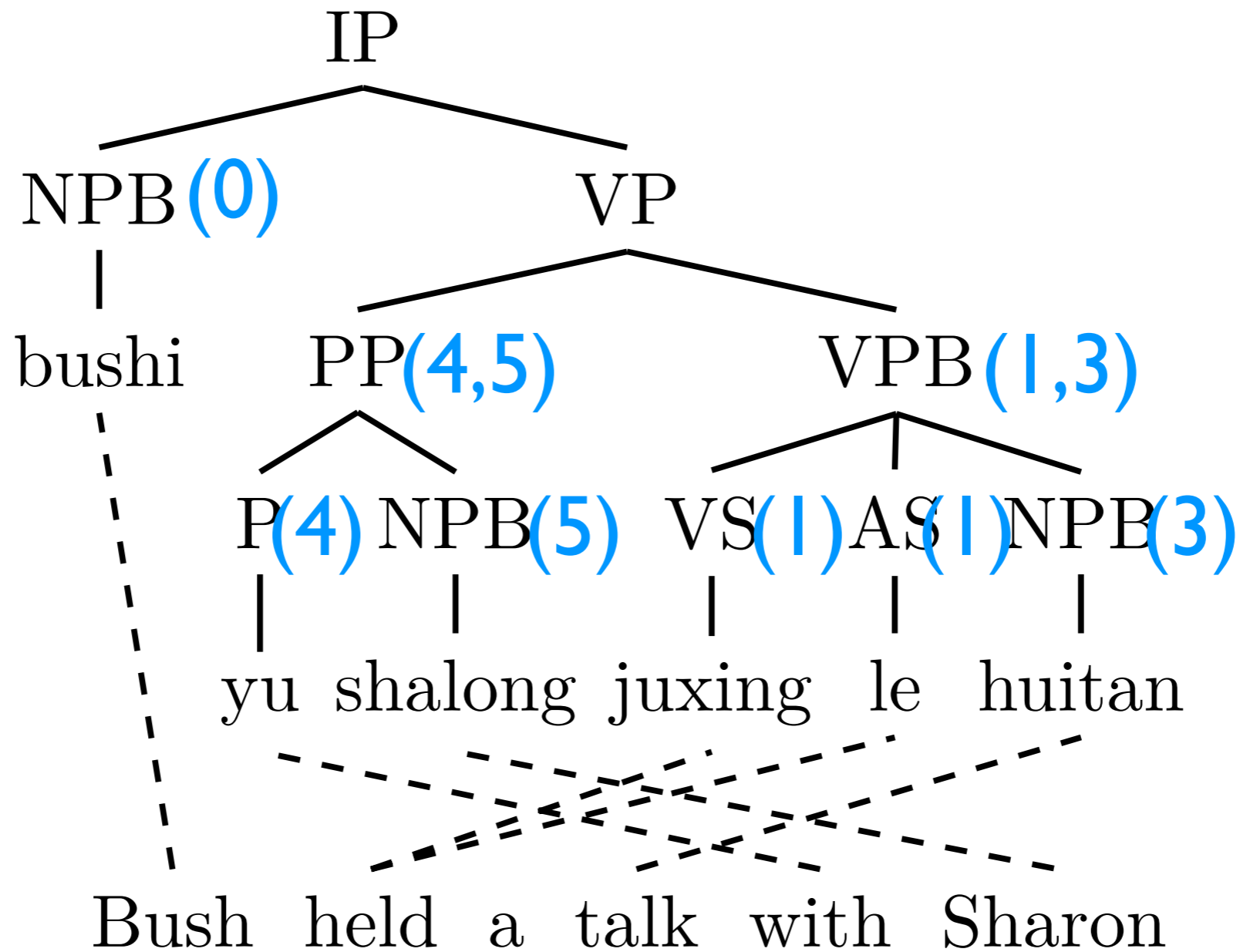
# Rule Extraction



(Galley et al., 2004)

- Compute “spans” by propagating alignment in bottom-up

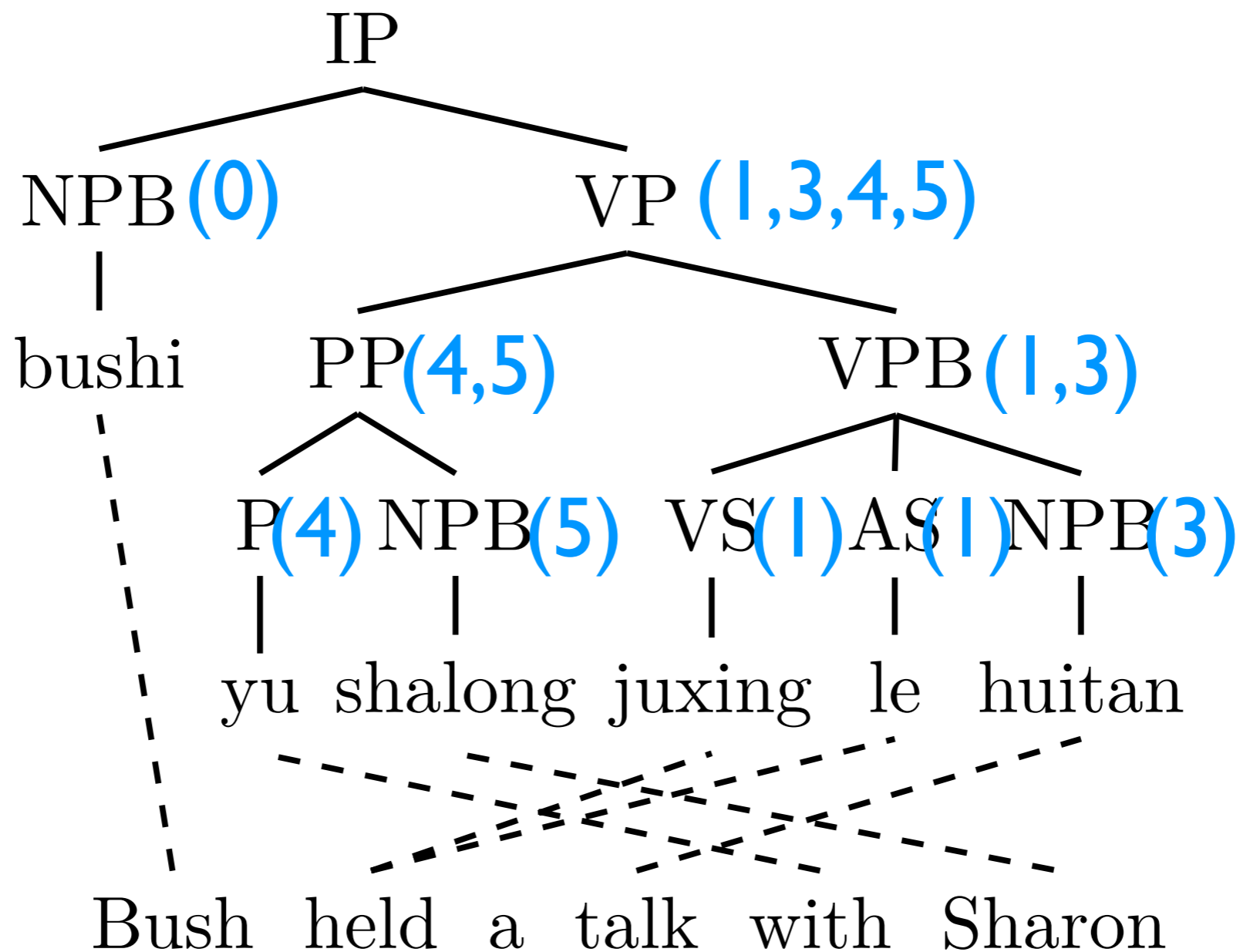
# Rule Extraction



(Galley et al., 2004)

- Compute “spans” by propagating alignment in bottom-up

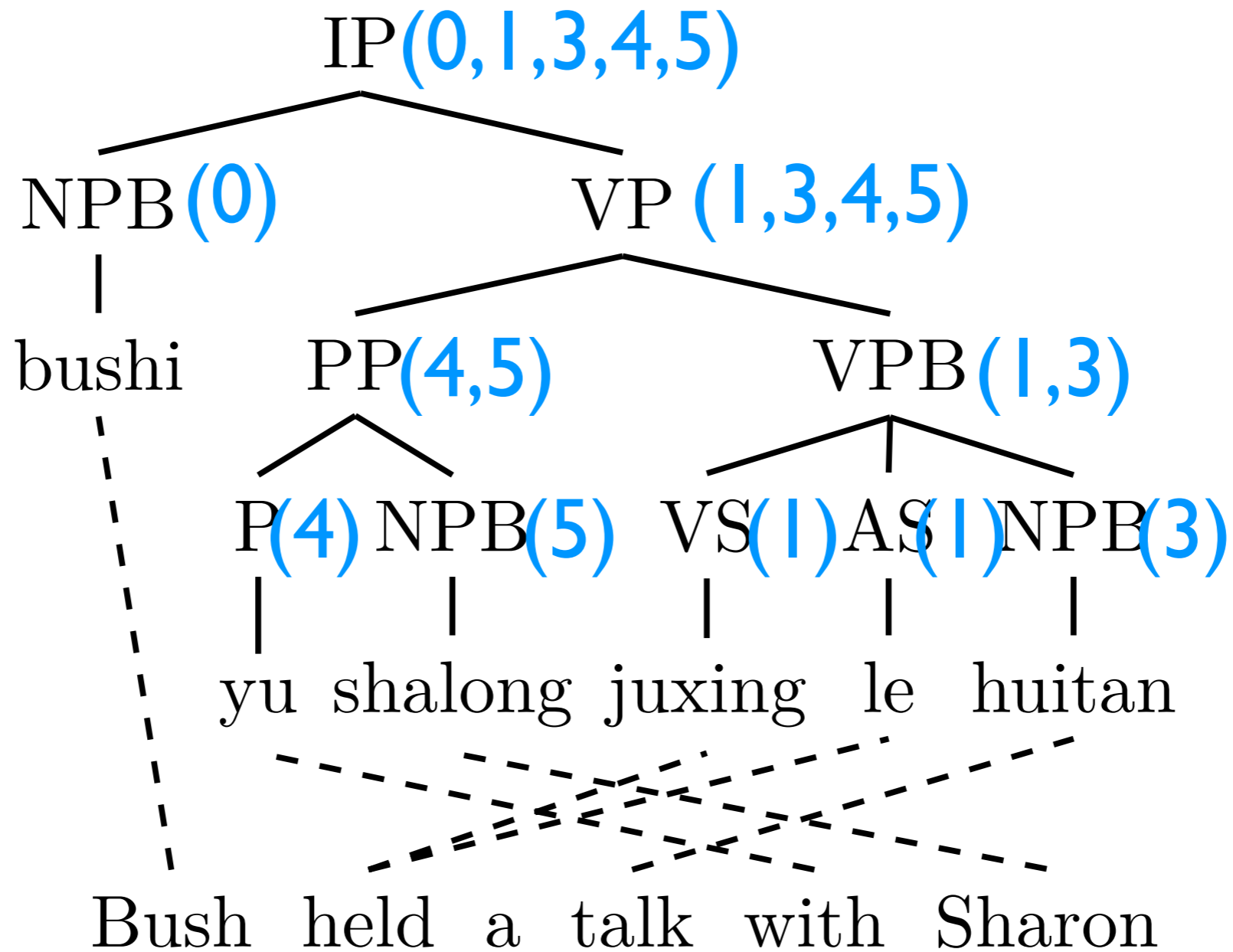
# Rule Extraction



(Galley et al., 2004)

- Compute “spans” by propagating alignment in bottom-up

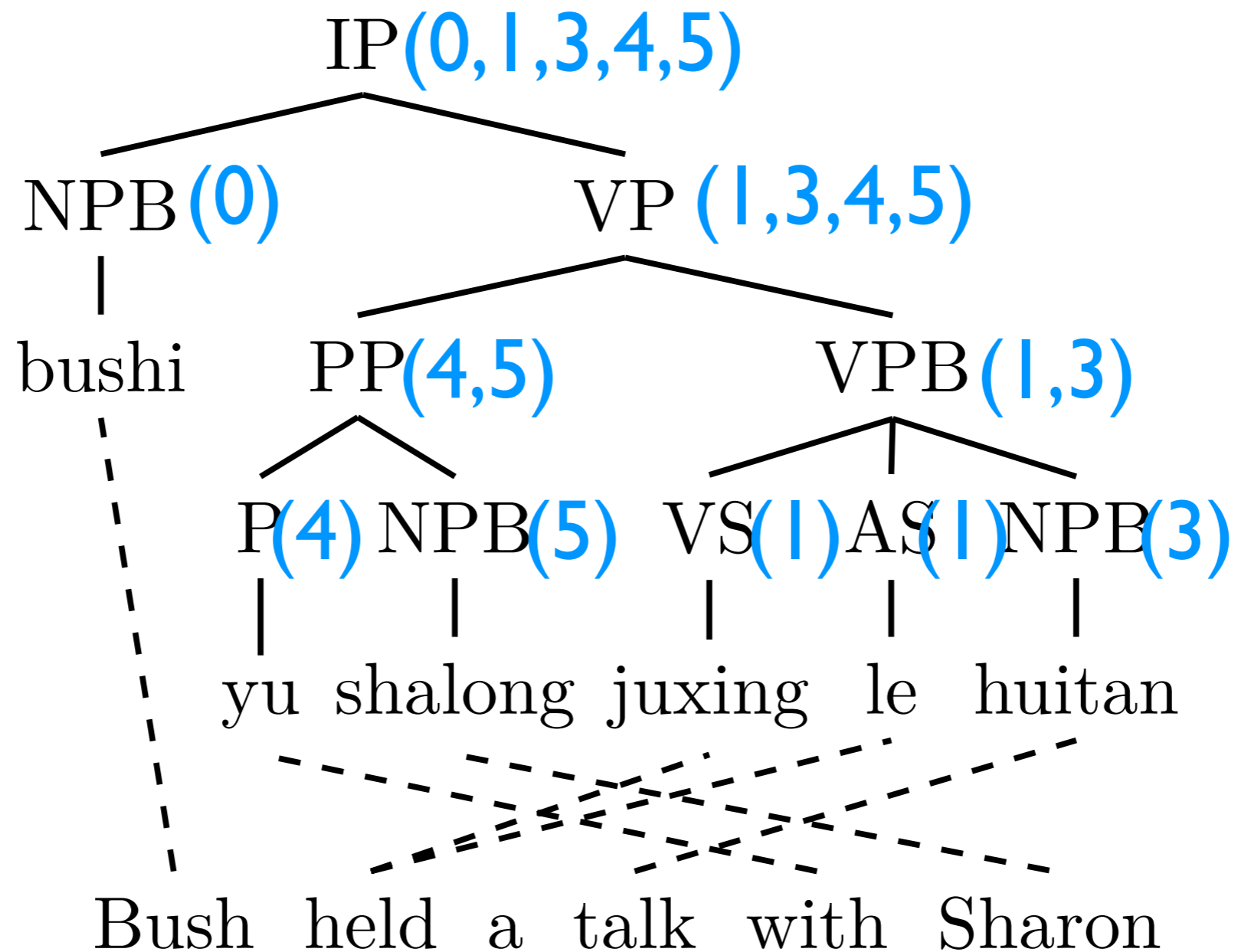
# Rule Extraction



(Galley et al., 2004)

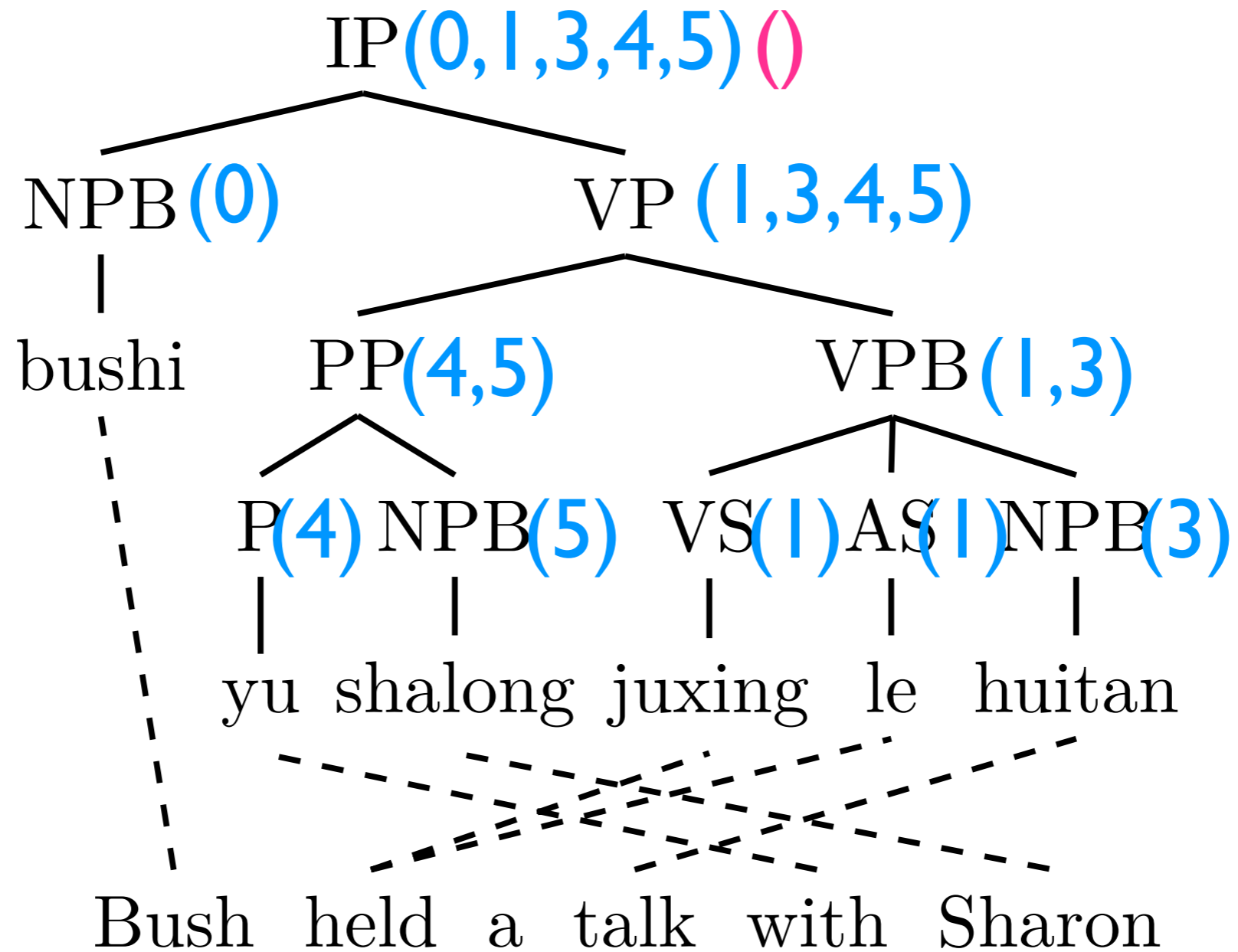
- Compute “spans” by propagating alignment in bottom-up

# Rule Extraction



- Compute “complements” in top-down

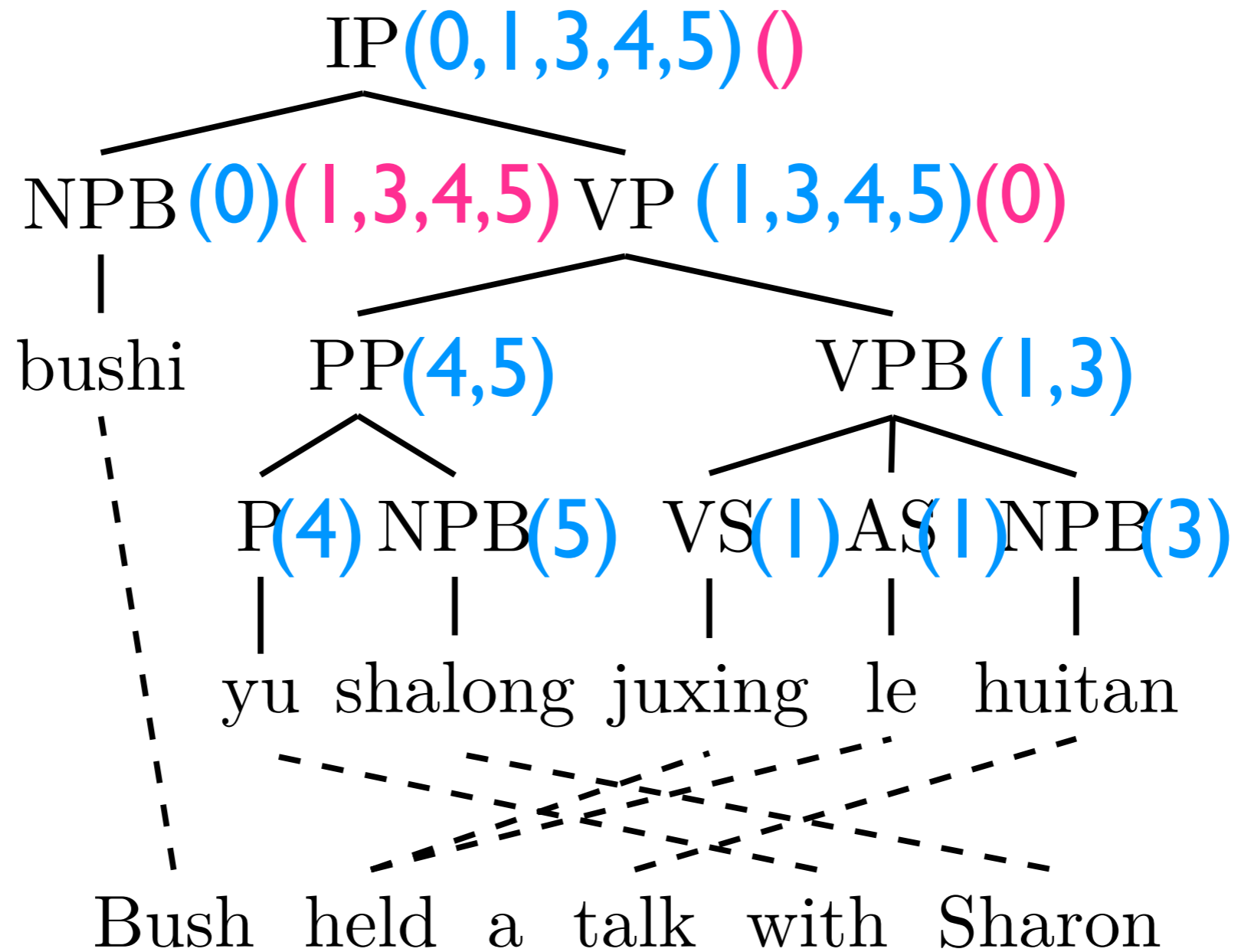
# Rule Extraction



- Compute “complements” in top-down

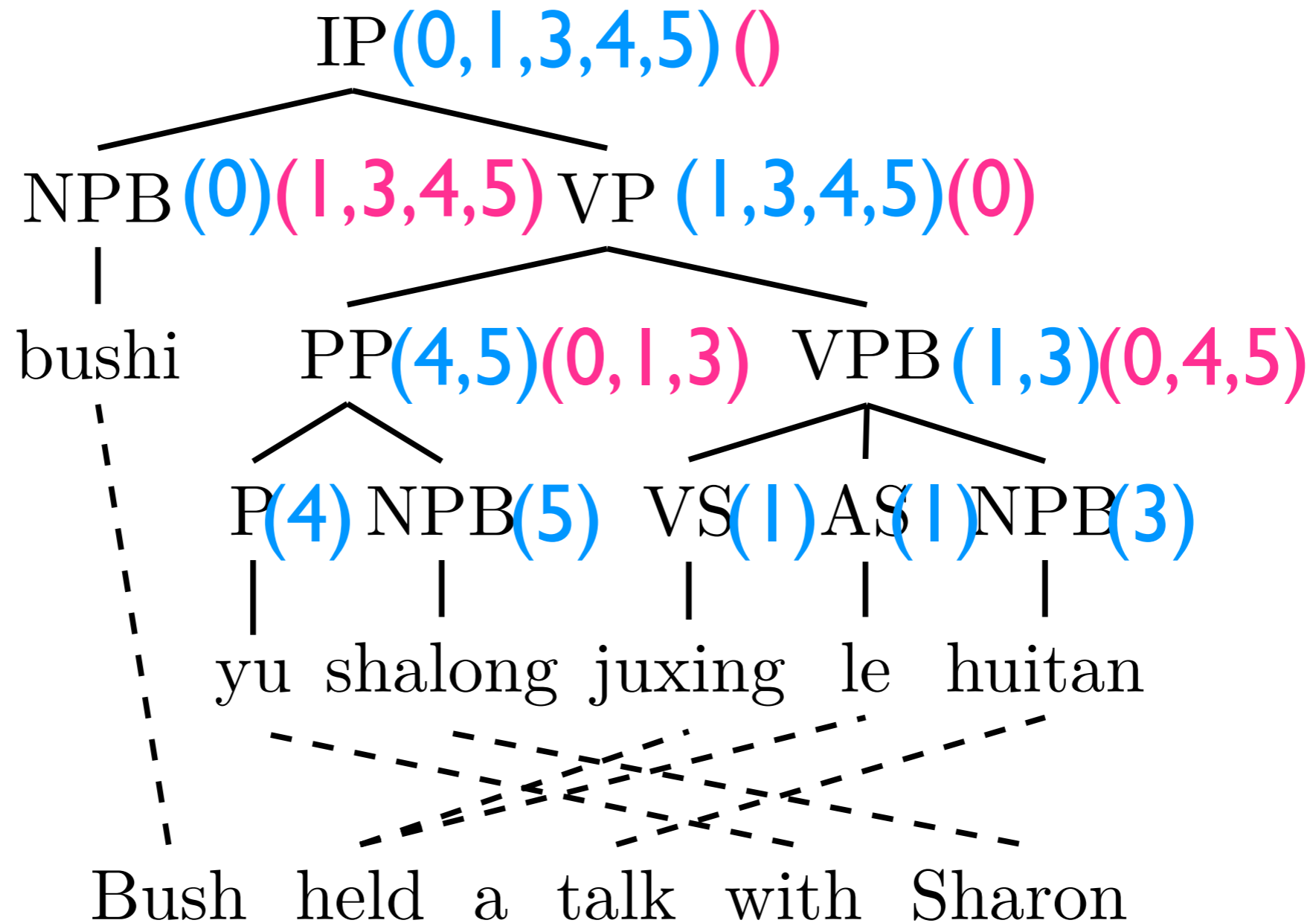


# Rule Extraction



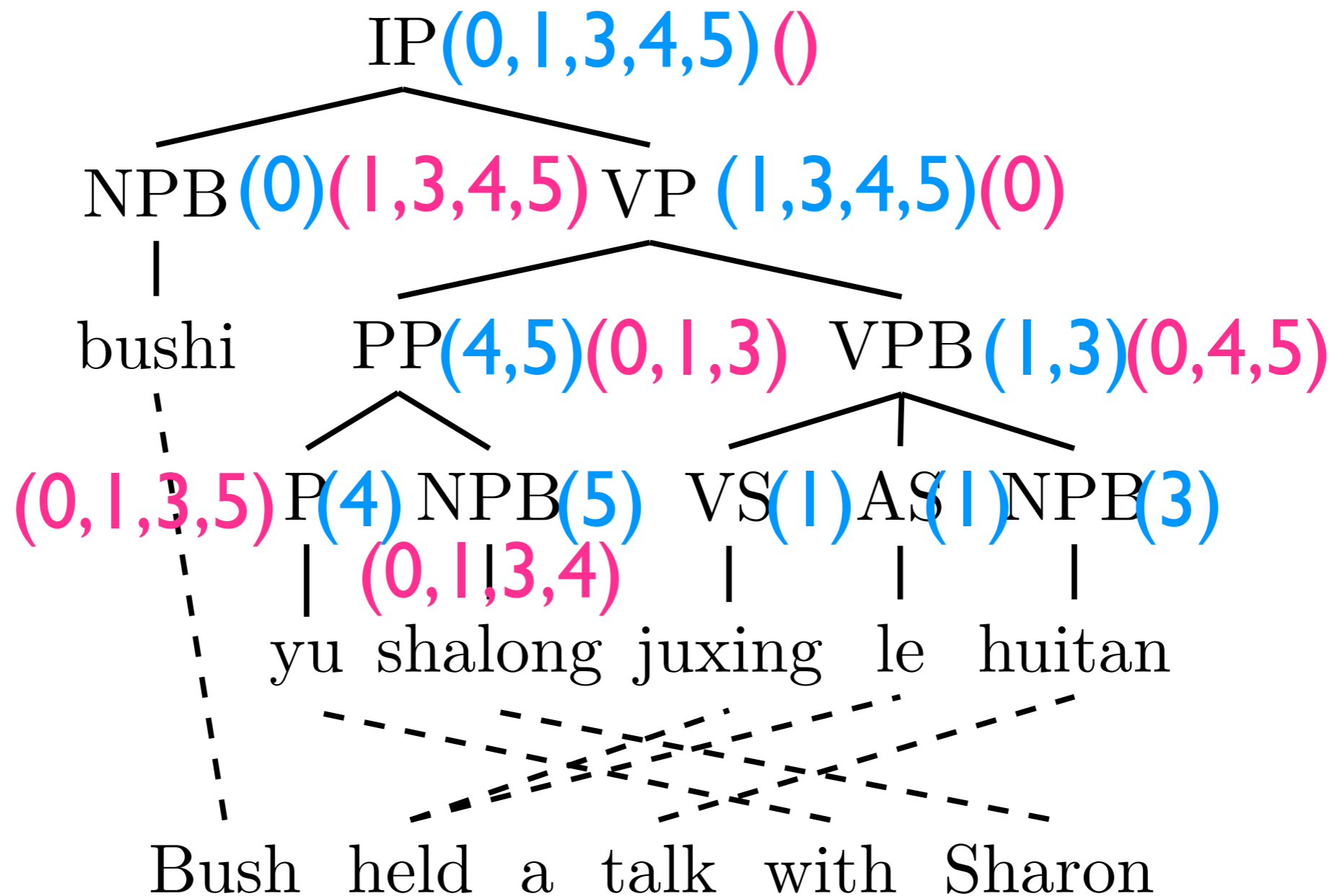
- Compute “complements” in top-down

# Rule Extraction



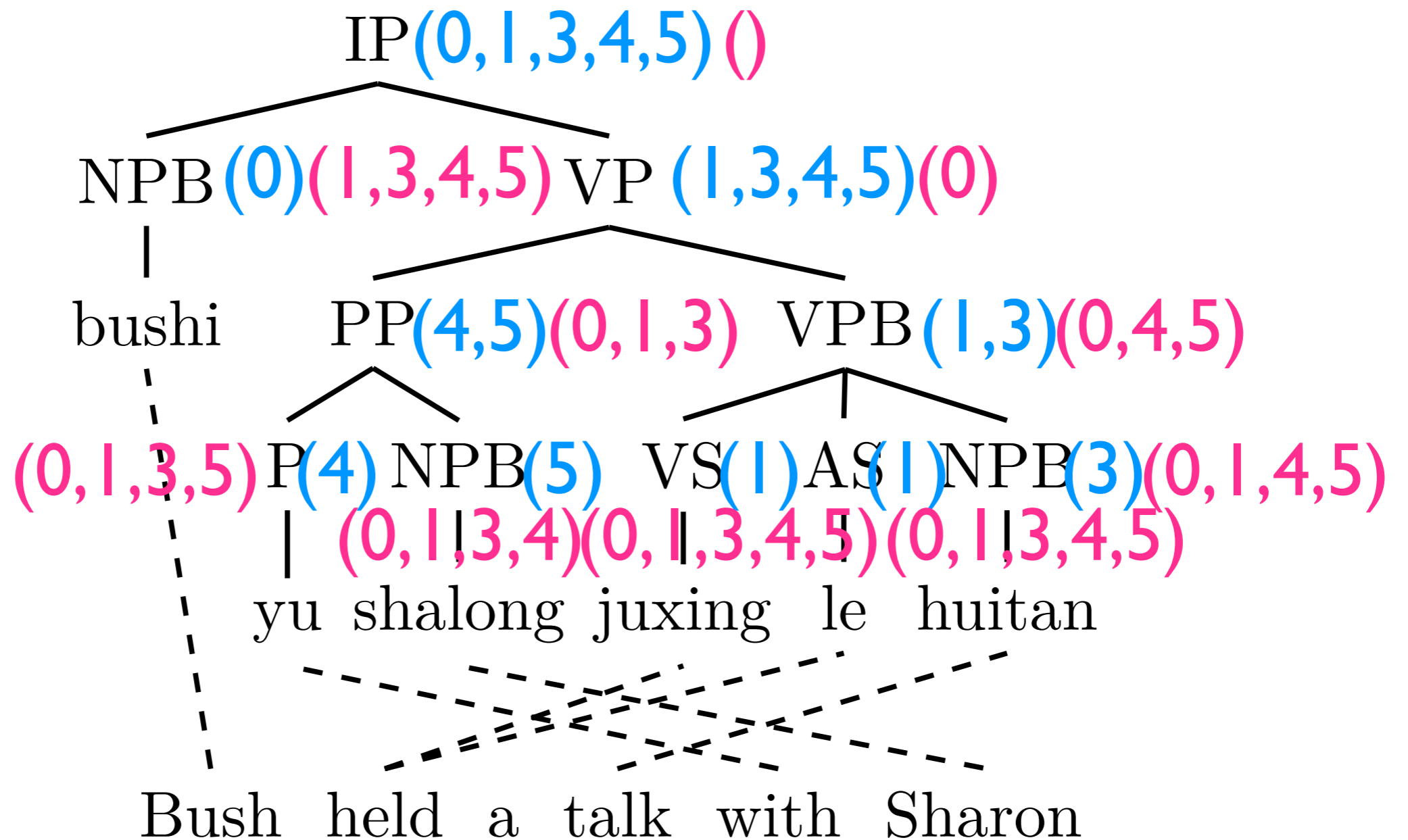
- Compute “complements” in top-down

# Rule Extraction



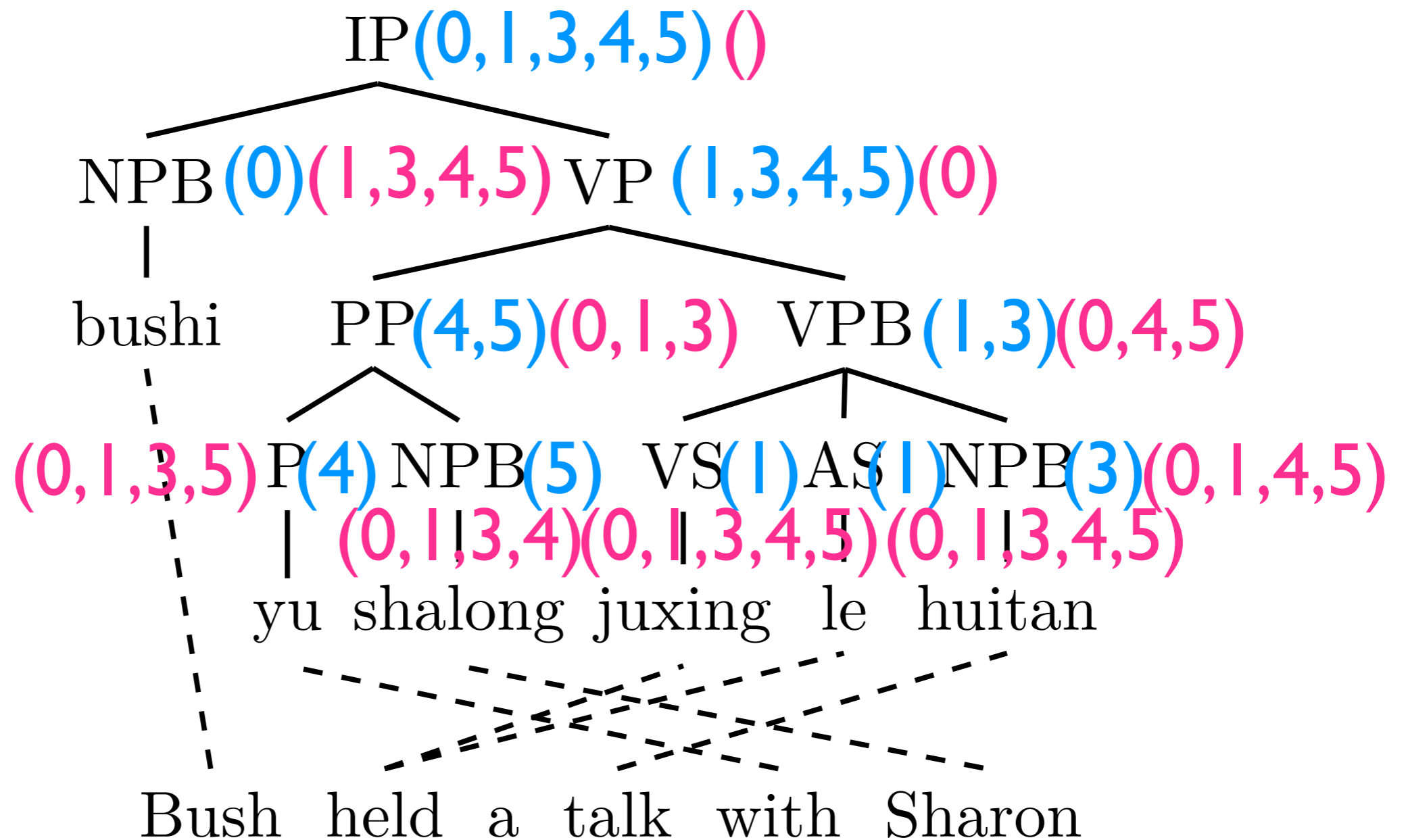
- Compute “complements” in top-down

# Rule Extraction



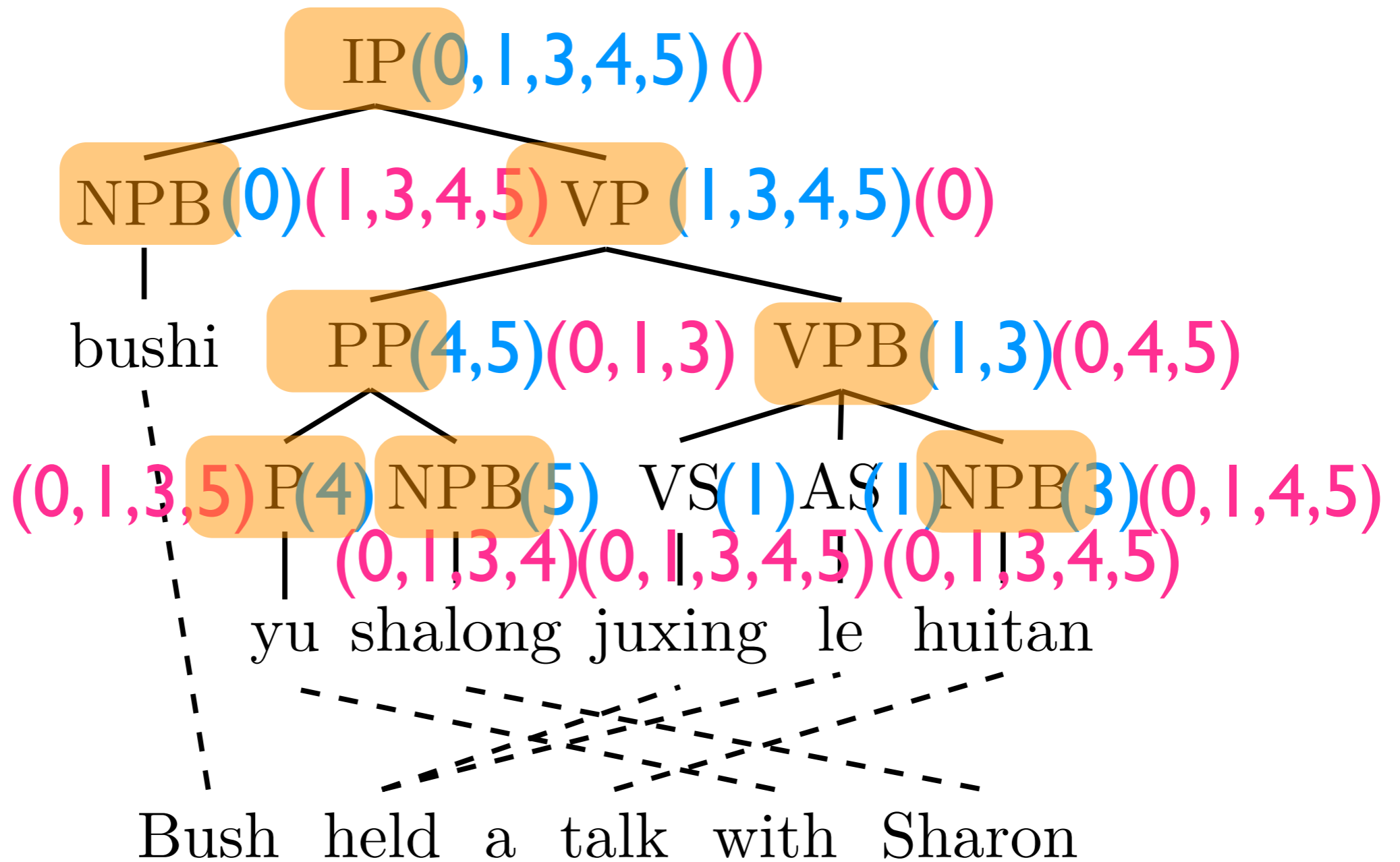
- Compute “complements” in top-down

# Rule Extraction



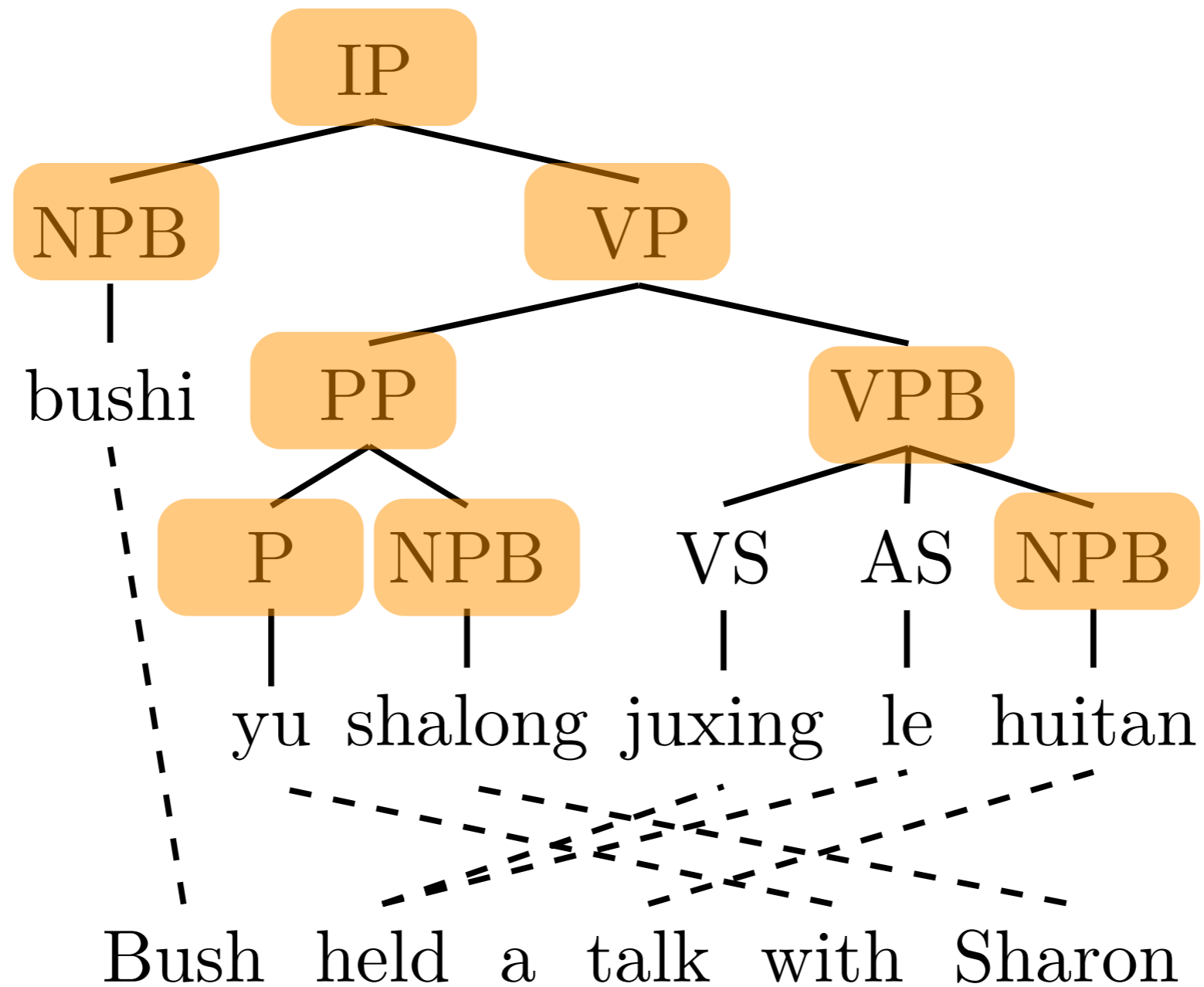
- Compute “frontiers”: The nodes in which the intersection of “spans” and “complements” is empty

# Rule Extraction



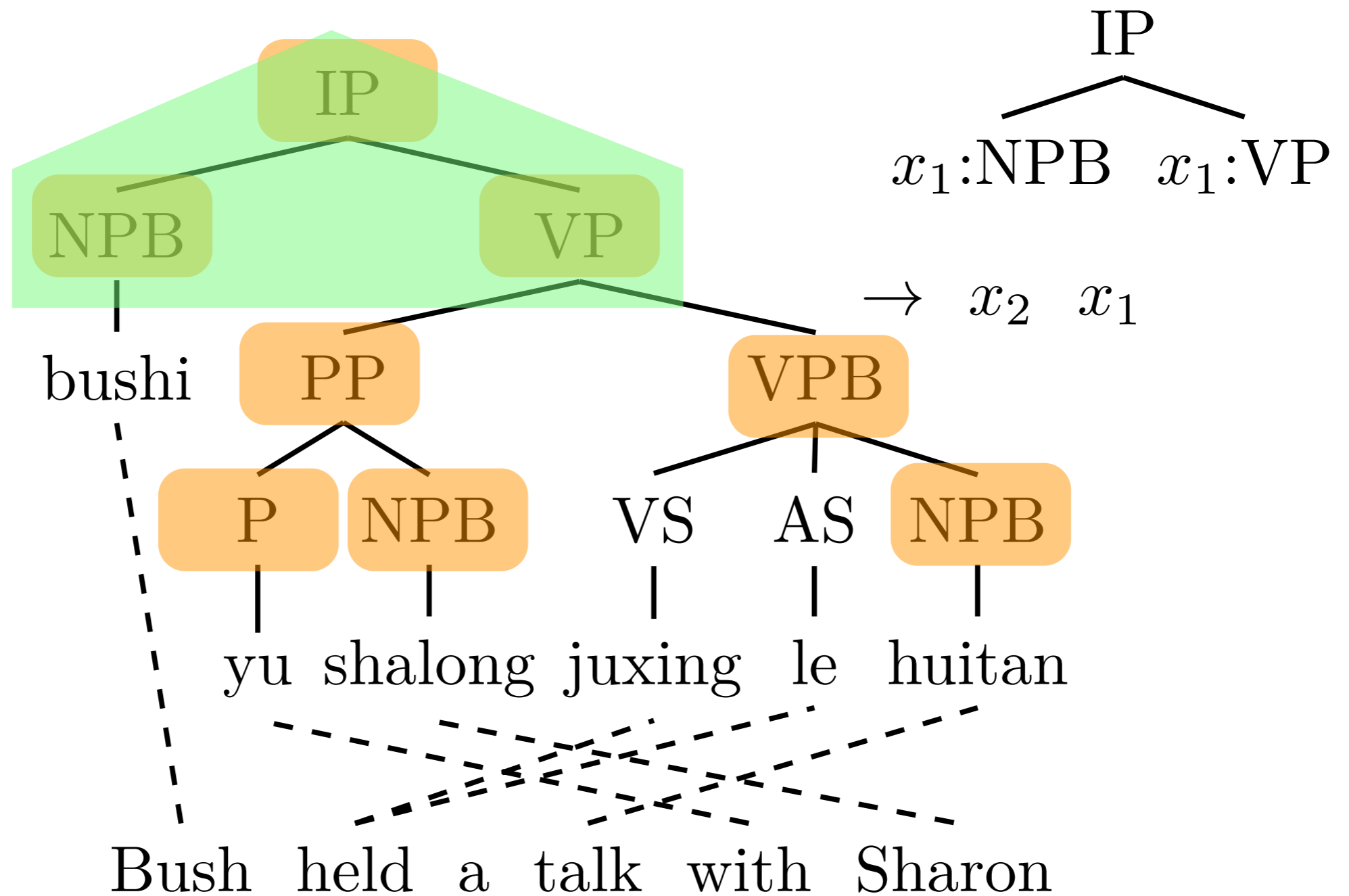
- Compute “frontiers”: The nodes in which the intersection of “spans” and “complements” is empty

# Rule Extraction



- Extract minimum rules using frontiers

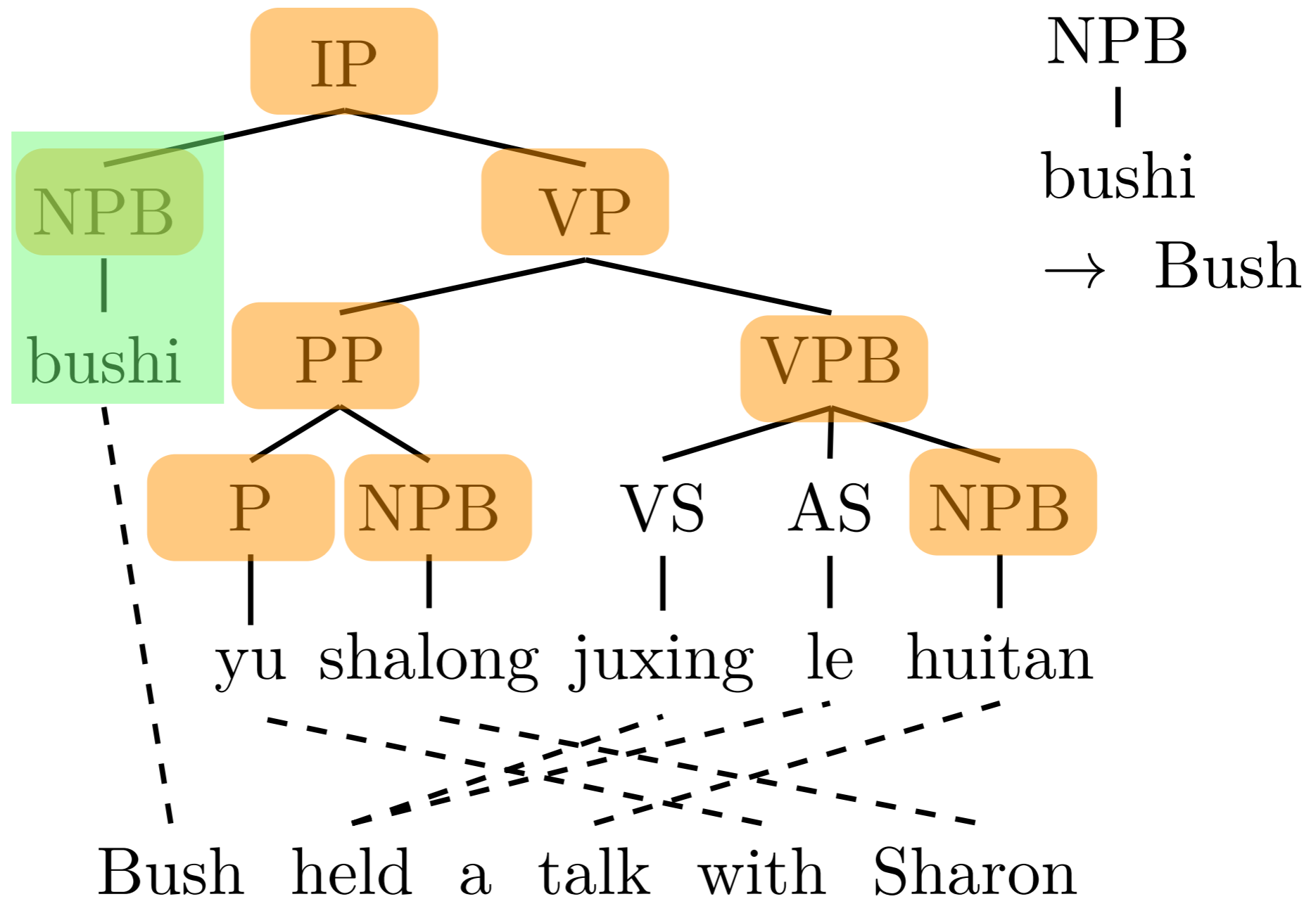
# Rule Extraction



- Extract minimum rules using frontiers

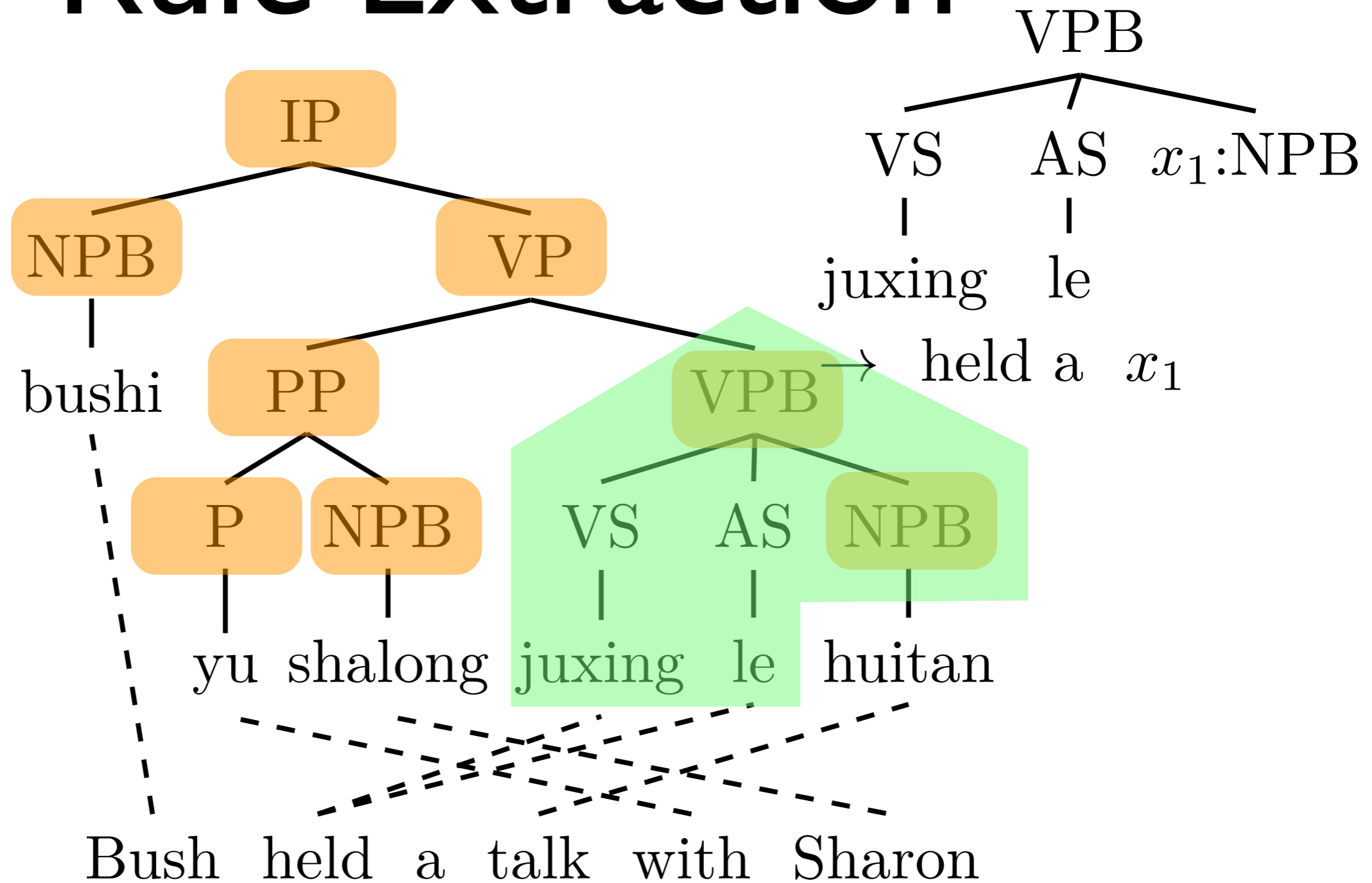


# Rule Extraction



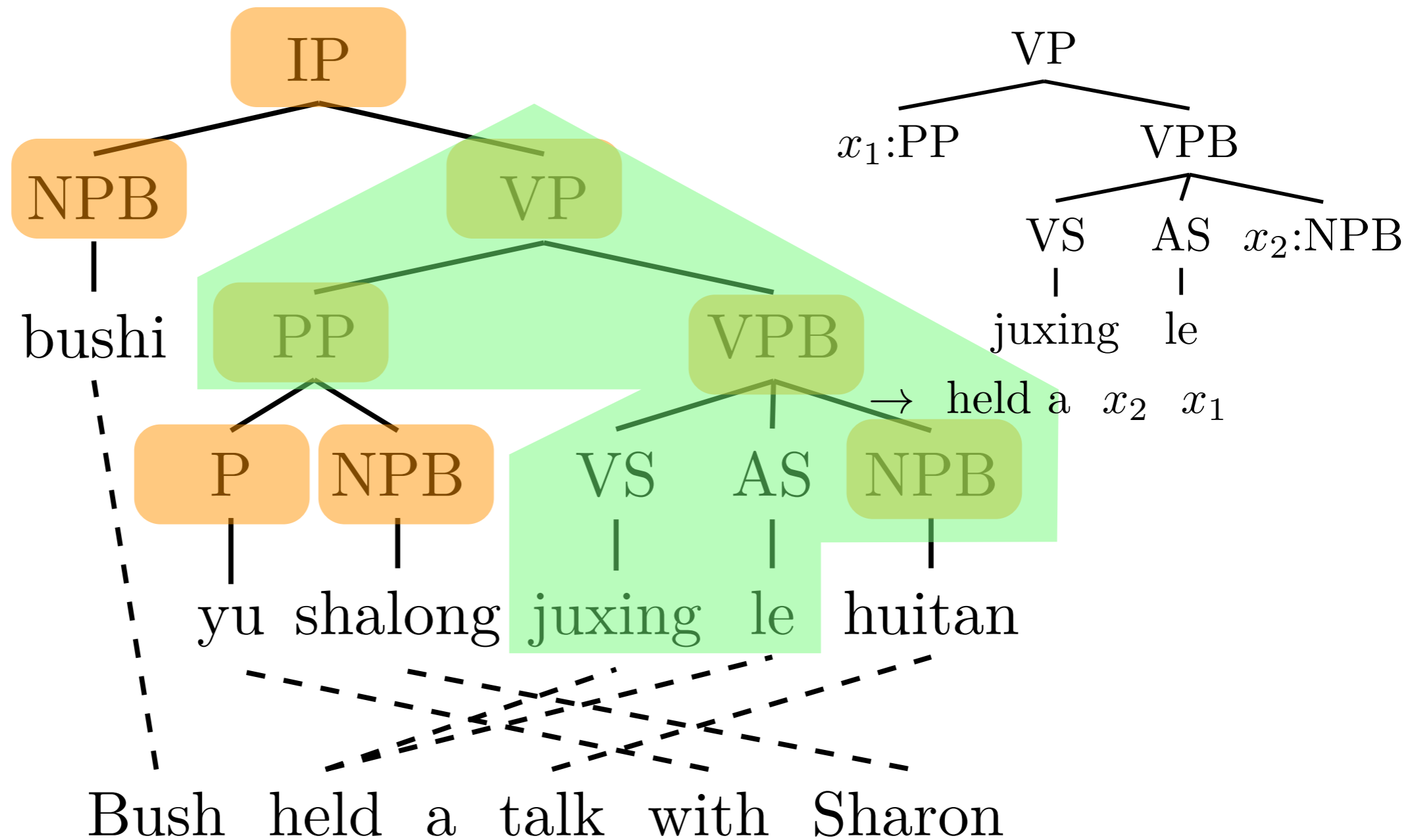
- Extract minimum rules using frontiers

# Rule Extraction



- Extract minimum rules using frontiers

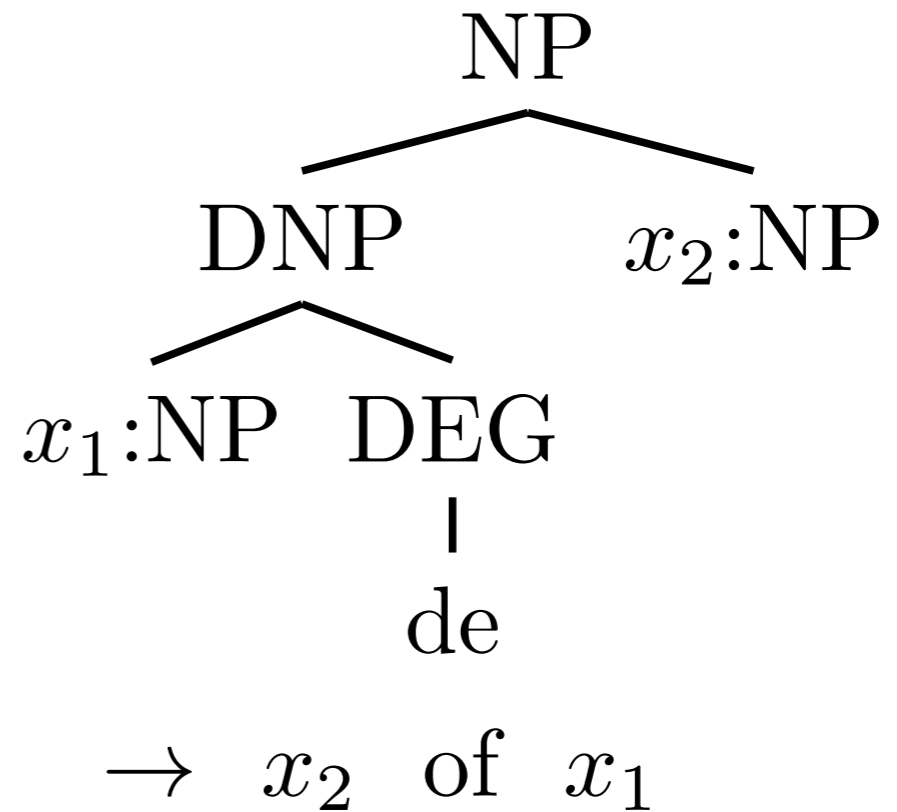
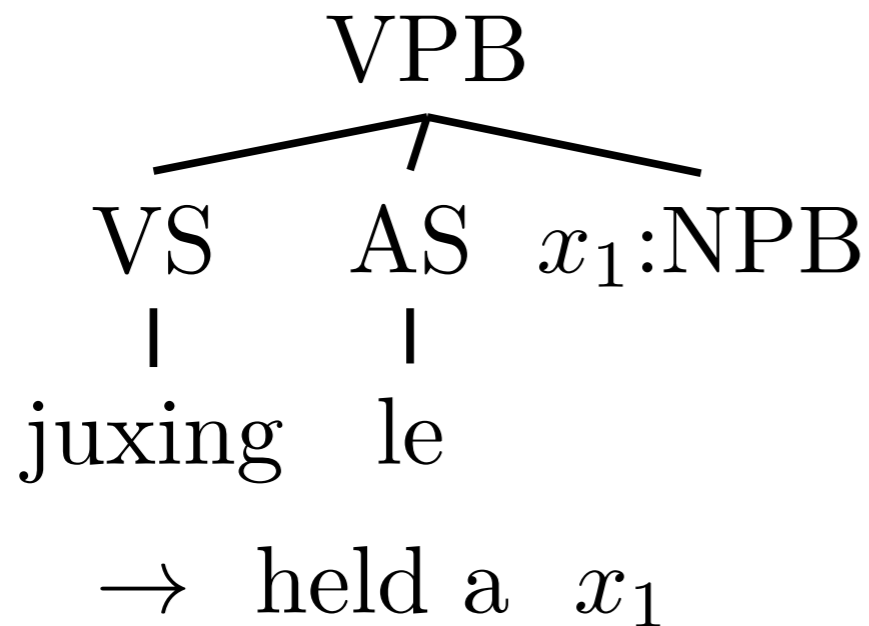
# Rule Extraction



(Galley et al., 2006)

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

# Decoding: String- $\{\text{String, Tree}\}$



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

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

(Galley et al., 2004)

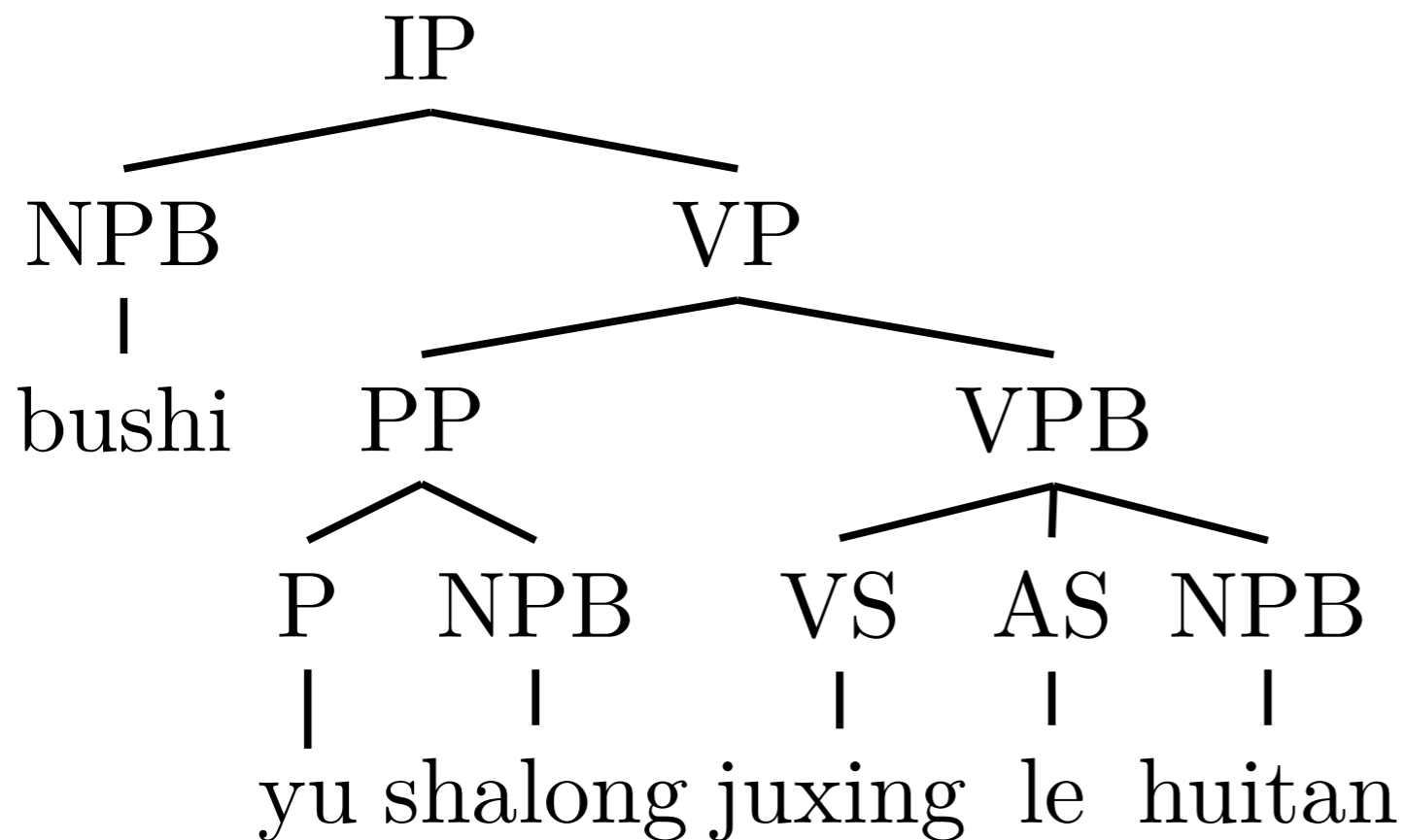
- Similar to SCFG decoding: Use the “collapsed” source side rule to perform CKY parsing
- Construct a translation forest using the target side

# Decoding: Tree- $\{\text{String}, \text{Tree}\}$

(Huang et al., 2006)

- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

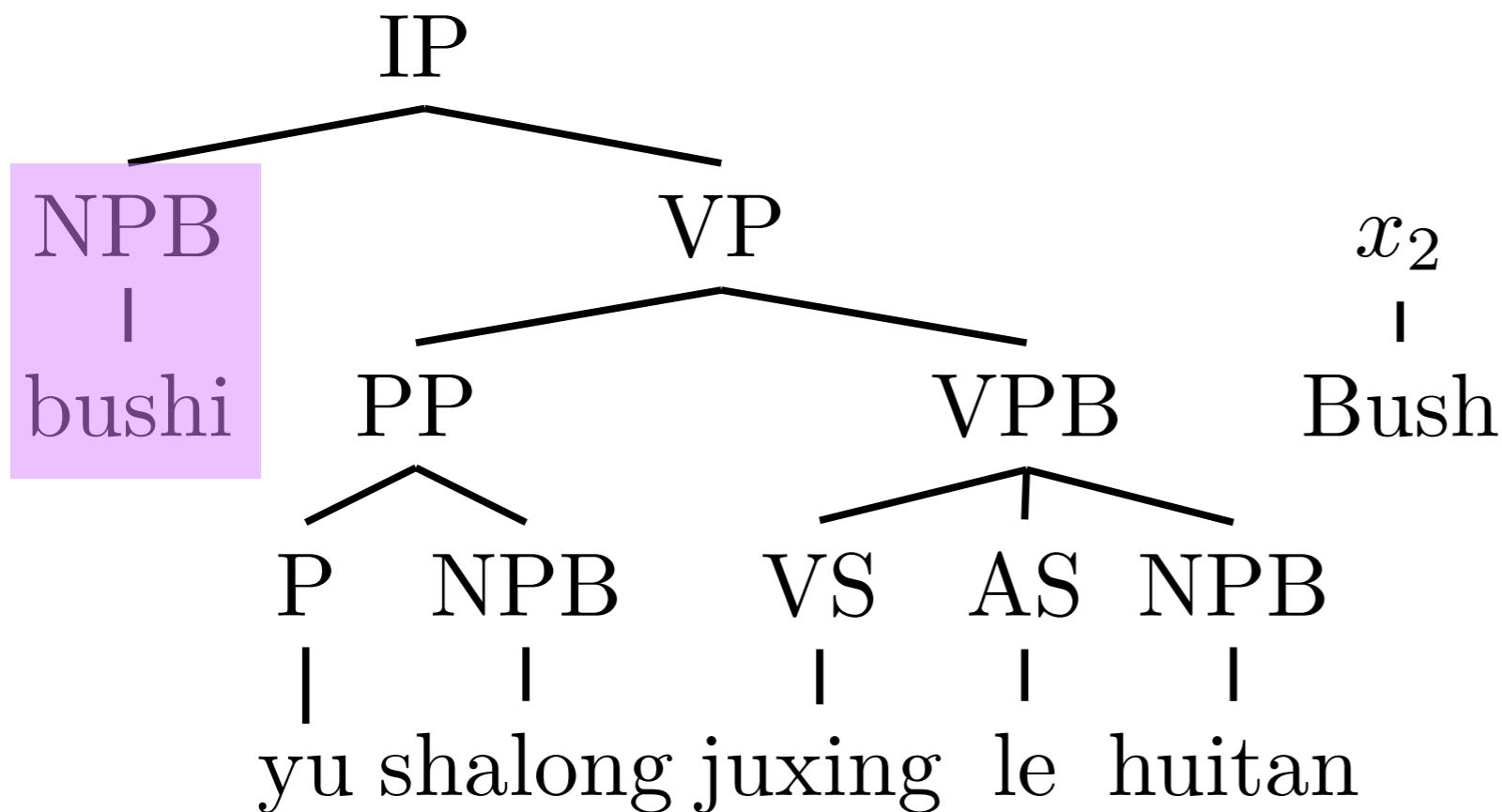
# Decoding: Tree- $\{\text{String}, \text{Tree}\}$



(Huang et al., 2006)

- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

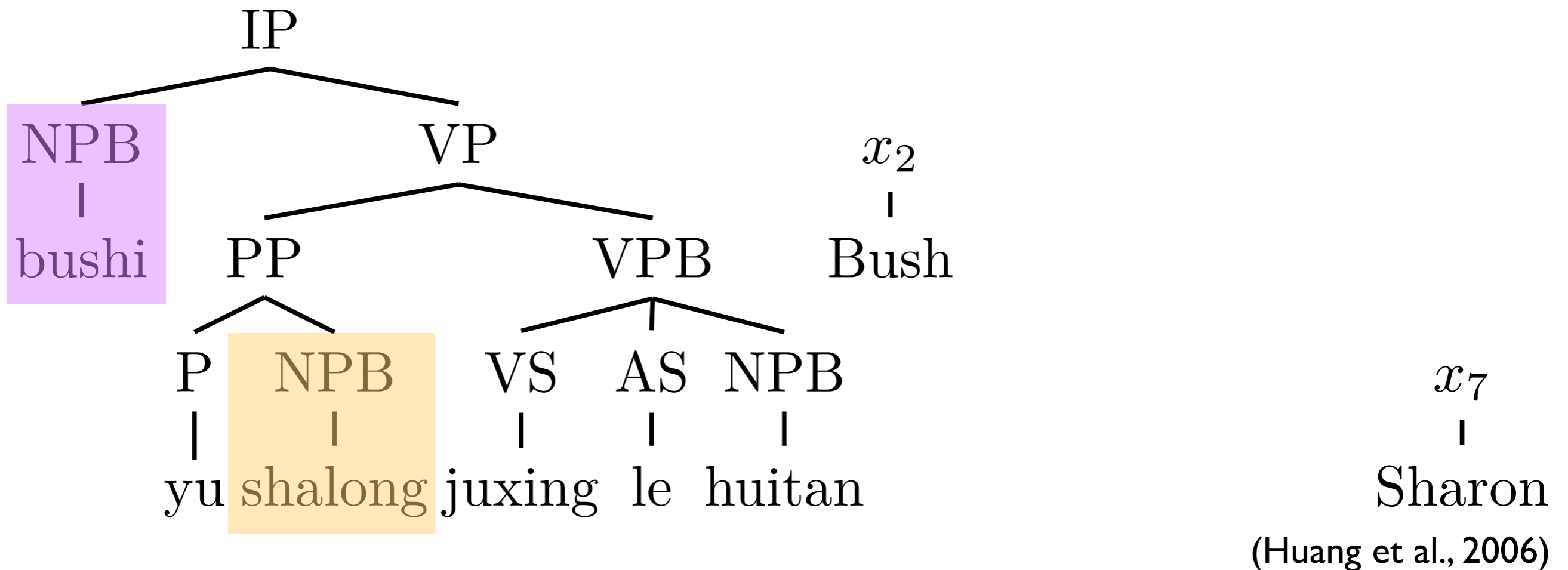
# Decoding: Tree- $\{$ String, Tree $\}$



(Huang et al., 2006)

- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

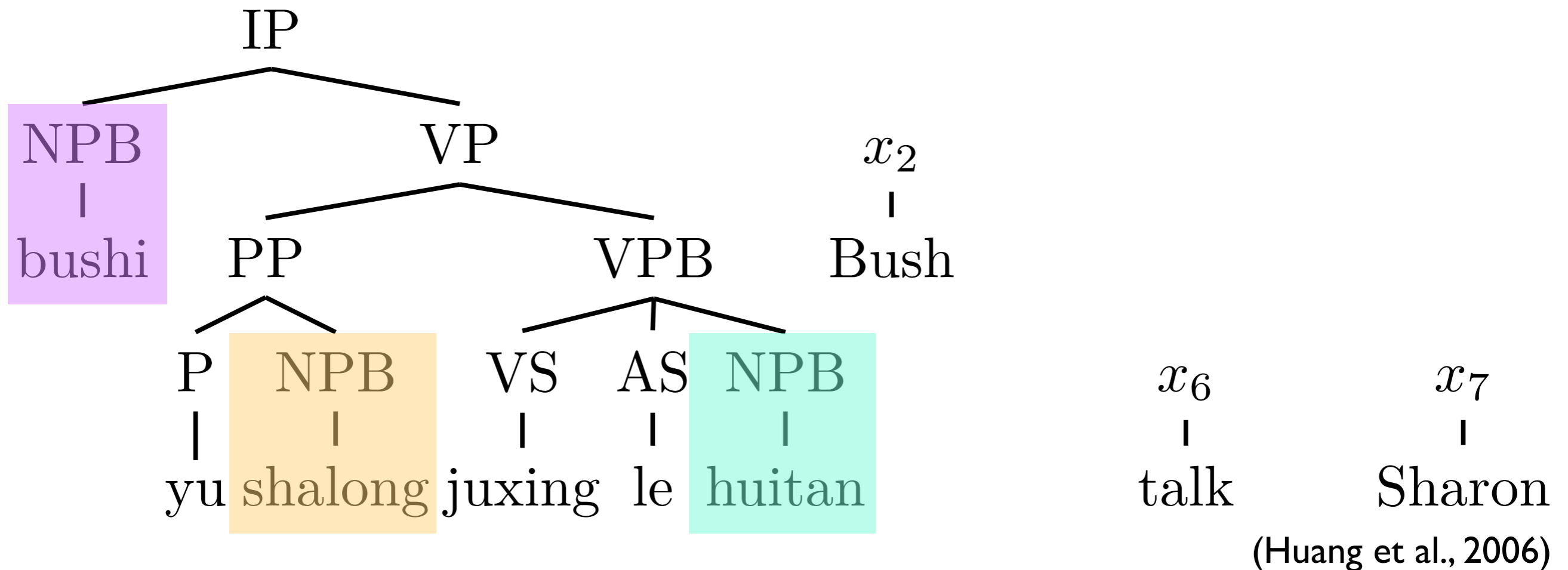
# Decoding: Tree- $\{$ String, Tree $\}$



- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

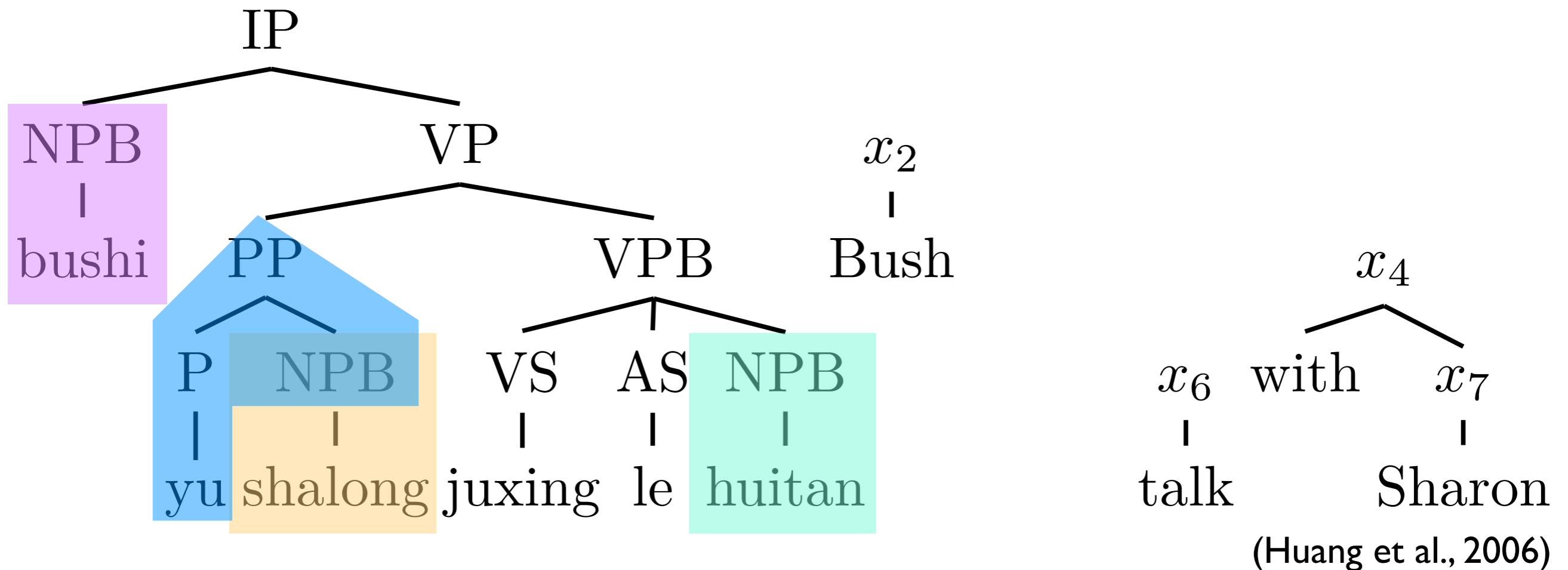


# Decoding: Tree- $\{$ String, Tree $\}$



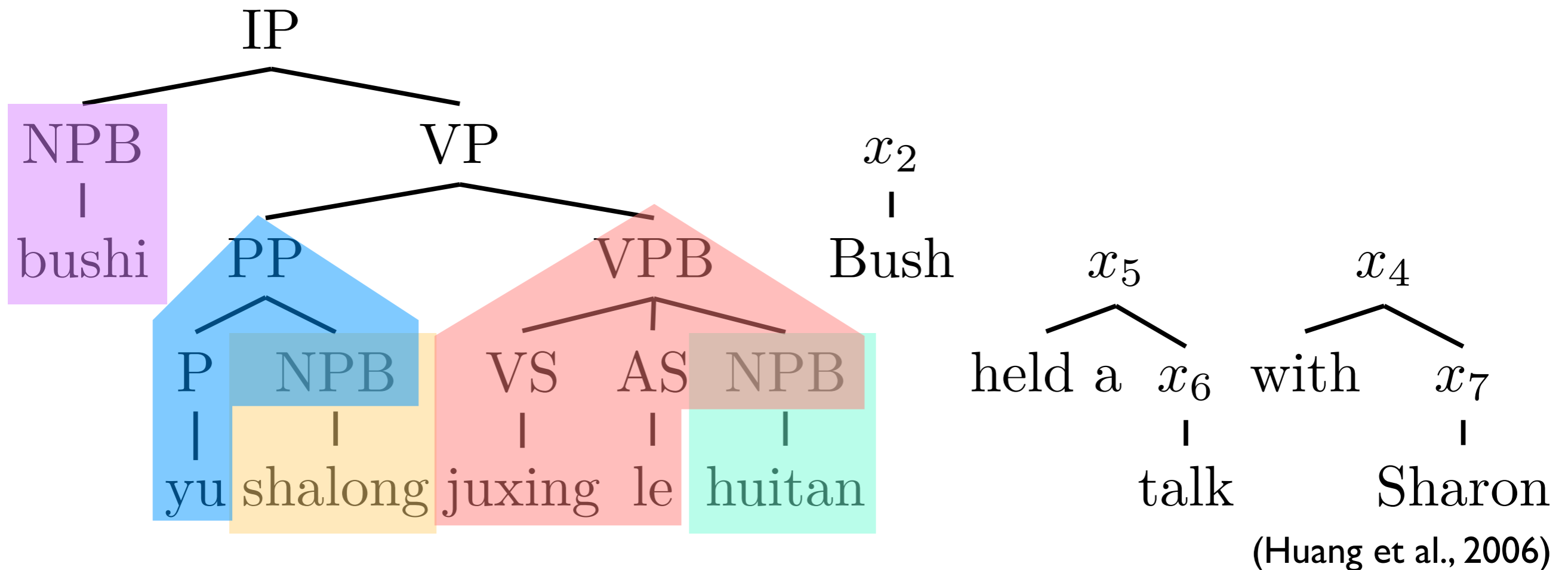
- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

# Decoding: Tree- $\{$ String, Tree $\}$



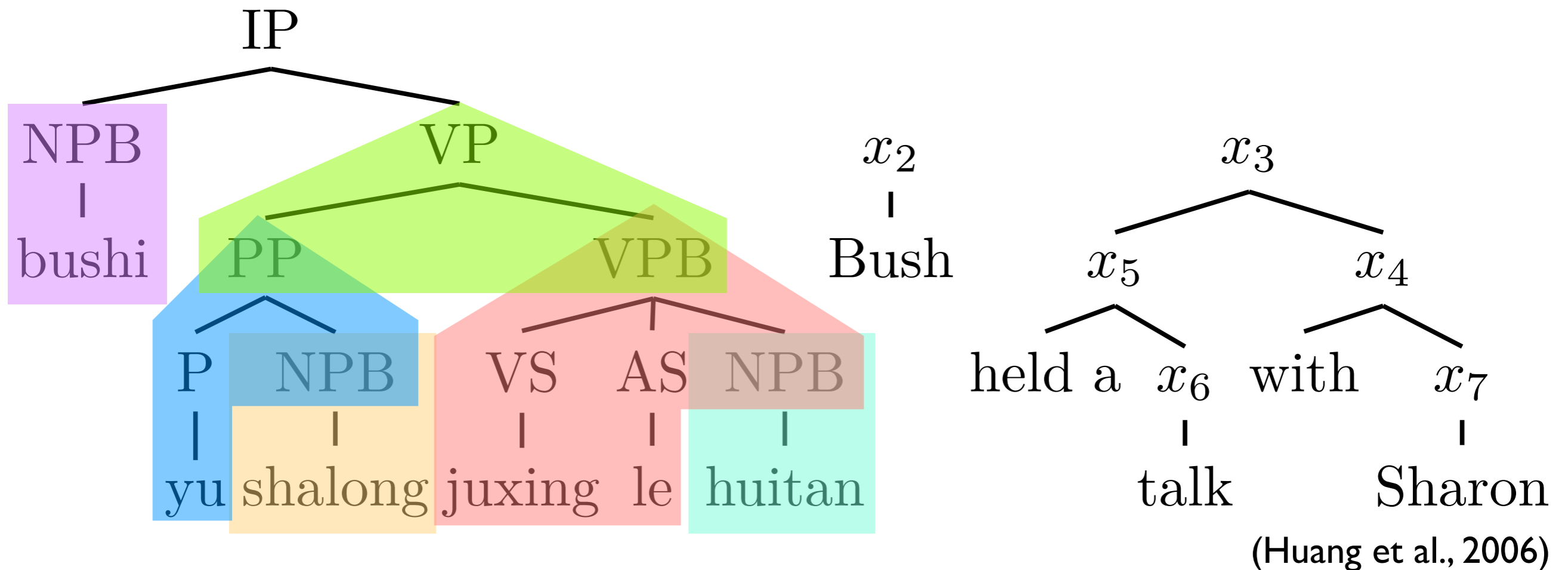
- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

# Decoding: Tree- $\{$ String, Tree $\}$



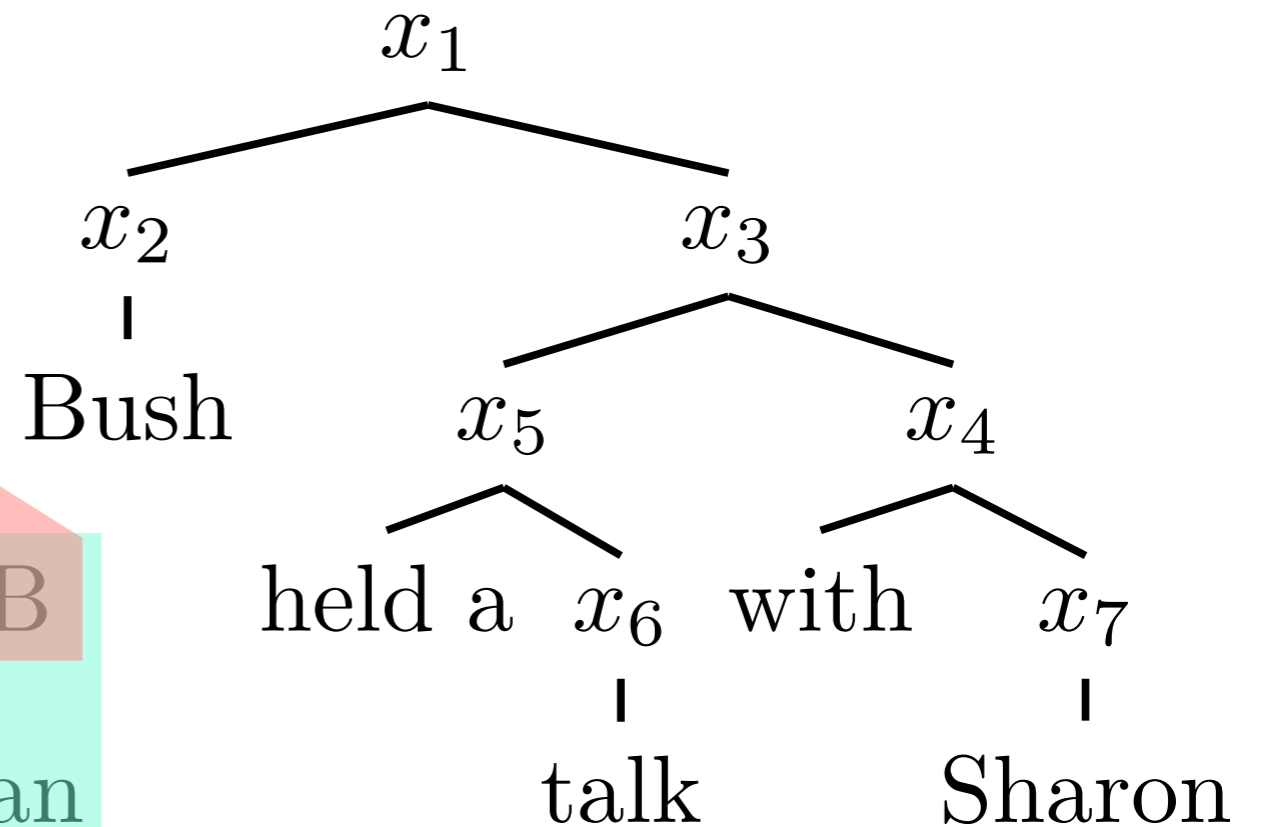
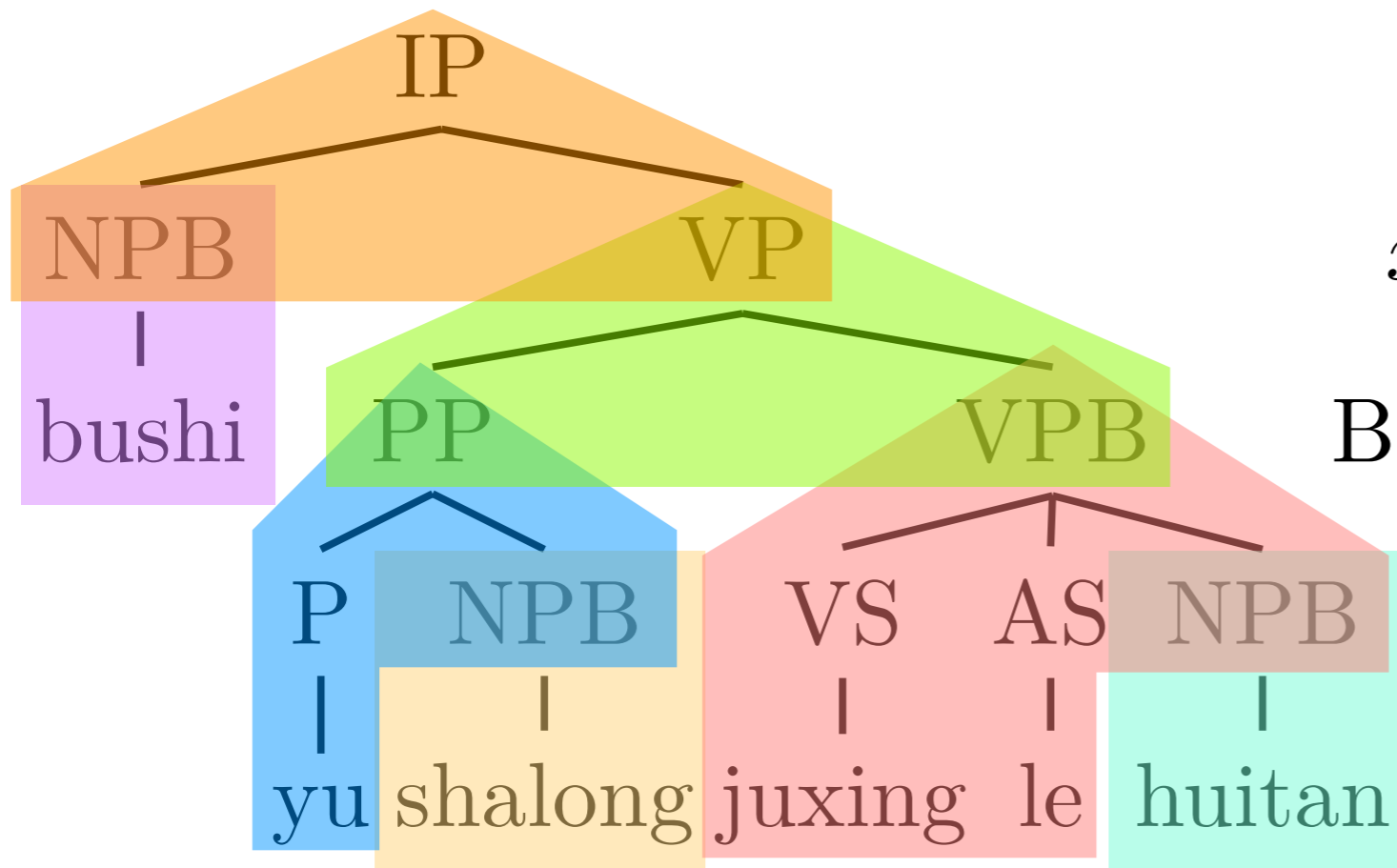
- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

# Decoding: Tree- $\{$ String, Tree $\}$



- First, an input sentence is parsed
- Input tree is transformed into a translation forest by tree rewriting

# 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

# 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)

# State-of-the-art (?)

		BLEU	
		dev	test
English-Chinese	pb	29.7	39.4
	hier	31.7	38.9
	bf2s	31.9	40.7**
English-Czech	<i>wmt best</i>	-	15.4
	pb	14.3	15.5
	hier	14.7	16.0
	bf2s	14.8	16.3*
English-French	<i>wmt best</i>	-	27.6
	pb	24.1	26.1
	hier	23.9	26.1
	bf2s	24.5	26.6**
English-German	<i>wmt best</i>	-	16.3
	pb	14.5	15.5
	hier	14.9	15.9
	bf2s	15.2	16.3**
English-Spanish	<i>wmt best</i>	-	28.4
	pb	24.1	27.9
	hier	24.2	28.4
	bf2s	24.9	28.9**



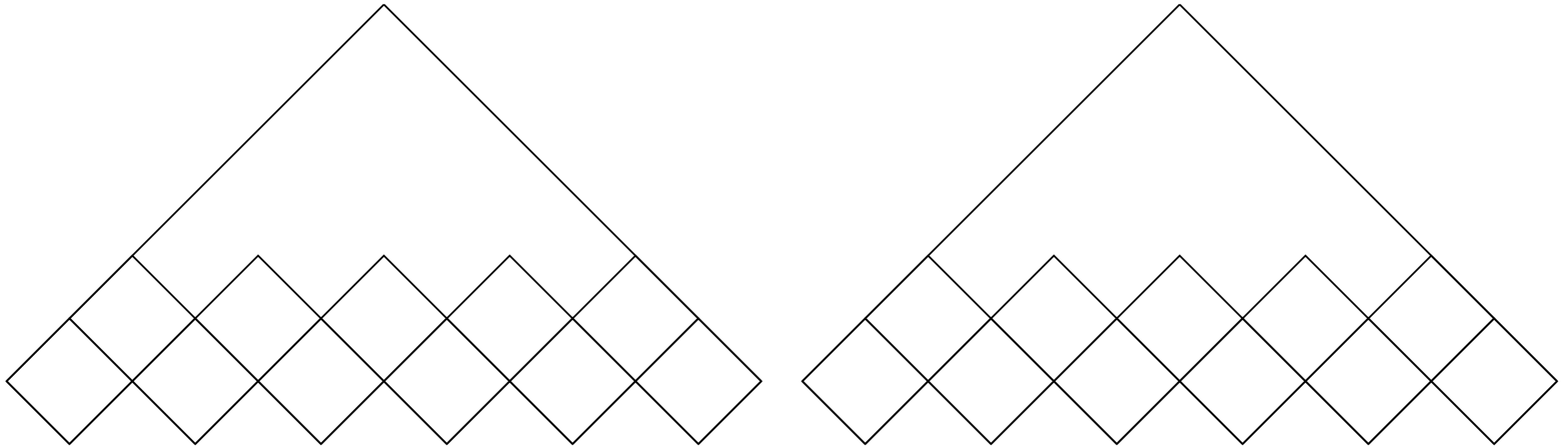
# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing

# Overview

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- **Tree-based SMT**
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - **Bitext parsing**

# Bitext Parsing



- Given bitext (and a synchronous grammar), compute the best paired derivation
- Bitext parsing takes  $O(N^3 M^3)$  for ITG (Wu, 1997)
- For each length  $n$  and  $m$ , for each position  $i$  and  $j$ , for each rule  $X \rightarrow \text{LHS}$ , for each split point  $k$  and  $l$

# ITG

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

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

$$X \rightarrow \langle f, e \rangle$$

- Inversion Transduction Grammar (ITG) (Wu, 1997), an instance of SCFG
- Frequently, we use an abbreviated form

# ITG

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

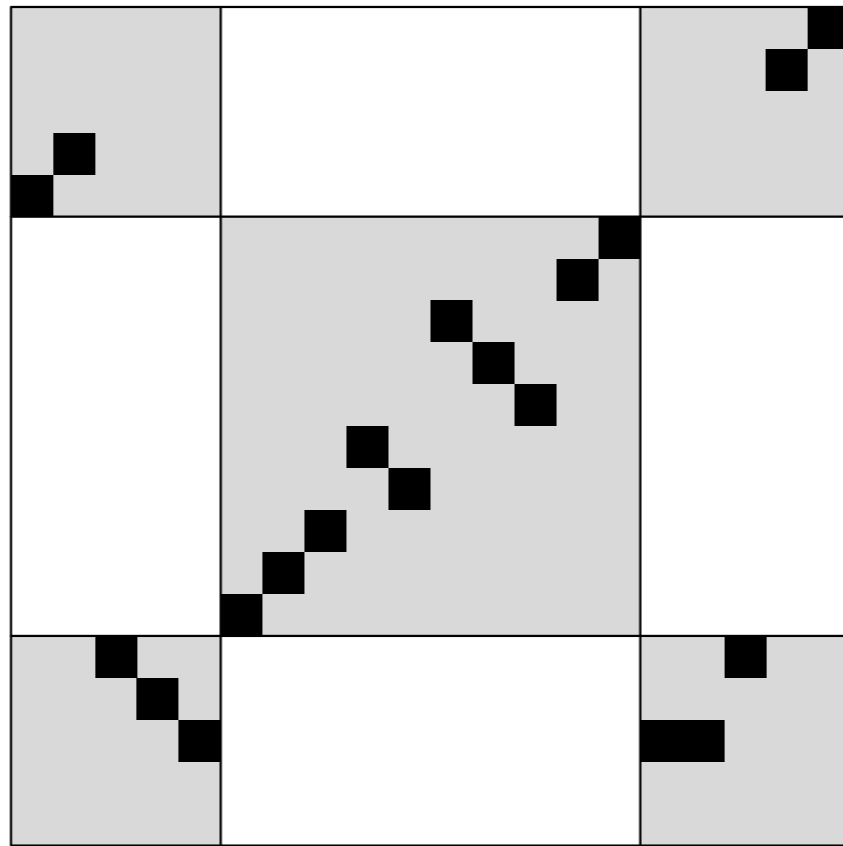
$$X \rightarrow \langle X_{\boxed{1}} X_{\boxed{2}}, X_{\boxed{2}} X_{\boxed{1}} \rangle$$

$$X \rightarrow \langle f, e \rangle$$

$$X \rightarrow [X X] \mid \langle X X \rangle \mid f/e$$

- Inversion Transduction Grammar (ITG) (Wu, 1997), an instance of SCFG
- Frequently, we use an abbreviated form

# Span Pruning

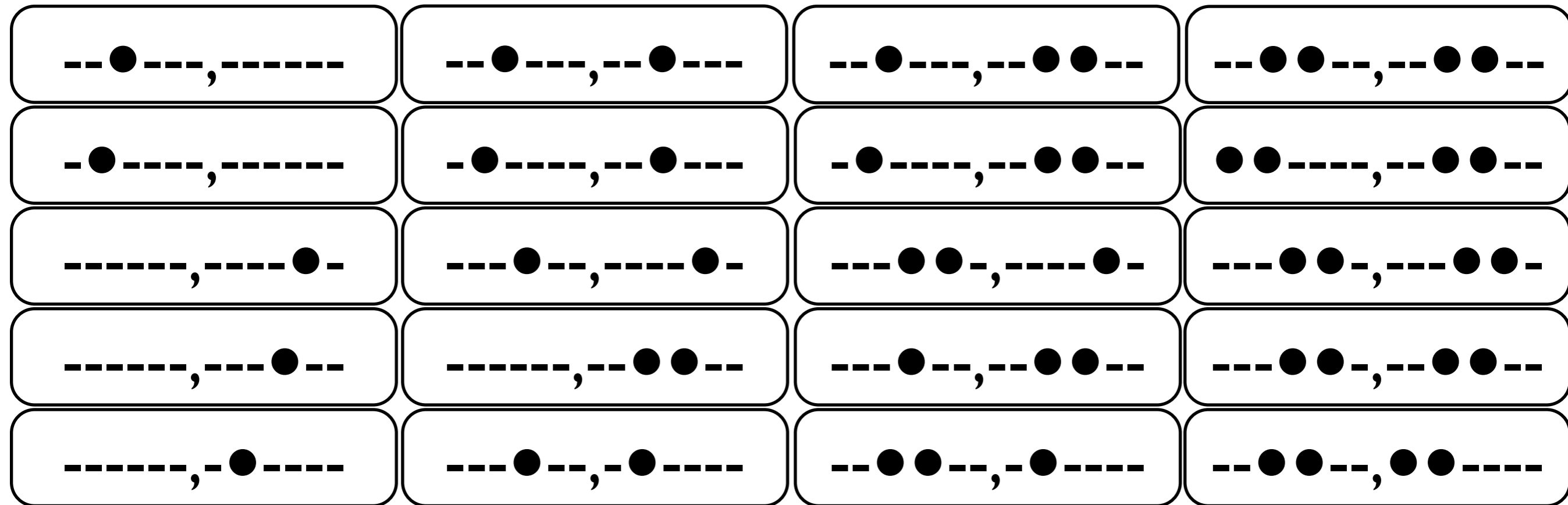


- You do not have to visit all the span pairs
- Use figure-of-merit to prune spans
  - $O(n^4)$  for a naive algorithm (Zhang and Gildea, 2005)
  - $O(n^3)$  for a DP-based algorithm (Zhang et al., 2008)

# Beam Pruning

- Re-organize the search space by the cardinality (Saers et al., 2009)
- Cardinality = # of source/target words parsed
- Prune by the cardinality: Complexity  $O(bn^3)$

# Beam Pruning



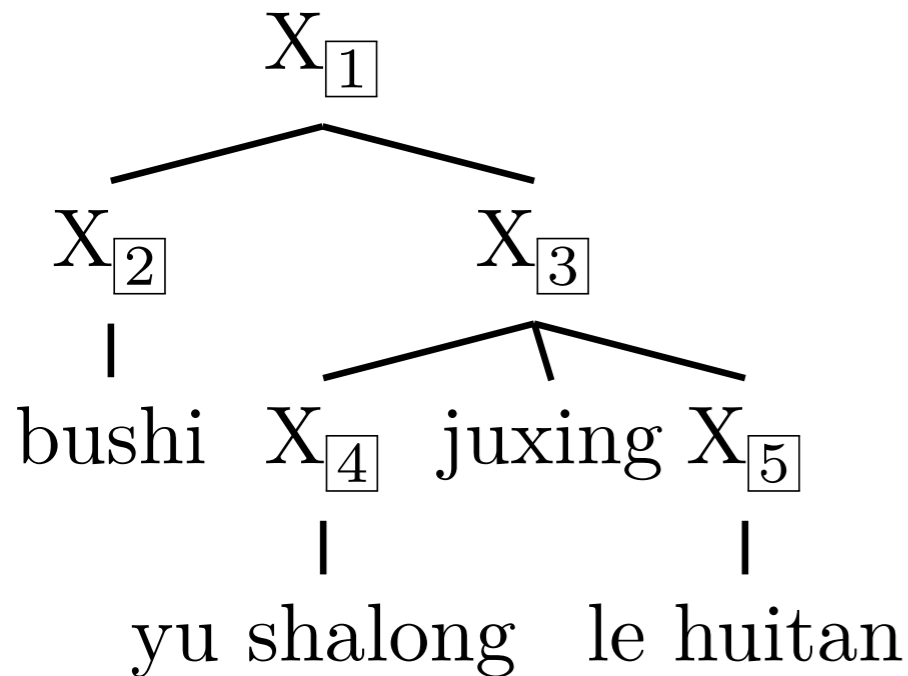
- Re-organize the search space by the cardinality (Saers et al., 2009)
- Cardinality = # of source/target words parsed
- Prune by the cardinality: Complexity  $O(bn^3)$



# Two Parse

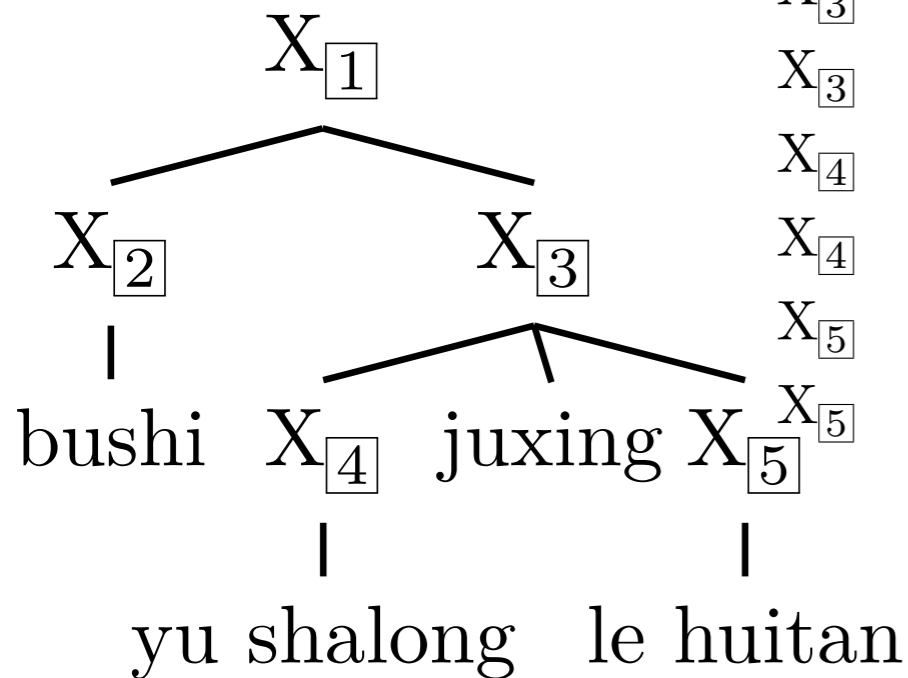
- In practice, we do not have to enumerate all the rules for the Hiero-like grammar (Dyer, 2010)
- First, parsing using the source side rules
- Then, parse the target by the “instantiated” rules

# Two Parse



- In practice, we do not have to enumerate all the rules for the Hiero-like grammar (Dyer, 2010)
- First, parsing using the source side rules
- Then, parse the target by the “instantiated” rules

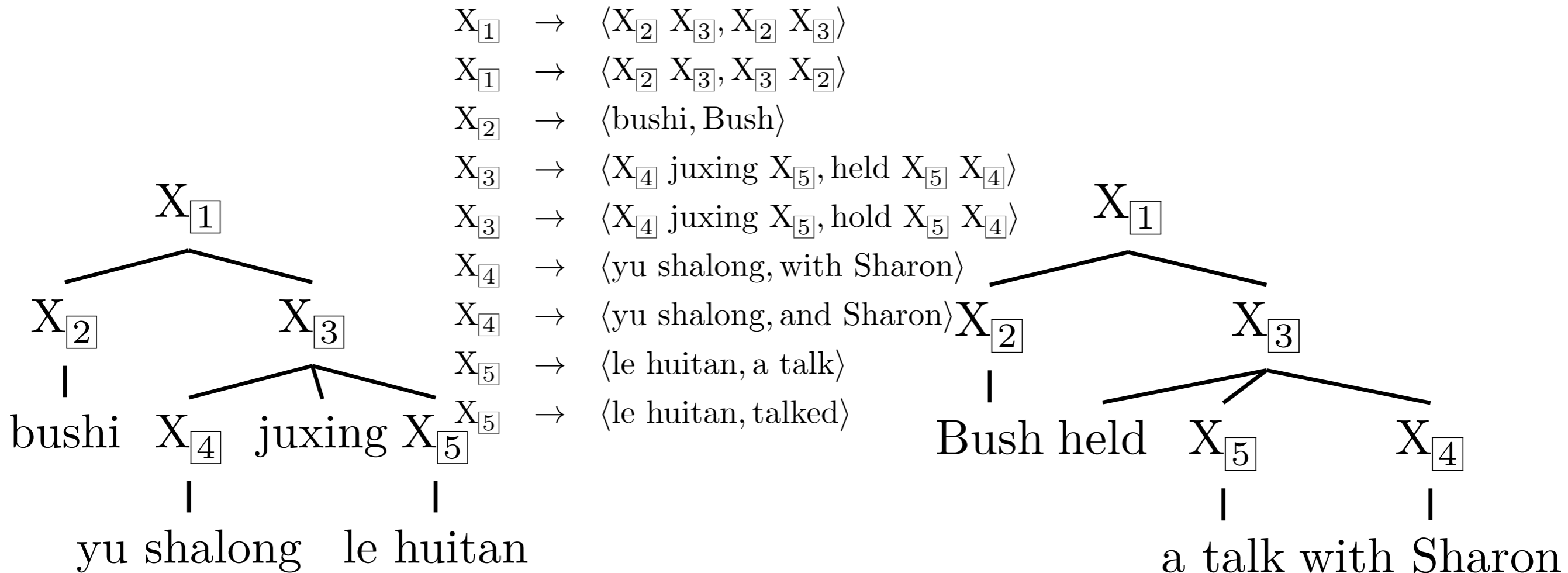
# Two Parse



- $X_1 \rightarrow \langle X_2 X_3, X_2 X_3 \rangle$
- $X_1 \rightarrow \langle X_2 X_3, X_3 X_2 \rangle$
- $X_2 \rightarrow \langle \text{bushi}, \text{Bush} \rangle$
- $X_3 \rightarrow \langle X_4 \text{ juxing } X_5, \text{held } X_5 X_4 \rangle$
- $X_3 \rightarrow \langle X_4 \text{ juxing } X_5, \text{hold } X_5 X_4 \rangle$
- $X_4 \rightarrow \langle \text{yu shalong}, \text{with Sharon} \rangle$
- $X_4 \rightarrow \langle \text{yu shalong}, \text{and Sharon} \rangle$
- $X_5 \rightarrow \langle \text{le huitan}, \text{a talk} \rangle$
- $X_5 \rightarrow \langle \text{le huitan}, \text{talked} \rangle$

- In practice, we do not have to enumerate all the rules for the Hiero-like grammar (Dyer, 2010)
- First, parsing using the source side rules
- Then, parse the target by the “instantiated” rules

# Two Parse



- In practice, we do not have to enumerate all the rules for the Hiero-like grammar (Dyer, 2010)

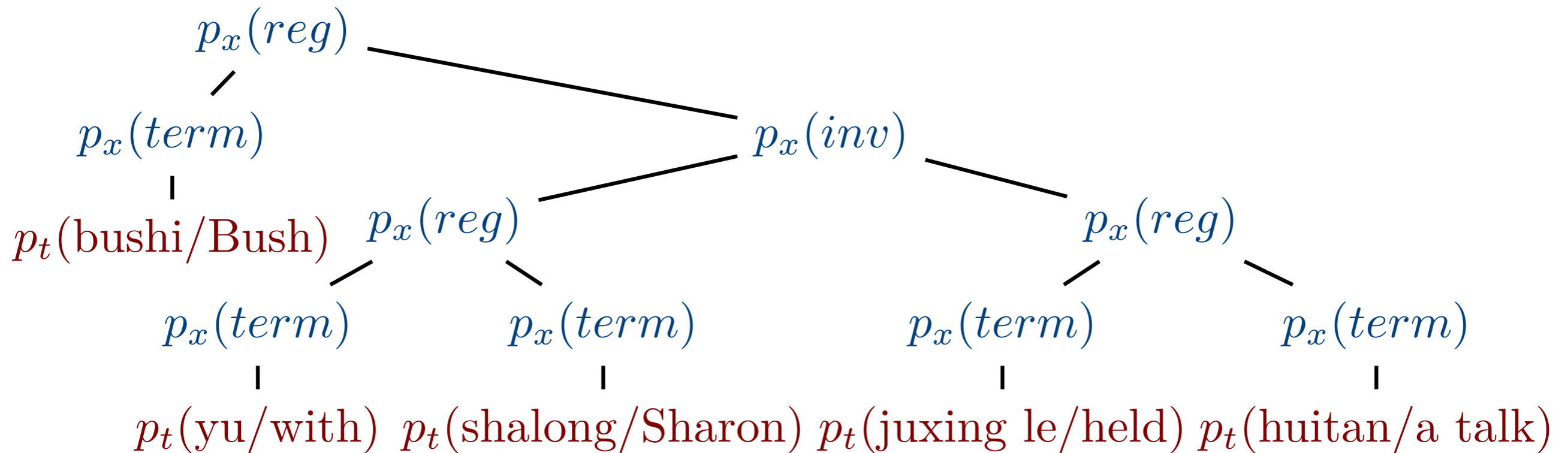
- First, parsing using the source side rules

- Then, parse the target by the “instantiated” rules

# ITG for Phrase Induction

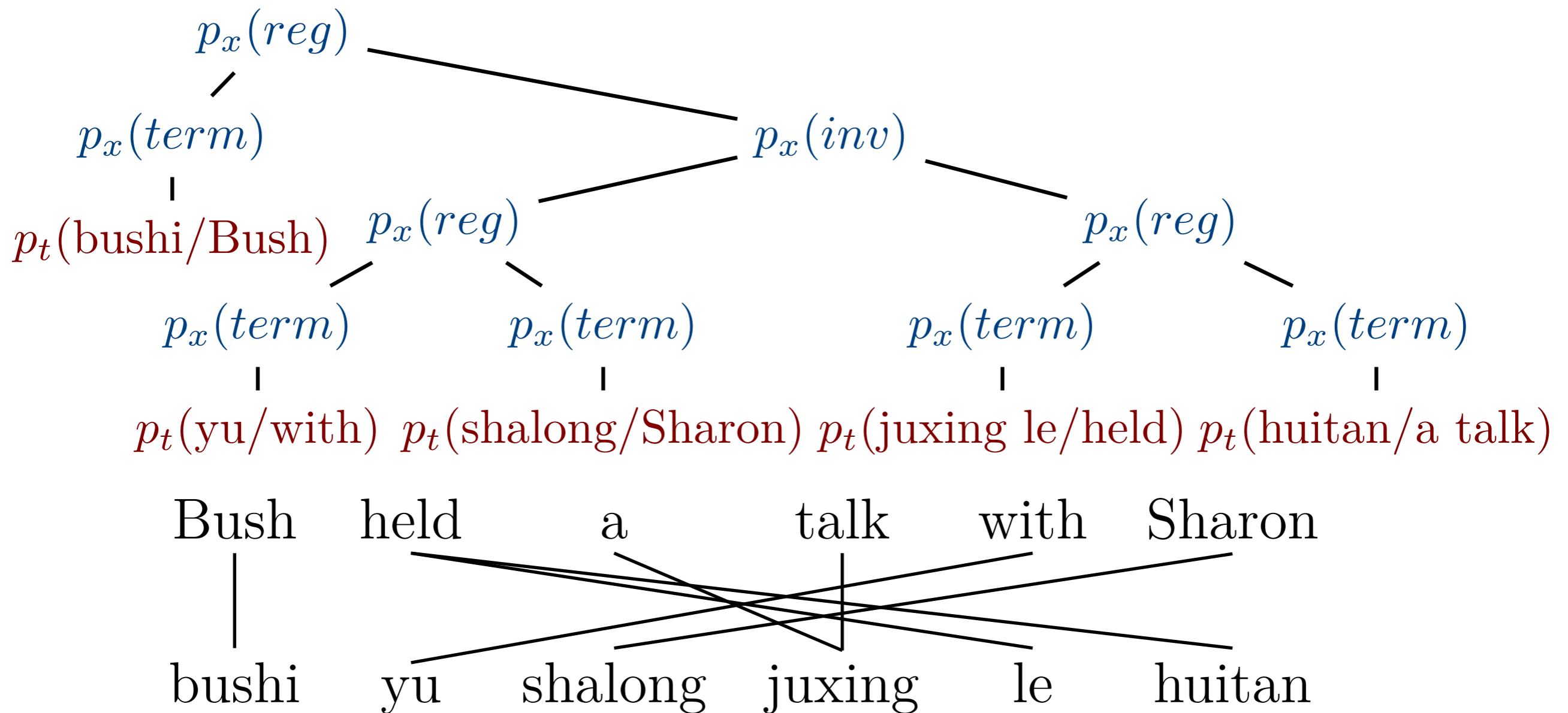
- Phrasal alignment by ITG (Cherry and Lin, 2007)

# ITG for Phrase Induction



- Phrasal alignment by ITG (Cherry and Lin, 2007)

# ITG for Phrase Induction



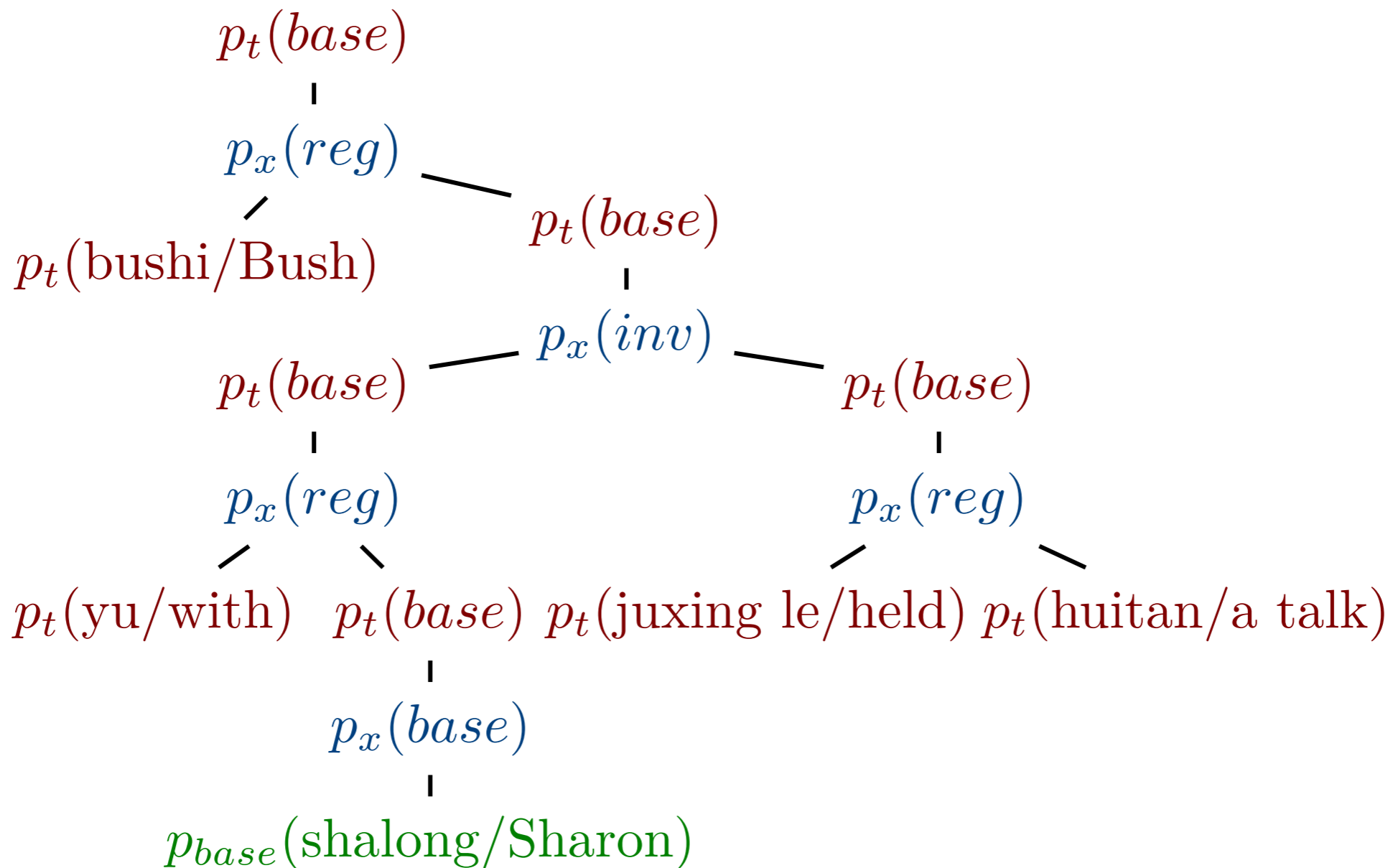
- Phrasal alignment by ITG (Cherry and Lin, 2007)

# Exhaustive ITG Phrases

- Recursively divide-and-conquer
- Multiple granularities are included in the model (Neubig et al., 2011)



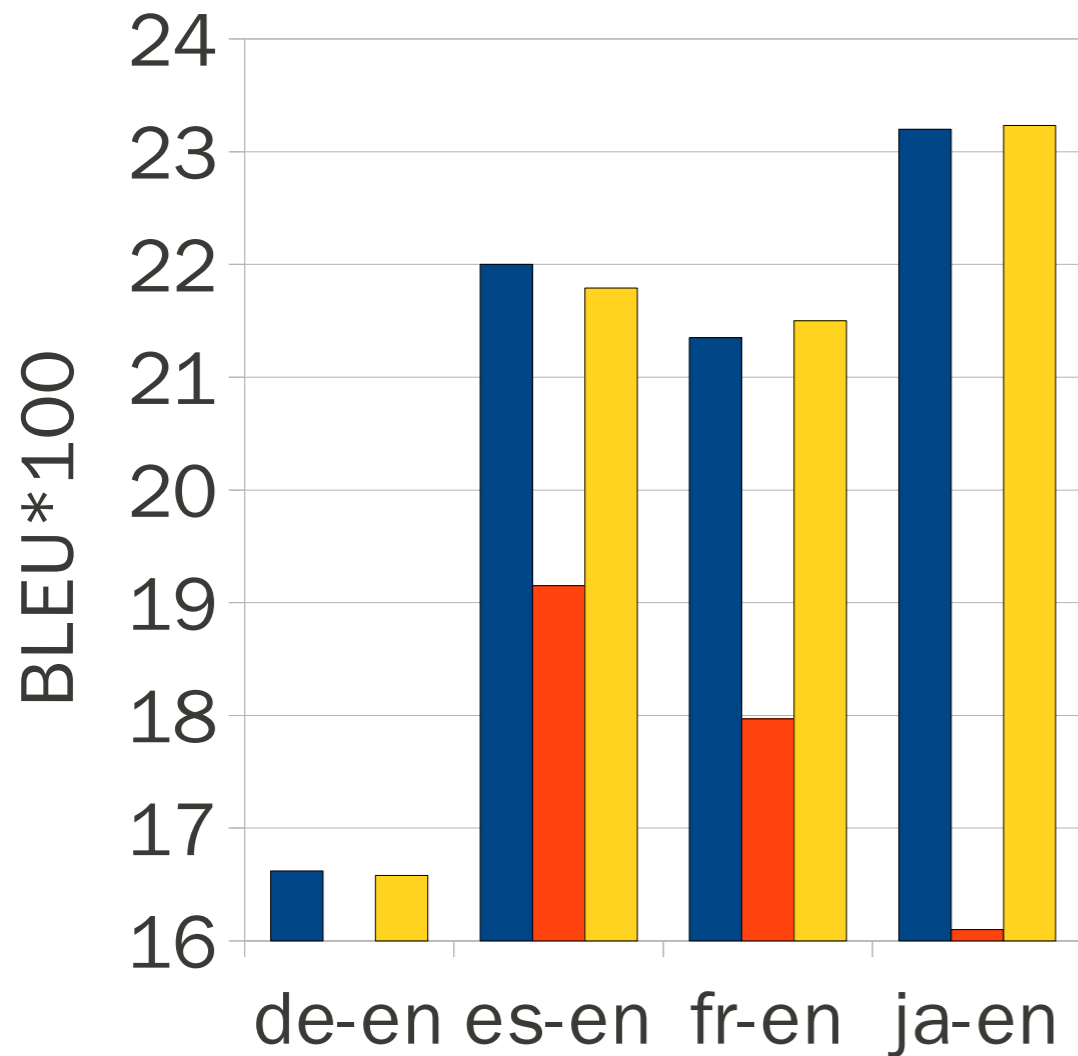
# Exhaustive ITG Phrases



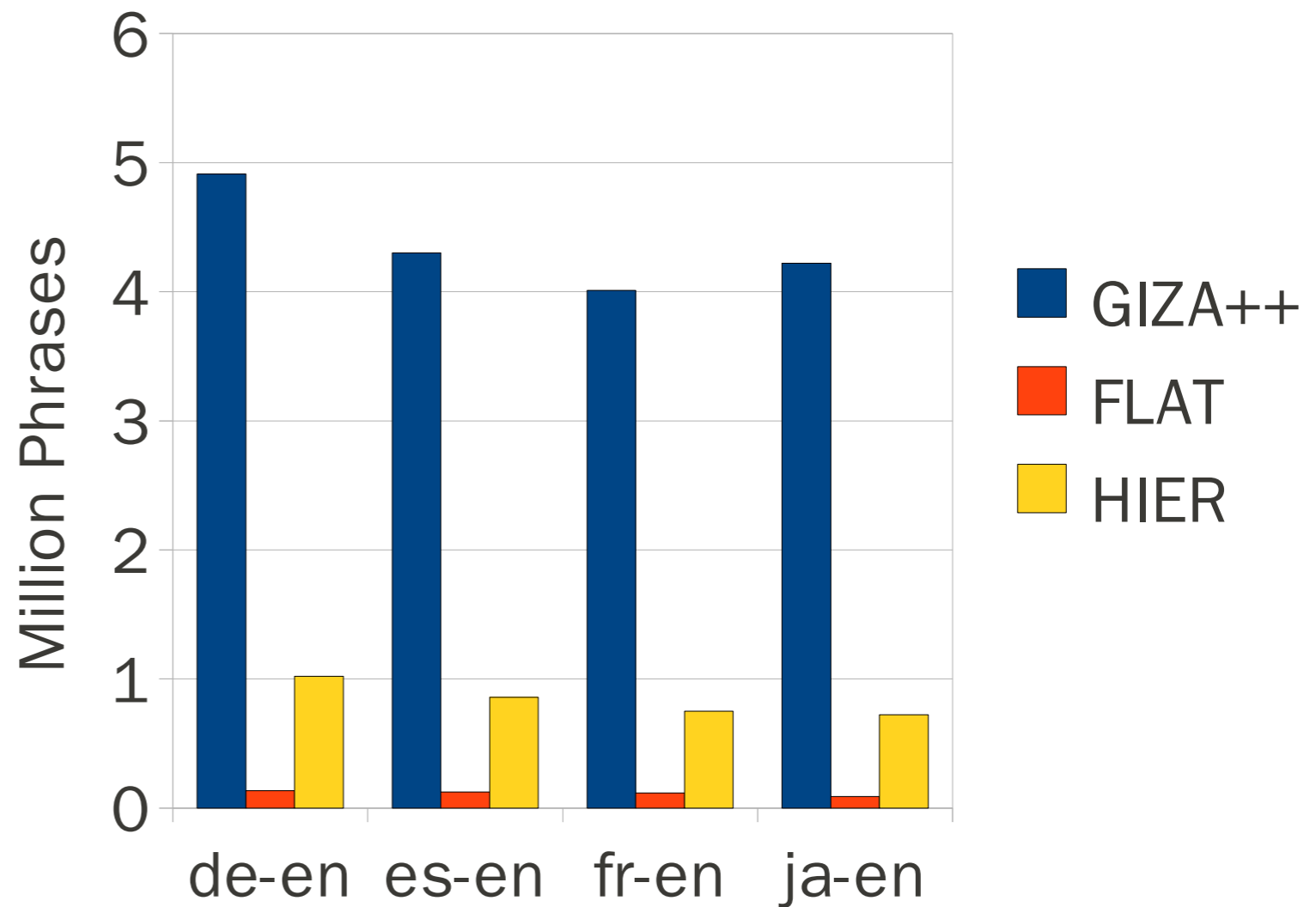
- Recursively divide-and-conquer
- Multiple granularities are included in the model (Neubig et al., 2011)

# Compact Model

Translation Accuracy



Phrase Table Size



(Neubig et al., 2011)

# Conclusion

- Bitext parsing with ITG
- Span pruning, beam pruning, two-parse
- ITG is a simple, yet powerful grammar for bilingual knowledge induction

# Summary

- Backgrounds
  - CFG, parsing, hypergraph, deductive system semirings
- Tree-based SMT
  - Synchronous-CFG
  - String-to-Tree, Tree-to-String
  - Bitext parsing

# Implementations

- synchronous-CFG
  - Cdec: <http://cdec-decoder.org>
  - Jane: <http://www-i6.informatik.rwth-aachen.de/jane/>
  - Joshua: <http://joshua.sourceforge.net>
  - Moses: <http://www.statmt.org/moses/>
- {Tree,String}-to-{Tree,String}: Usually, closed source, like my private implementation:-)
- ITG-model induction
  - Pialign: <http://www.phontron.com/pialign/>

# References

- P. F. Brown, J. Cocke, S. D. Pietra, V. J. D. Pietra, F. Jelinek, J. D. Lafferty, R. L. Mercer, and P. S. Roossin, "A statistical approach to machine translation," *Computational Linguistics*, vol. 16, no. 2, pp. 79--85, 1990.
- P. Koehn, F. J. Och, and D. Marcu, "Statistical phrase-based translation," in *Proc. of HLT-NAACL 2003*, (Edmonton), pp. 48--54, May-June 2003.
- D. Chiang, "Hierarchical phrase-based translation," *Comput. Linguist.*, vol. 33, no. 2, pp. 201--228, 2007.
- M. Galley, M. Hopkins, K. Knight, and D. Marcu, "What's in a translation rule?," in *HLT-NAACL 2004: Main Proceedings* (D. M. Susan Dumais and S. Roukos, eds.), (Boston, Massachusetts, USA), pp. 273--280, Association for Computational Linguistics, May 2 - May 7 2004.
- A. Zollmann and A. Venugopal, "Syntax augmented machine translation via chart parsing," in *StatMT '06: Proceedings of the Workshop on Statistical Machine Translation*, (Morristown, NJ, USA), pp. 138--141, Association for Computational Linguistics, 2006.
- L. Huang, K. Knight, and A. Joshi, "Statistical syntax-directed translation with extended domain of locality," in *In Proc. AMTA 2006*, pp. 66--73, 2006.

# References

- Y. Liu, Q. Liu, and S. Lin, "Tree-to-string alignment template for statistical machine translation," in *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, (Sydney, Australia), pp. 609--616, Association for Computational Linguistics, July 2006.
- C. Quirk, A. Menezes, and C. Cherry, "Dependency treelet translation: syntactically informed phrasal smt," in *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, (Morristown, NJ, USA), pp. 271--279, Association for Computational Linguistics, 2005.
- L. Shen, J. Xu, and R. Weischedel, "A new string-to-dependency machine translation algorithm with a target dependency language model," in *Proceedings of ACL-08: HLT*, (Columbus, Ohio), pp. 577--585, Association for Computational Linguistics, June 2008.
- Y. Ding and M. Palmer, "Machine translation using probabilistic synchronous dependency insertion grammars," in *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, (Morristown, NJ, USA), pp. 541--548, Association for Computational Linguistics, 2005.
- Y. Liu, Y. Lu, and Q. Liu, "Improving tree-to-tree translation with packed forests," in *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, (Suntec, Singapore), pp. 558--566, Association for Computational Linguistics, August 2009.
- D. Klein and C. D. Manning, "Parsing and hypergraphs," in *In IWPT*, pp. 123--134, 2001.

# References

- S. M. Shieber, Y. Schabes, and O. C. N. Pereira, "Principles and implementation of deductive parsing," *Journal of Logic Programming*, 1995.
- L. Huang and D. Chiang, "Better k-best parsing," in *Proceedings of the Ninth International Workshop on Parsing Technology*, (Vancouver, British Columbia), pp. 53--64, Association for Computational Linguistics, October 2005.
- L. Huang and D. Chiang, "Forest rescoring: Faster decoding with integrated language models," in *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, (Prague, Czech Republic), pp. 144--151, Association for Computational Linguistics, June 2007.
- A. Gesmundo and J. Henderson, "Faster Cube Pruning," in *Proceedings of the seventh International Workshop on Spoken Language Translation (IWSLT)* (M. Federico, I. Lane, M. Paul, and F. Yvon, eds.), pp. 267--274, 2010.
- M. Galley, J. Graehl, K. Knight, D. Marcu, S. DeNeefe, W. Wang, and I. Thayer, "Scalable inference and training of context-rich syntactic translation models," in *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, (Sydney, Australia), pp. 961-968, Association for Computational Linguistics, July 2006.
- H. Mi and L. Huang, "Forest-based translation rule extraction," in *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, (Honolulu, Hawaii), pp. 206--214, Association for Computational Linguistics, October 2008.
- H. Mi, L. Huang, and Q. Liu, "Forest-based translation," in *Proceedings of ACL-08: HLT*, (Columbus, Ohio), pp. 192--199, Association for Computational Linguistics, June 2008.



# References

- H. Zhang, L. Fang, P. Xu, and X. Wu, "Binarized forest to string translation," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, (Portland, Oregon, USA), pp. 835--845, Association for Computational Linguistics, June 2011.
- D. Chiang, "Learning to translate with source and target syntax," in *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, (Uppsala, Sweden), pp. 1443--1452, Association for Computational Linguistics, July 2010.
- J. Xie, H. Mi, and Q. Liu, "A novel dependency-to-string model for statistical machine translation," in *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, (Edinburgh, Scotland, UK.), pp. 216-226, Association for Computational Linguistics, July 2011.
- D. Wu, "Stochastic inversion transduction grammars and bilingual parsing of parallel corpora," *Comput. Linguist.*, vol. 23, no. 3, pp. 377--403, 1997.
- H. Zhang and D. Gildea, "Stochastic lexicalized inversion transduction grammar for alignment," in *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, ACL '05, (Stroudsburg, PA, USA), pp. 475--482, Association for Computational Linguistics, 2005.
- H. Zhang, C. Quirk, R. C. Moore, and D. Gildea, "Bayesian learning of non-compositional phrases with synchronous parsing," in *Proceedings of ACL-08: HLT*, (Columbus, Ohio), pp. 97--105, Association for Computational Linguistics, June 2008.

# References

- M. Saers, J. Nivre, and D. Wu, ``Learning stochastic bracketing inversion transduction grammars with a cubic time biparsing algorithm," in *Proceedings of the 11th International Conference on Parsing Technologies (IWPT'09)*, (Paris, France), pp. 29--32, Association for Computational Linguistics, October 2009.
- C. Dyer, ``Two monolingual parses are better than one (synchronous parse)," in *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, (Los Angeles, California), pp. 263--266, Association for Computational Linguistics, June 2010.
- C. Cherry and D. Lin, ``Inversion transduction grammar for joint phrasal translation modeling," in *Proceedings of SSST, NAACL-HLT 2007 / AMTA Workshop on Syntax and Structure in Statistical Translation*, (Rochester, New York), pp. 17--24, Association for Computational Linguistics, April 2007.
- G. Neubig, T. Watanabe, E. Sumita, S. Mori, and T. Kawahara, ``An unsupervised model for joint phrase alignment and extraction," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, (Portland, Oregon, USA), pp. 632-641, Association for Computational Linguistics, June 2011.