

Tree-based and Forest-Based Translation



Liang Huang



Joint work with Kevin Knight (ISI), Aravind Joshi (Penn), Haitao Mi and Qun Liu (ICT)

UC Berkeley, Feb 6, 2009

Translation is hard!



zi zhu zhong duan
自 助 终 端

self help terminal device

help oneself terminating machine

(ATM, “self-service terminal”)

Translation is hard!



Translation is hard!



Translation is hard!



Translation is hard!



or even...



or even...



clear evidence that MT is used in real life.

How do people translate?

1. understand the source language sentence
2. generate the target language translation

布什 与 沙龙 举行 了 会谈

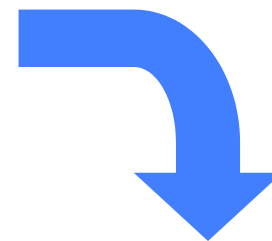
Bùshí yu Shalóng jùxíng le huìtán

Bush and/
with Sharon hold [*past.*] meeting

How do people translate?

1. understand the source language sentence
2. generate the target language translation

布什 与 沙龙 举行 了 会谈
Bùshí yu Shalóng jùxíng le huìtán
Bush and/ with Sharon hold [*past.*] meeting



How do people translate?

1. understand the source language sentence
2. generate the target language translation

布什 与 沙龙 举行 了 会谈
Bùshí yu Shalóng jùxíng le huìtán
Bush and/ with Sharon hold [past.] meeting



“Bush held a meeting with Sharon”

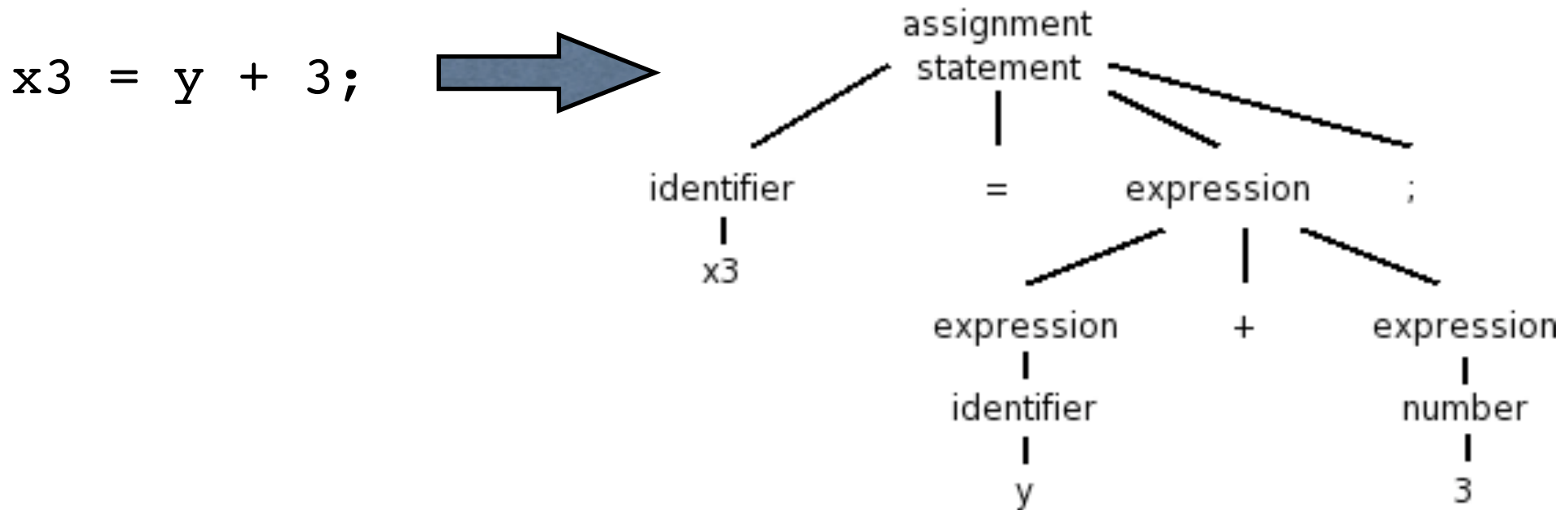
How do compilers translate?

1. parse high-level language program into a syntax tree
2. generate intermediate or machine code accordingly

```
x3 = y + 3;
```

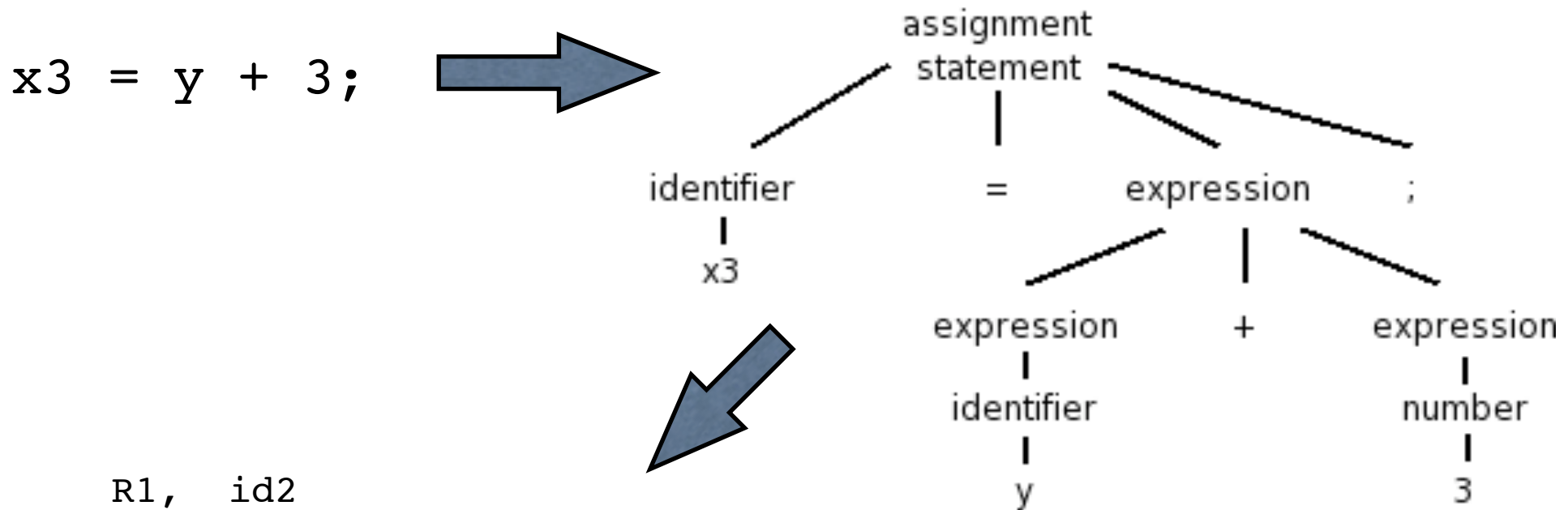

How do compilers translate?

1. parse high-level language program into a syntax tree
2. generate intermediate or machine code accordingly



How do compilers translate?

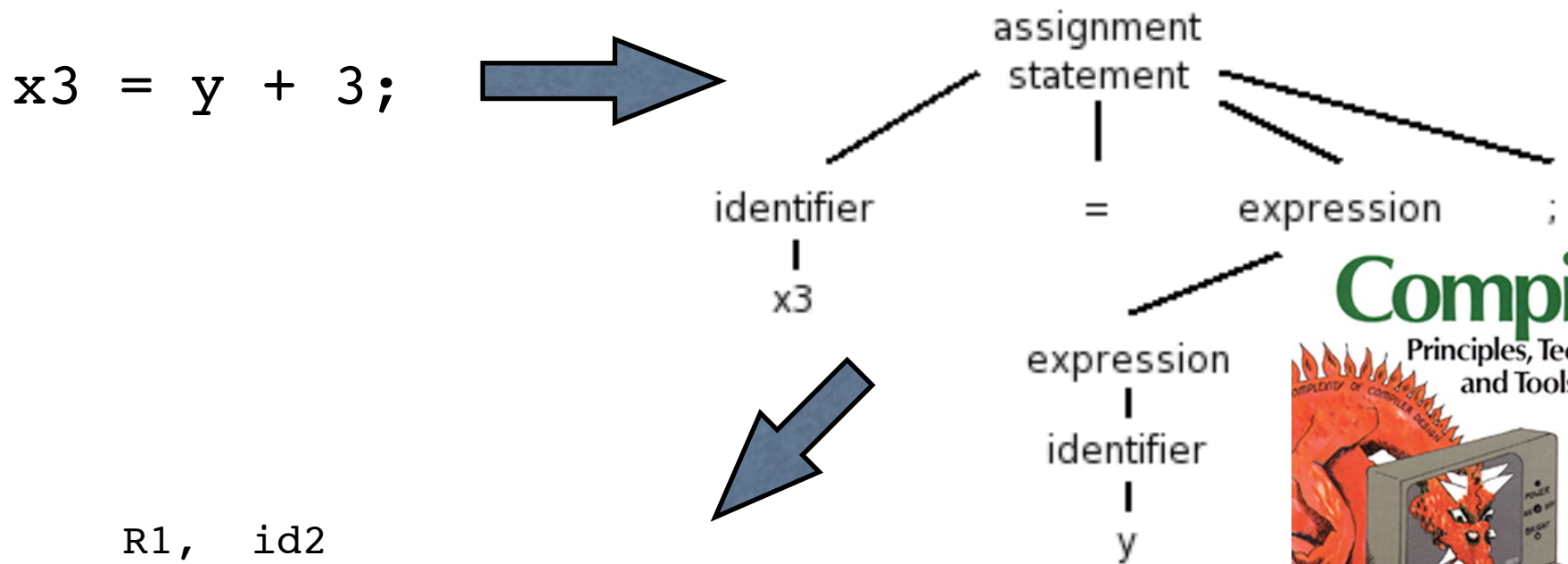
1. parse high-level language program into a syntax tree
2. generate intermediate or machine code accordingly



```
LD    R1, id2
ADDF  R1, R1, #3.0 // add float
RTOI  R2, R1      // real to int
ST    id1, R2
```

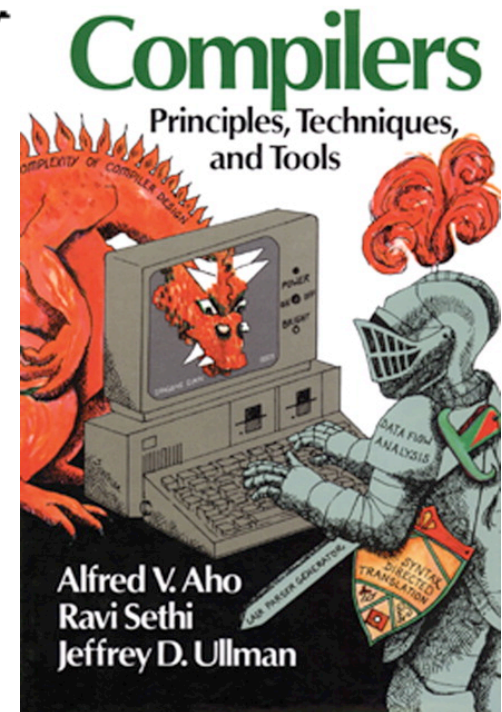

How do compilers translate?

1. parse high-level language program into a syntax tree
2. generate intermediate or machine code accordingly



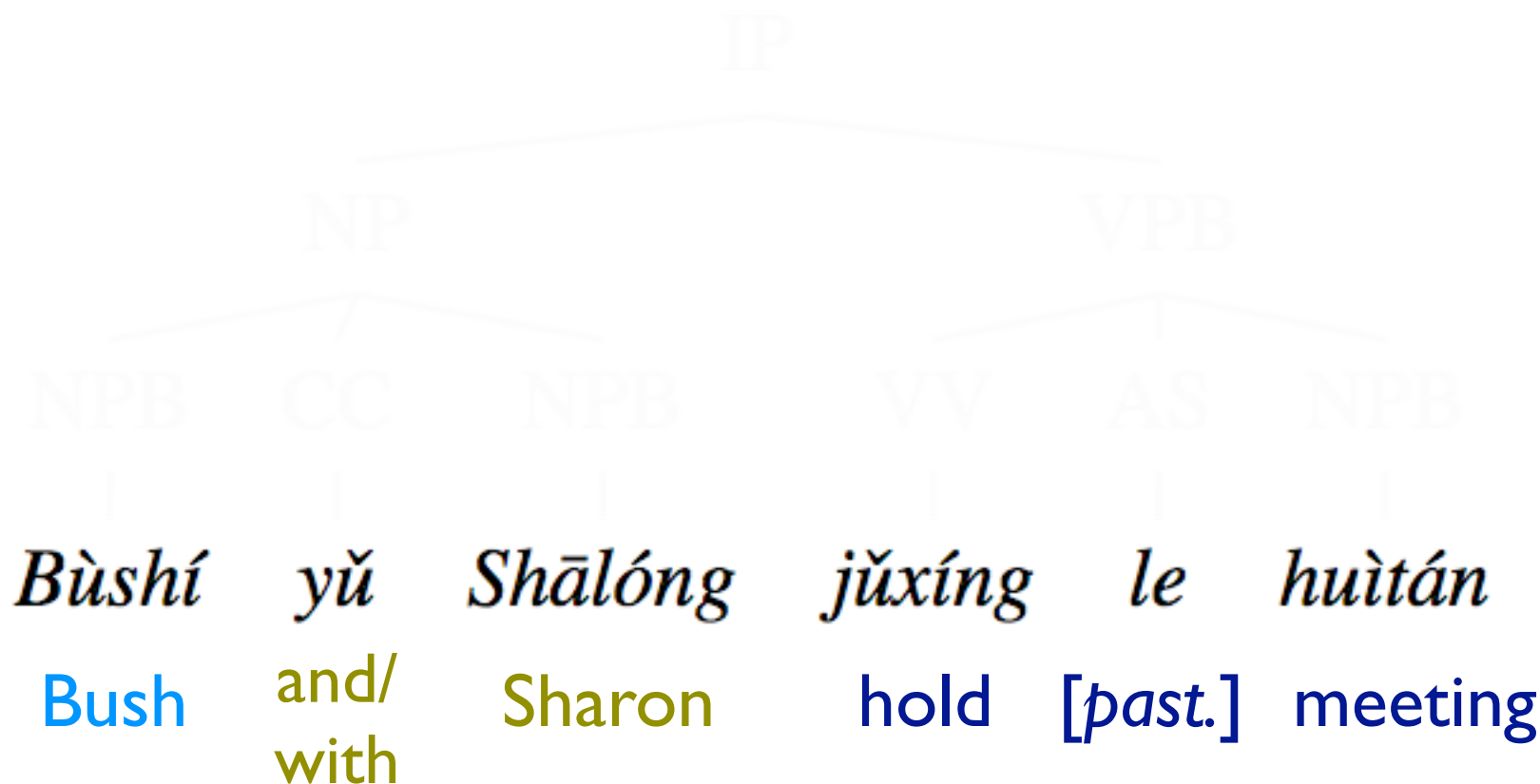
```
LD    R1, id2
ADDF  R1, R1, #3.0 // add float
RTOI  R2, R1      // real to int
ST    id1, R2
```

syntax-directed translation (~1960)



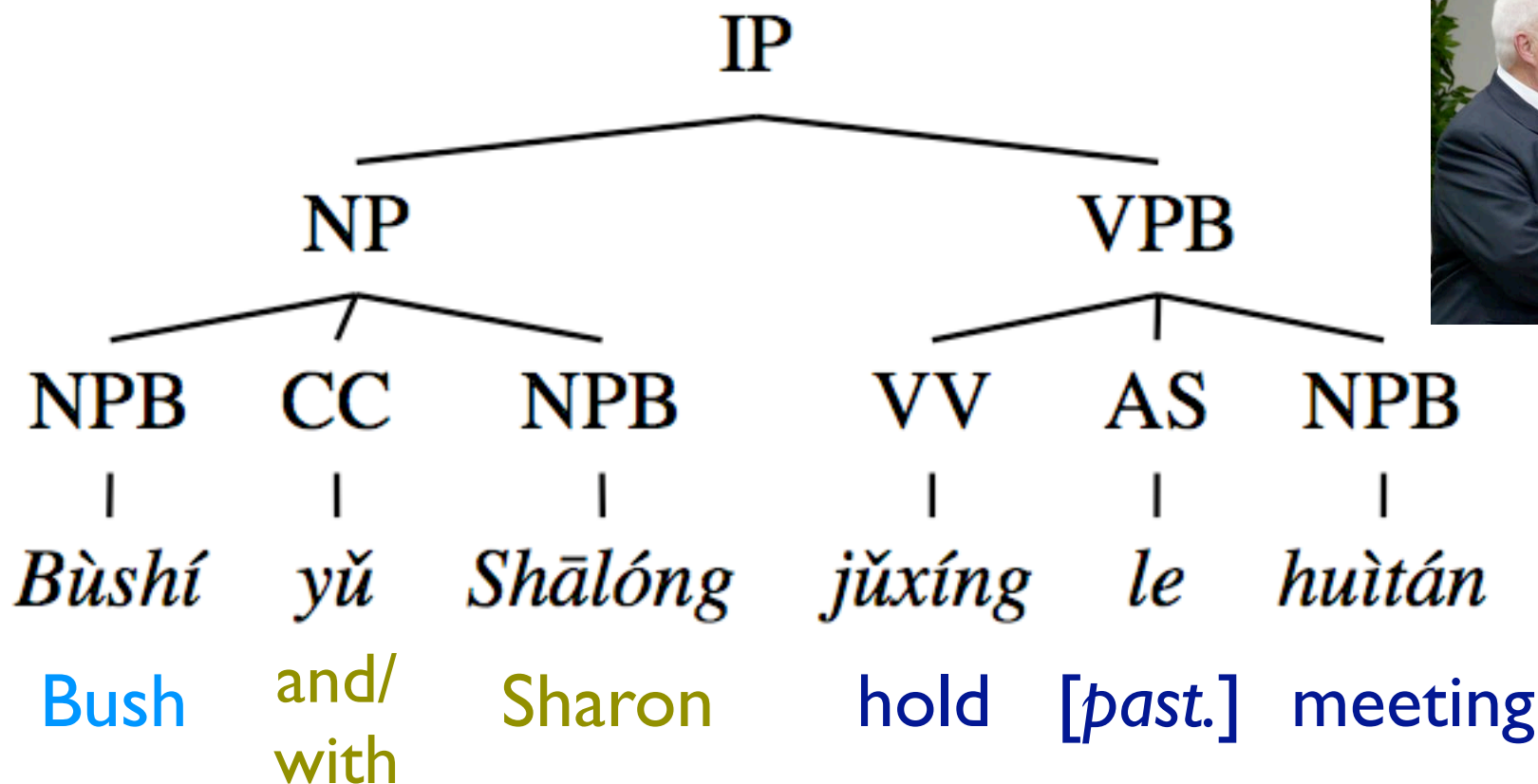
Syntax-Directed Machine Translation

1. parse the source-language sentence into a tree
2. recursively convert it into a target-language sentence



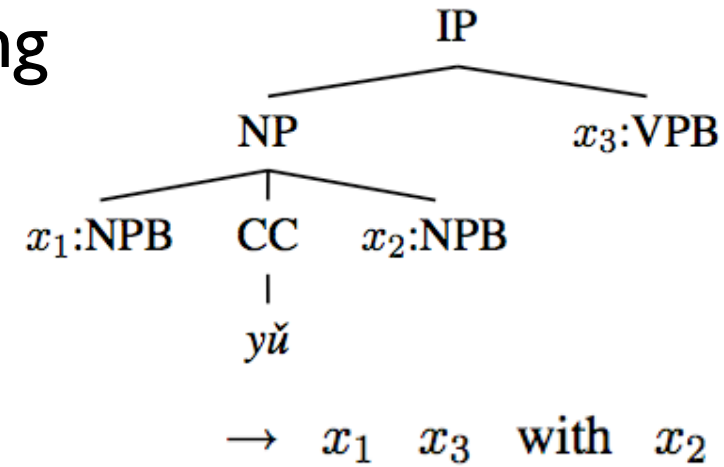
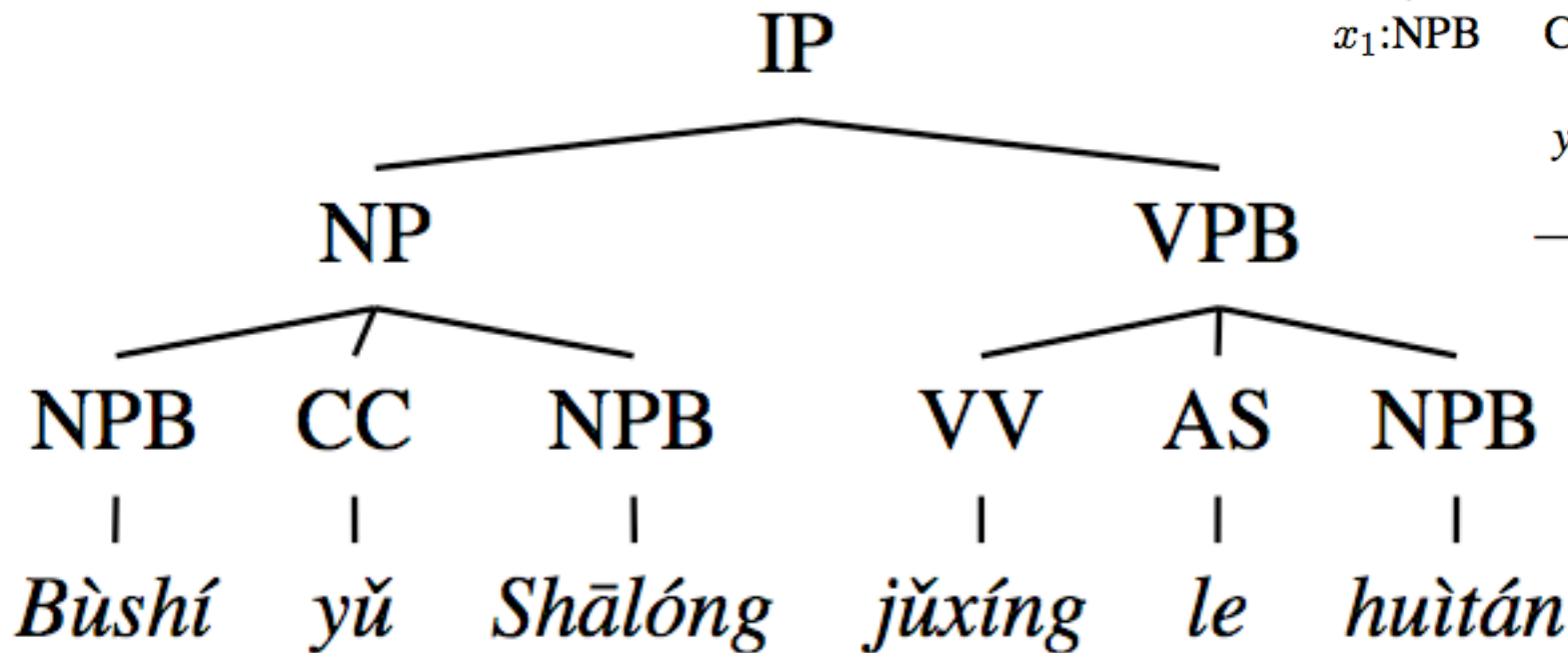
Syntax-Directed Machine Translation

1. parse the source-language sentence into a tree
2. recursively convert it into a target-language sentence



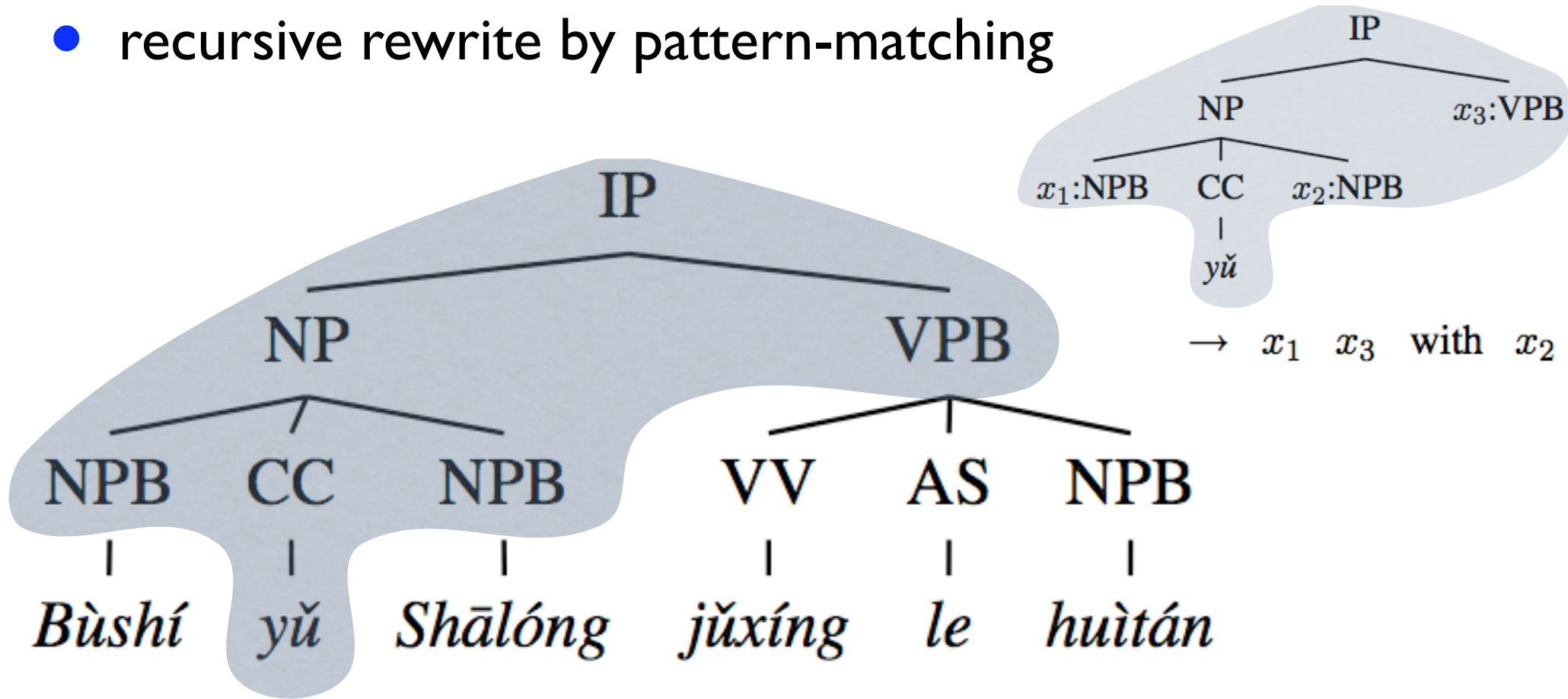
Syntax-Directed Machine Translation

- recursive rewrite by pattern-matching



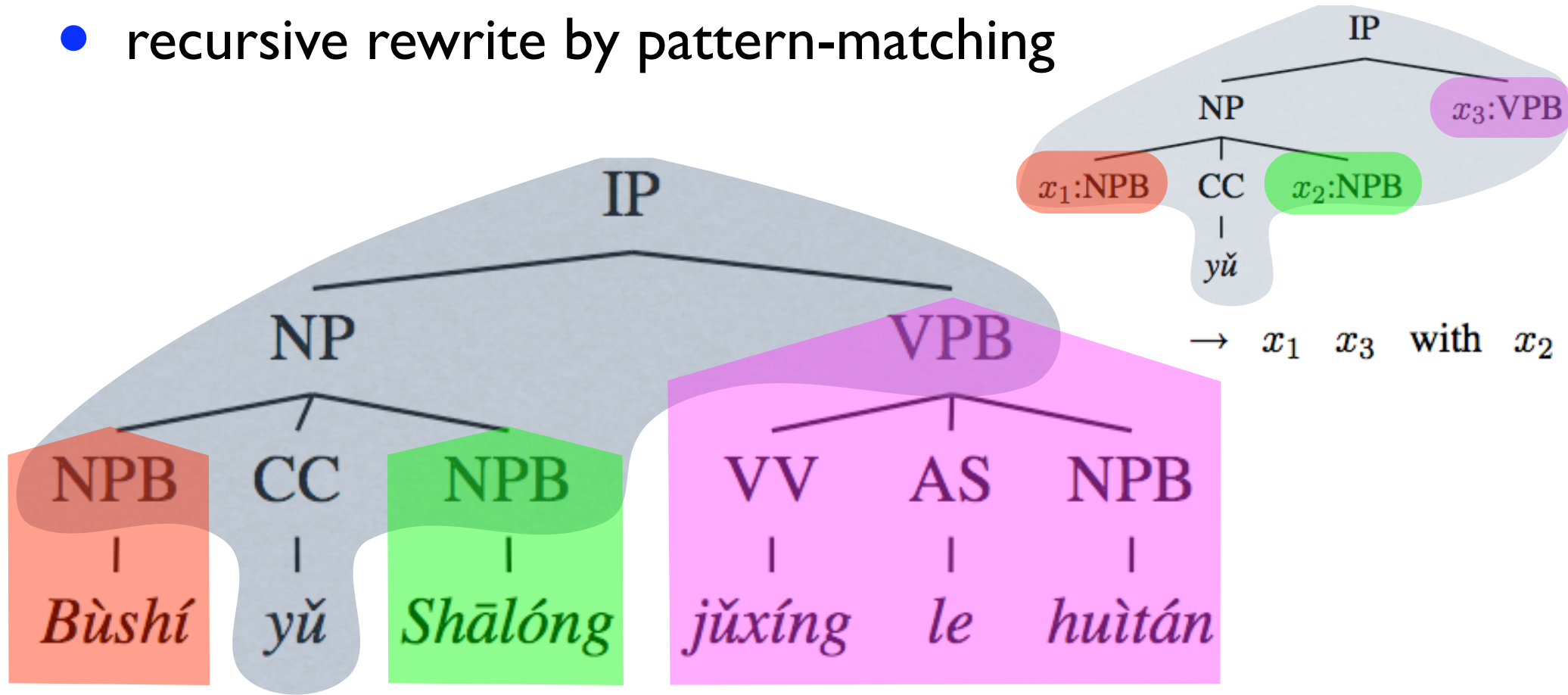
Syntax-Directed Machine Translation

- recursive rewrite by pattern-matching



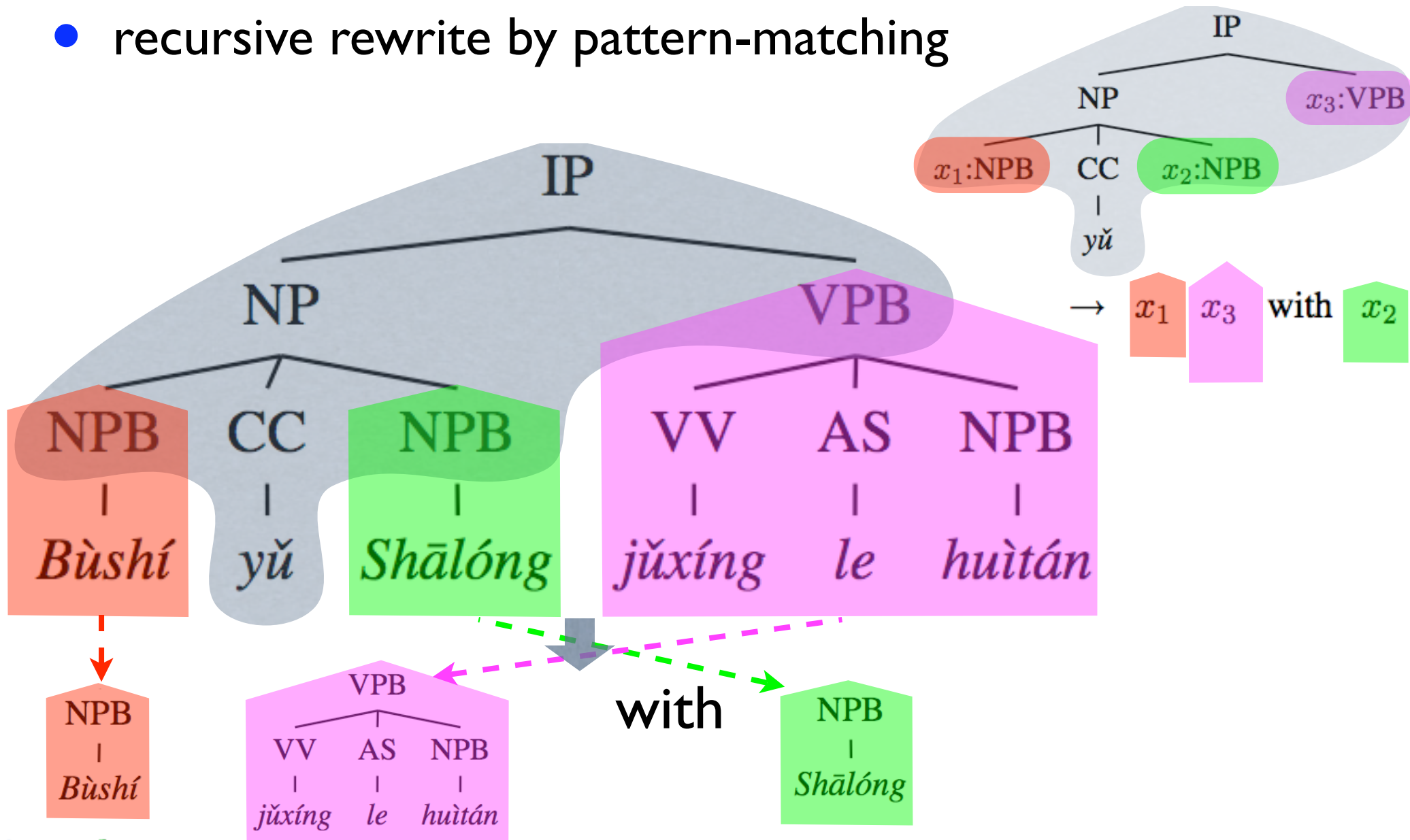
Syntax-Directed Machine Translation

- recursive rewrite by pattern-matching



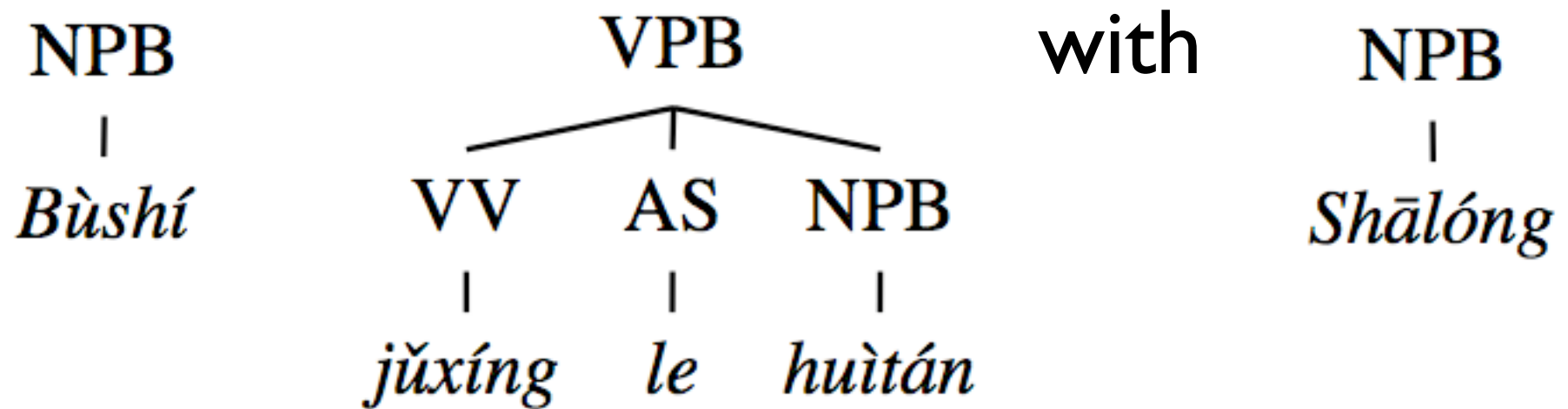
Syntax-Directed Machine Translation

- recursive rewrite by pattern-matching



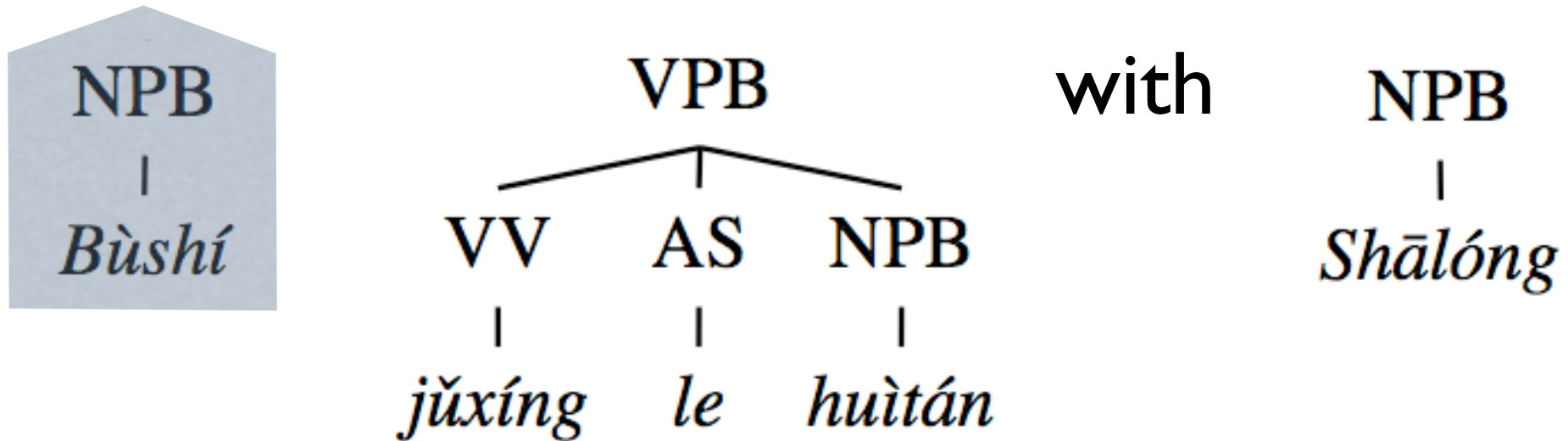
Syntax-Directed Machine Translation?

- recursively solve unfinished subproblems



Syntax-Directed Machine Translation?

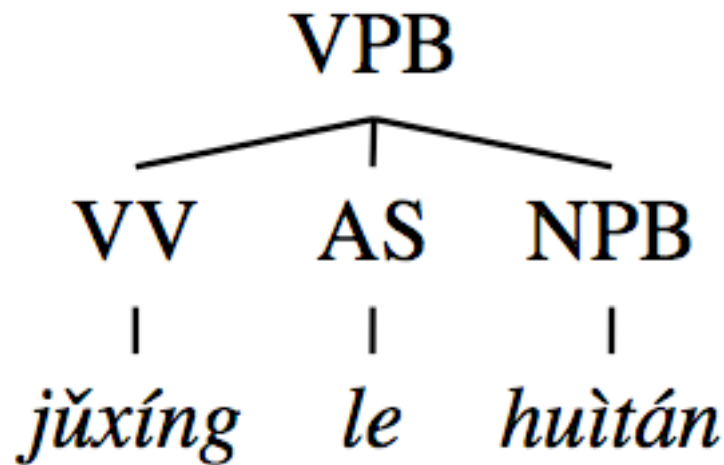
- recursively solve unfinished subproblems



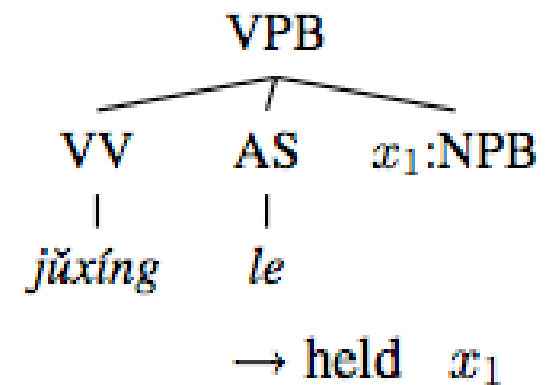
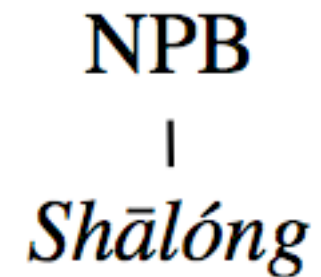
Syntax-Directed Machine Translation?

- recursively solve unfinished subproblems

Bush



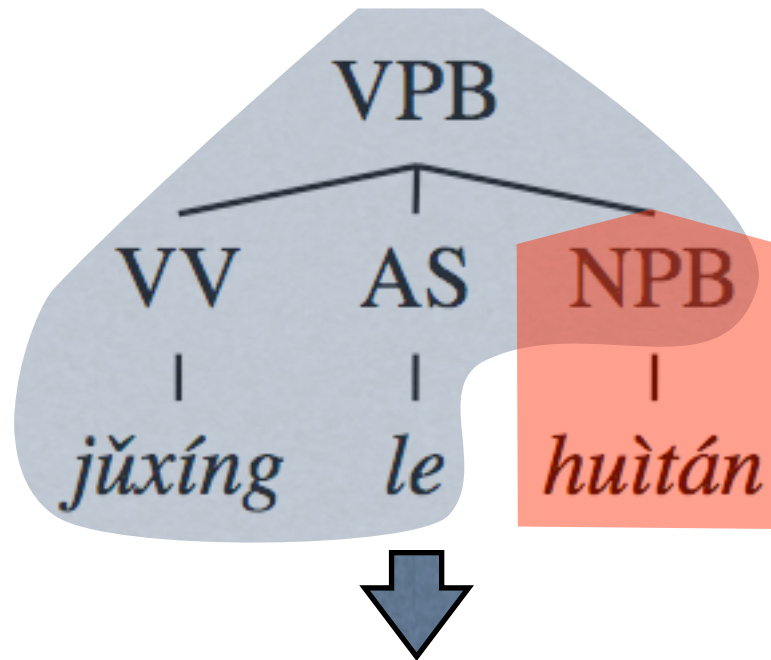
with



Syntax-Directed Machine Translation?

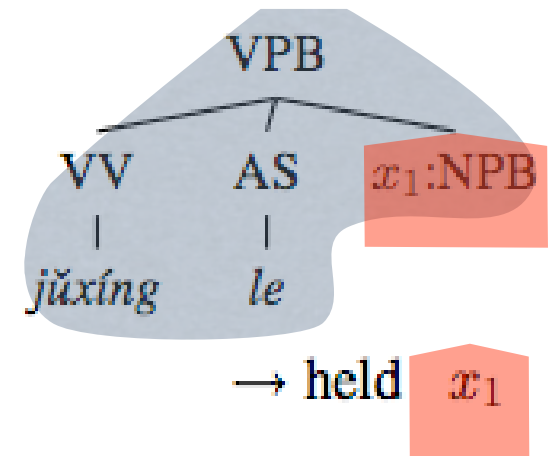
- recursively solve unfinished subproblems

Bush



with

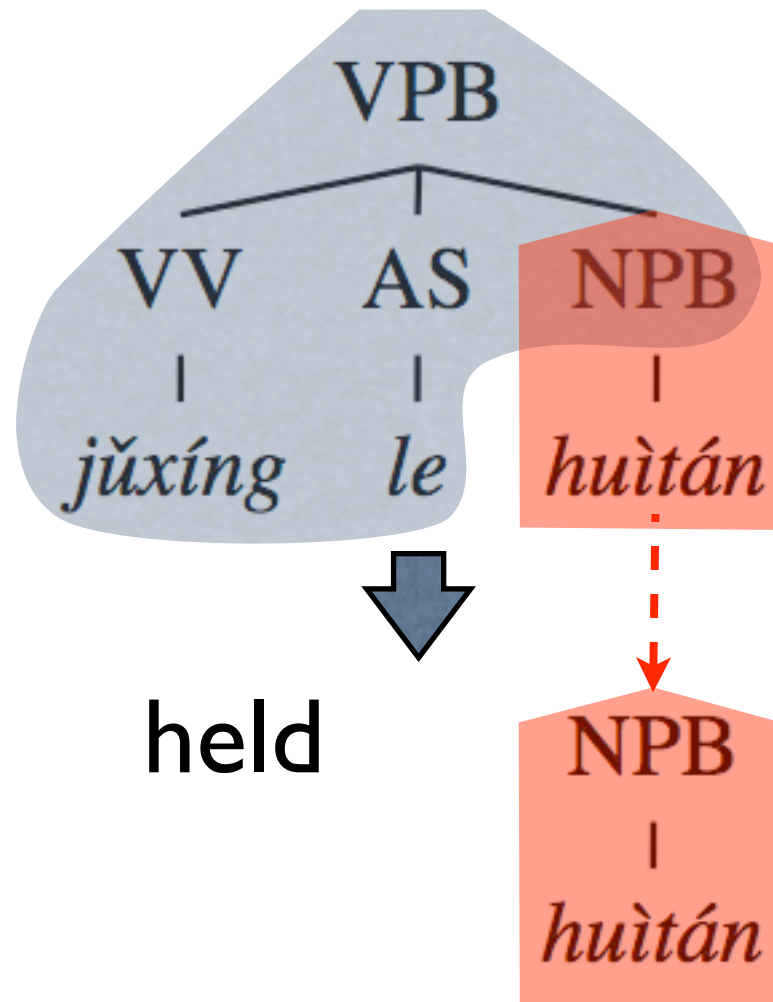
NPB
|
Shānlóng



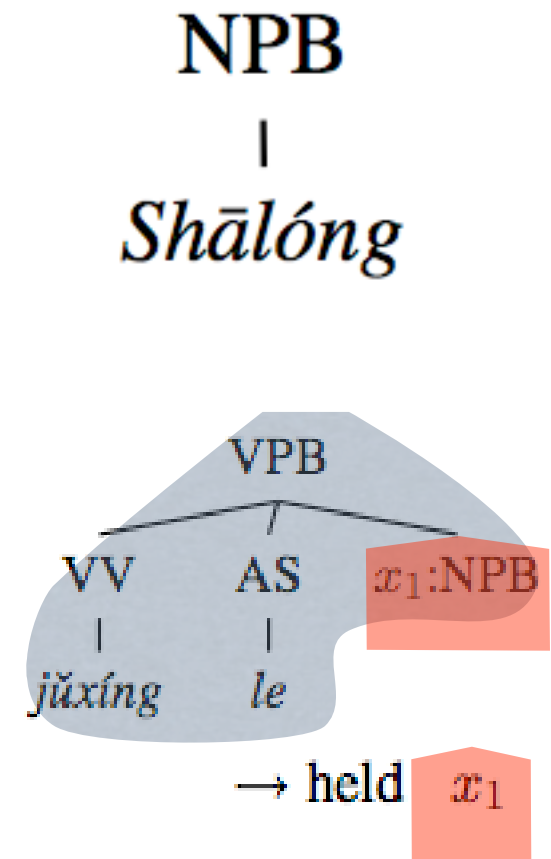
Syntax-Directed Machine Translation?

- recursively solve unfinished subproblems

Bush



with



Syntax-Directed Machine Translation?

- continue pattern-matching

Bush

held

NPB

with

NPB

|

huìtán

|

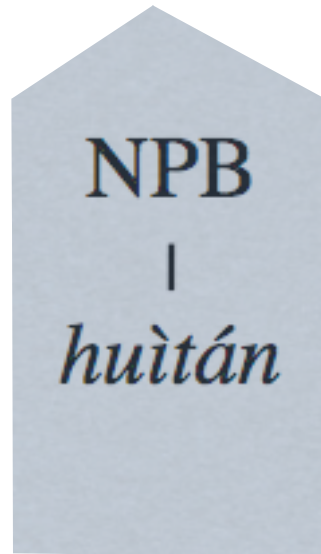
Shānlóng

Syntax-Directed Machine Translation?

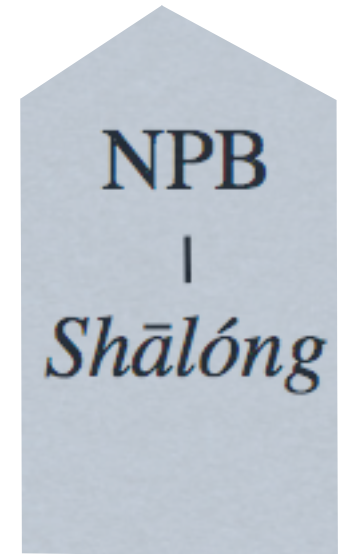
- continue pattern-matching

Bush

held



with



a meeting

Sharon

Syntax-Directed Machine Translation?

- continue pattern-matching

Bush held a meeting with Sharon

Syntax-Directed Machine Translation?

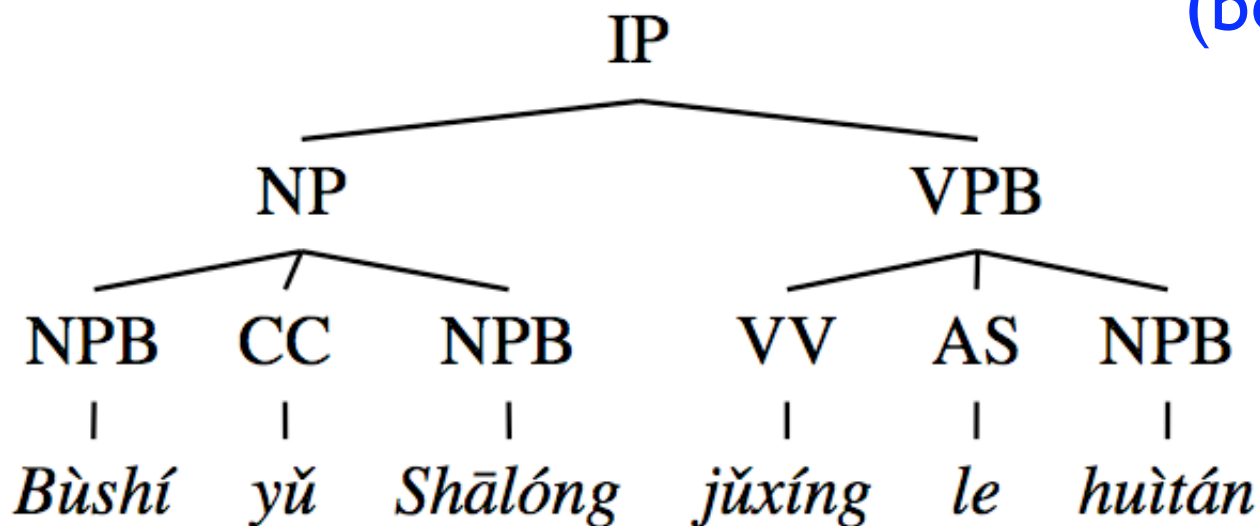
- continue pattern-matching

Bush held a meeting with Sharon

Pros: simple, fast, and expressive

- simple architecture: separate parsing and translation
- efficient linear-time dynamic programming
 - “soft decision” at each node on which rule to use
 - (trivial) depth-first traversal with memoization
- expressive multi-level rules for syntactic divergence

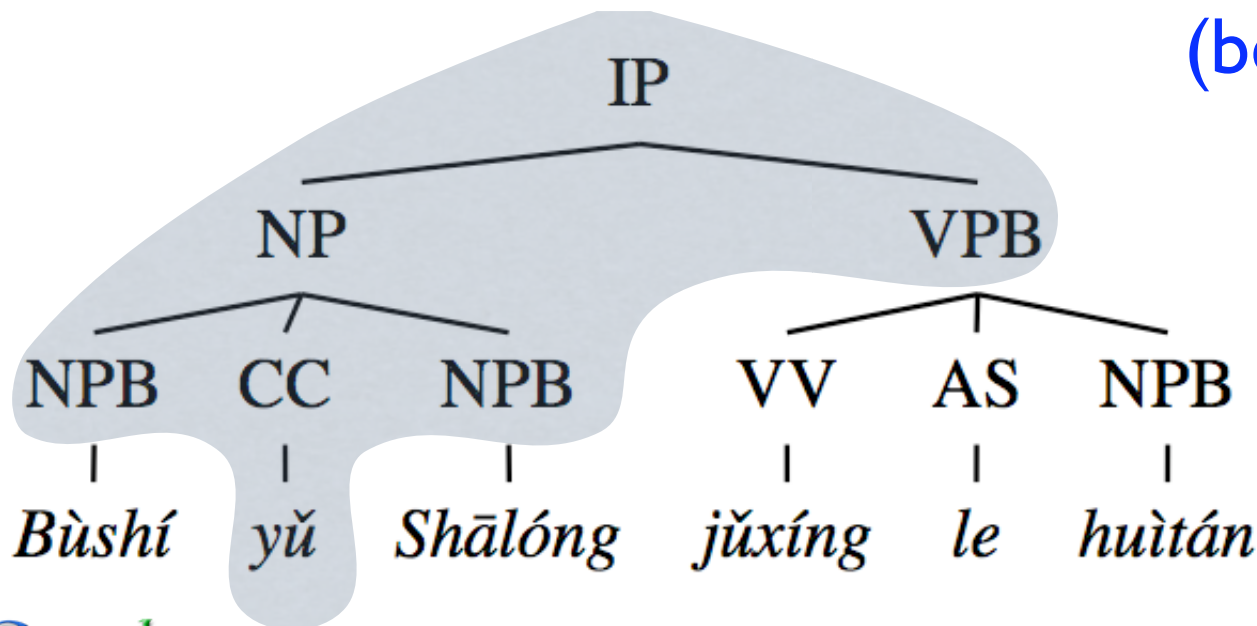
(beyond CFG)



Pros: simple, fast, and expressive

- simple architecture: separate parsing and translation
- efficient linear-time dynamic programming
 - “soft decision” at each node on which rule to use
 - (trivial) depth-first traversal with memoization
- expressive multi-level rules for syntactic divergence

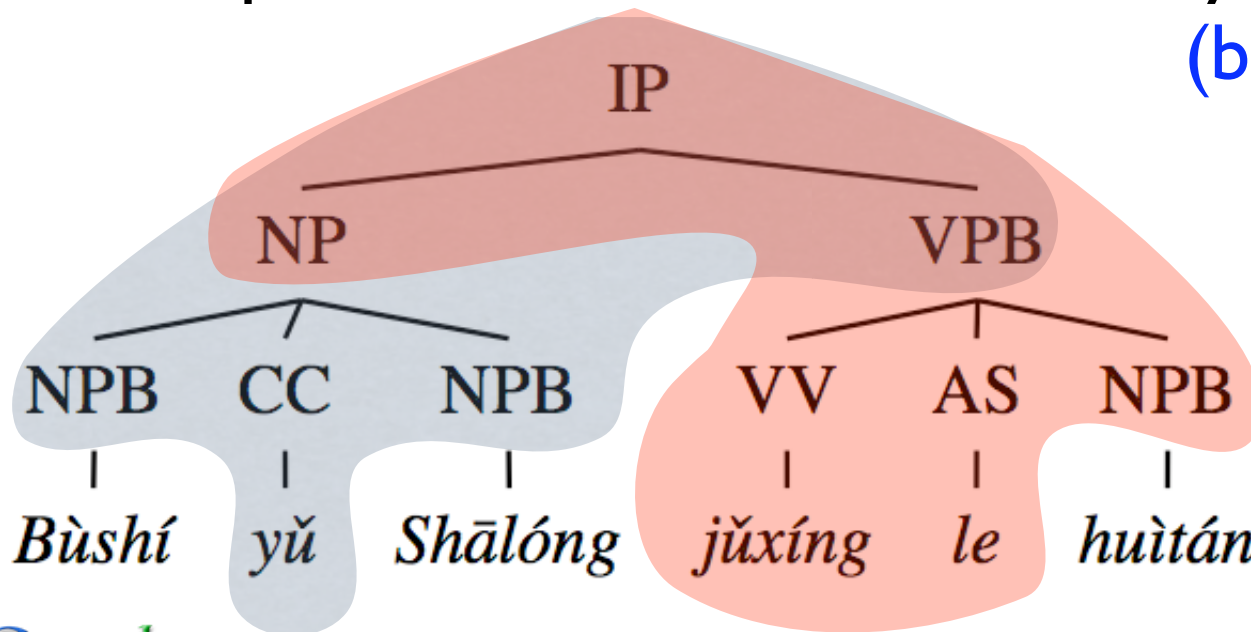
(beyond CFG)



Pros: simple, fast, and expressive

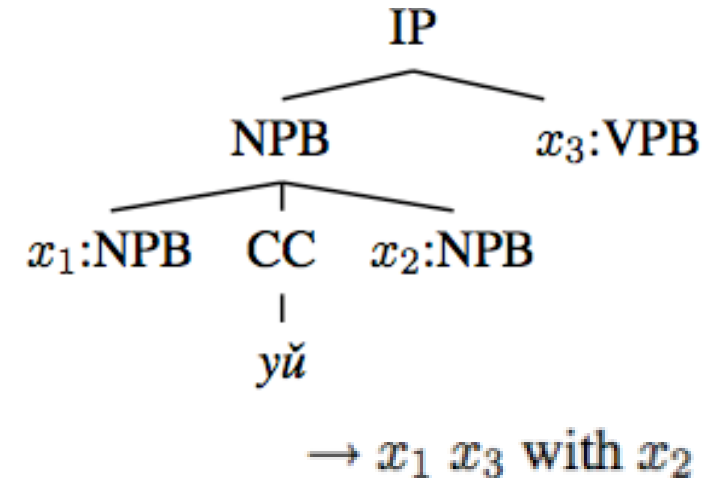
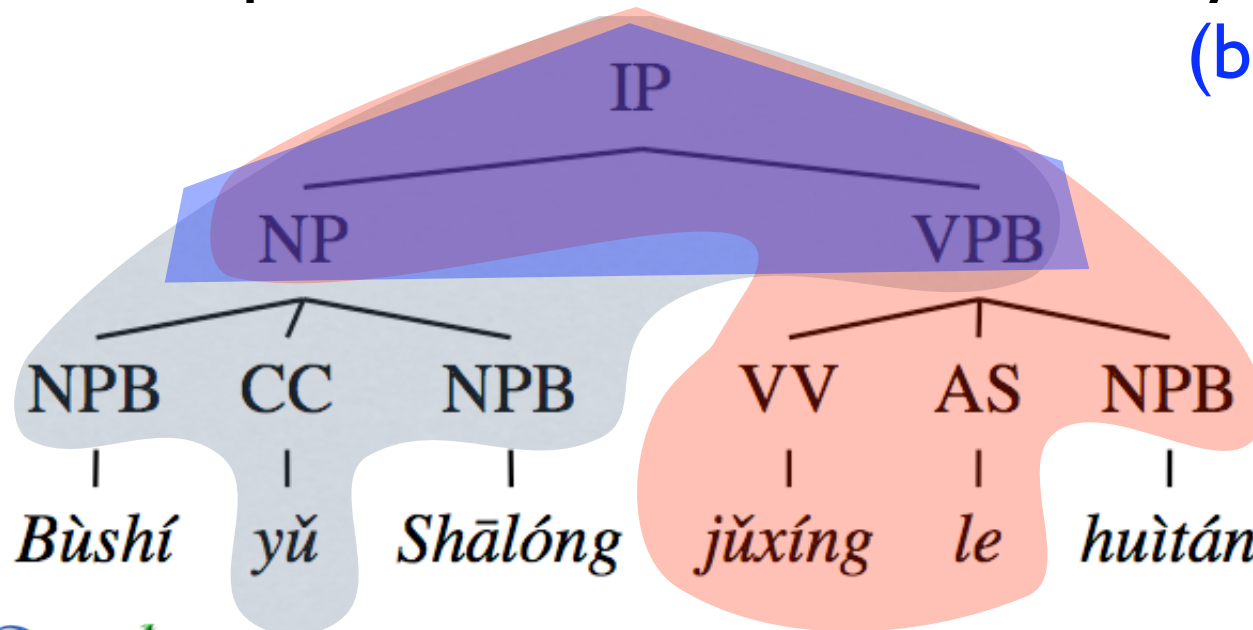
- simple architecture: separate parsing and translation
- efficient linear-time dynamic programming
 - “soft decision” at each node on which rule to use
 - (trivial) depth-first traversal with memoization
- expressive multi-level rules for syntactic divergence

(beyond CFG)



Pros: simple, fast, and expressive

- simple architecture: separate parsing and translation
- efficient linear-time dynamic programming
 - “soft decision” at each node on which rule to use
 - (trivial) depth-first traversal with memoization
- expressive multi-level rules for syntactic divergence (beyond CFG)



Cons: Parsing Errors

- ambiguity is a fundamental problem in natural languages
 - probably will never have perfect parsers (unlike compiling)
- parsing errors affect translation quality!



Cons: Parsing Errors

- ambiguity is a fundamental problem in natural languages
 - probably will never have perfect parsers (unlike compiling)
- parsing errors affect translation quality!



emergency exit
or “safe exports”?

Cons: Parsing Errors

- ambiguity is a fundamental problem in natural languages
- probably will never have perfect parsers (unlike compiling)
- parsing errors affect translation quality!



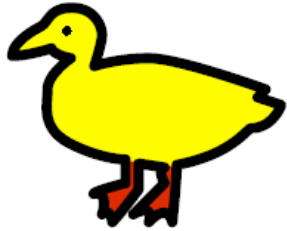
emergency exit
or “safe exports”?



mind your head
or “meet cautiously”?

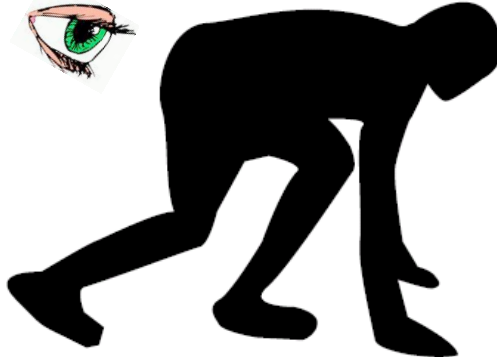
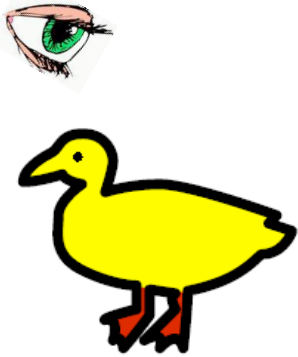
Exponential Explosion of Ambiguity

I saw her duck.



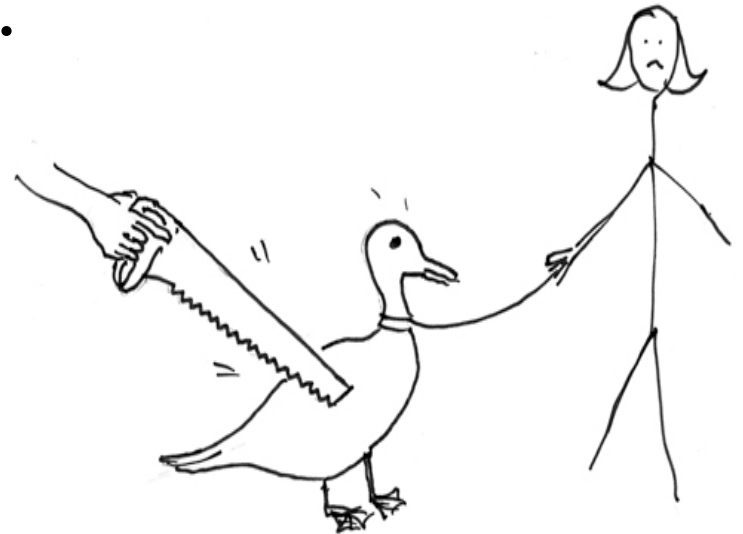
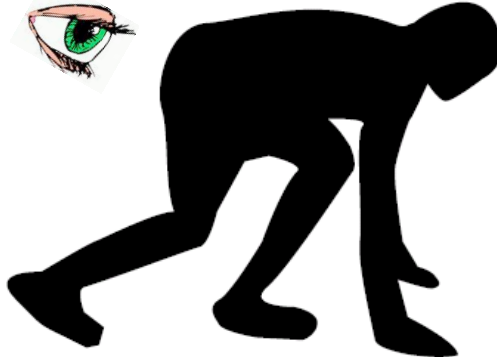
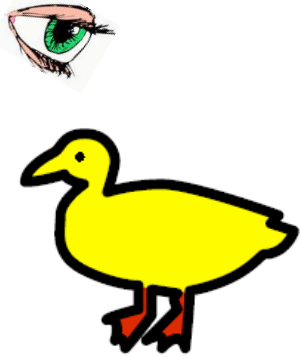
Exponential Explosion of Ambiguity

I saw her duck.



Exponential Explosion of Ambiguity

I saw her duck.



Exponential Explosion of Ambiguity

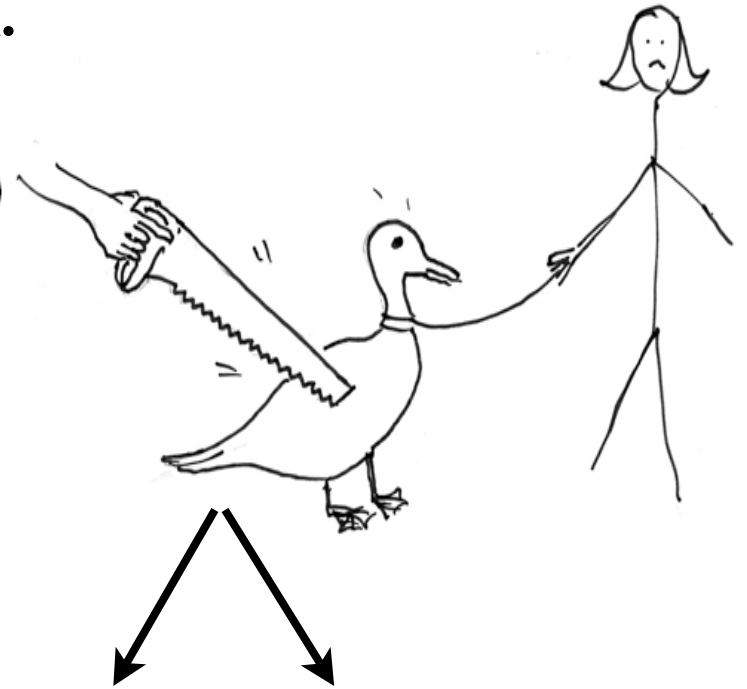
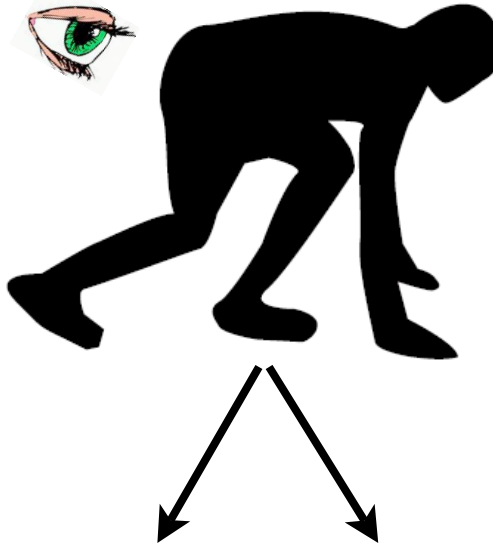
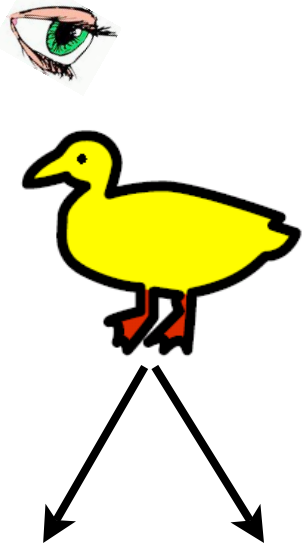
I saw her duck.



- how about...
 - I saw her duck with a telescope.
 - I saw her duck with a telescope in the garden...

Exponential Explosion of Ambiguity

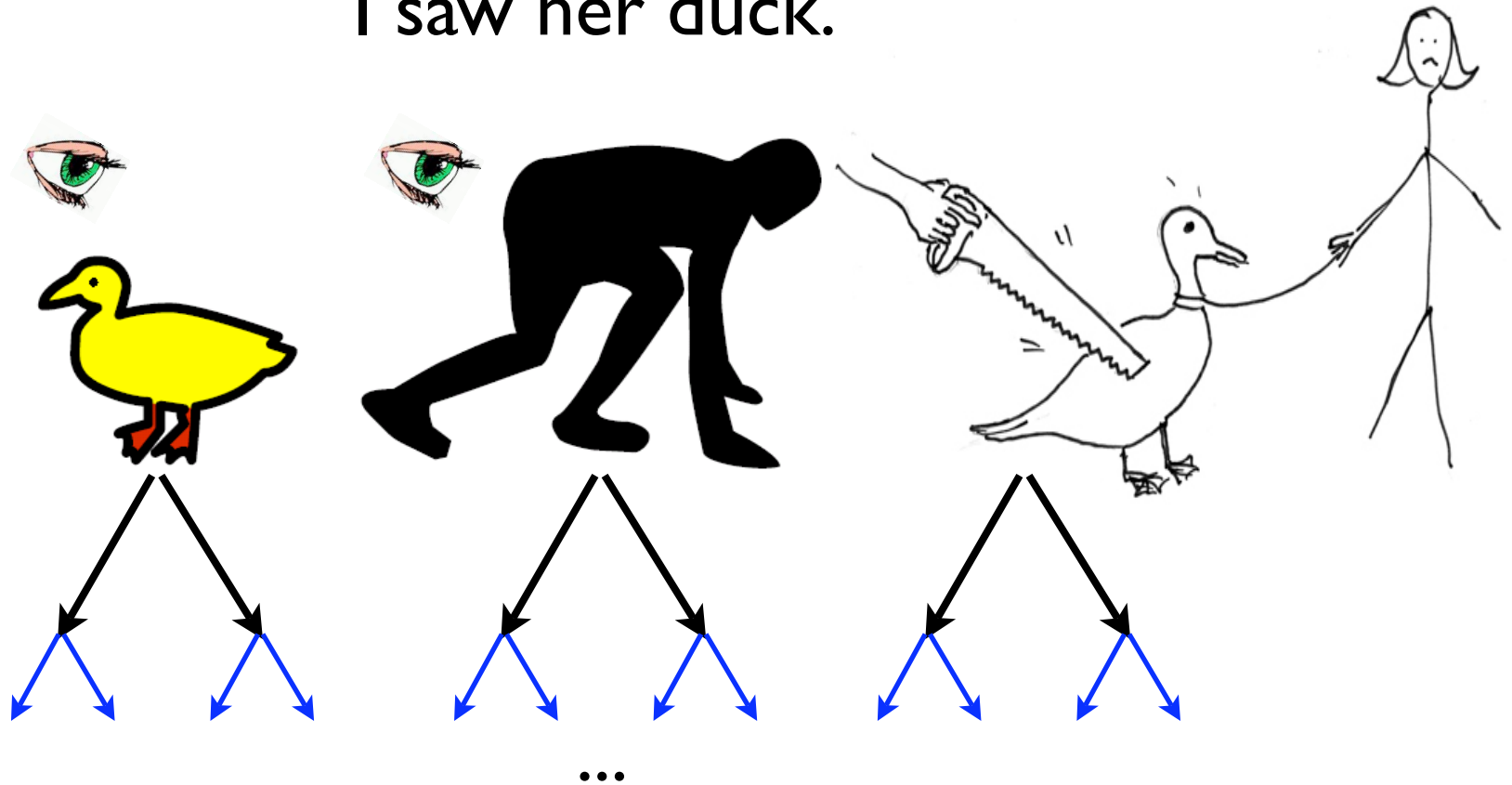
I saw her duck.



- how about...
 - I saw her duck with a telescope.
 - I saw her duck with a telescope in the garden...

Exponential Explosion of Ambiguity

I saw her duck.



- how about...
 - I saw her duck with a telescope.
 - I saw her duck with a telescope in the garden...

Tackling Ambiguities in Translation

Tackling Ambiguities in Translation

- simplest idea: take top- k trees rather than 1-best parse
 - but only covers tiny fraction of the exponential space
 - and these k -best trees are very similar
 - e.g., 50-best trees \sim 5-6 binary ambiguities ($2^5 < 50 < 2^6$)
 - very inefficient to translate on these very similar trees

Tackling Ambiguities in Translation

- simplest idea: take top- k trees rather than 1-best parse
 - but only covers tiny fraction of the exponential space
 - and these k -best trees are very similar
 - e.g., 50-best trees \sim 5-6 binary ambiguities ($2^5 < 50 < 2^6$)
 - very inefficient to translate on these very similar trees
- most ambitious idea: combining parsing and translation
 - start from the input string, rather than 1-best tree
 - essentially considering all trees (search space too big)

Tackling Ambiguities in Translation

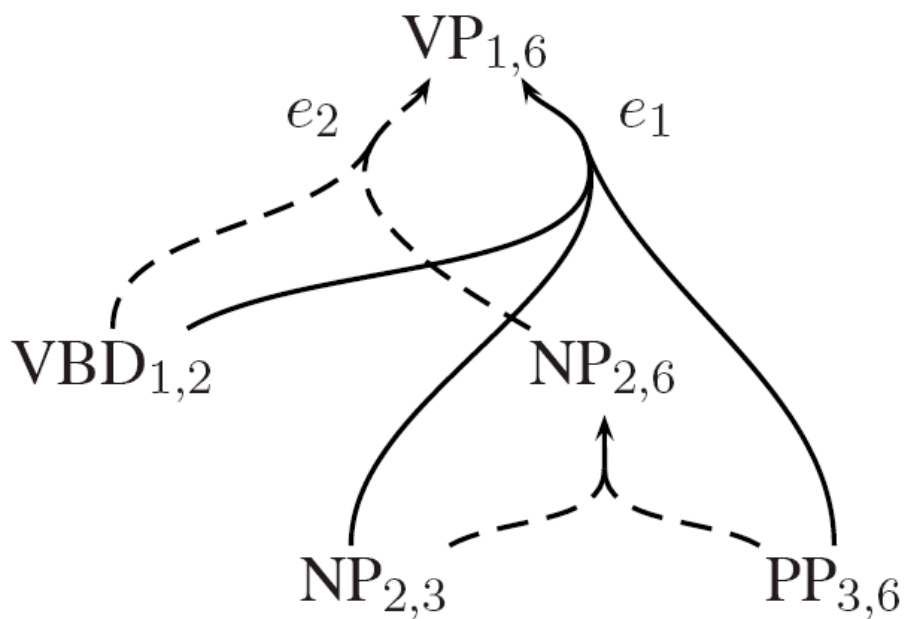
- simplest idea: take top- k trees rather than 1-best parse
 - but only covers tiny fraction of the exponential space
 - and these k -best trees are very similar
 - e.g., 50-best trees \sim 5-6 binary ambiguities ($2^5 < 50 < 2^6$)
 - very inefficient to translate on these very similar trees
- most ambitious idea: combining parsing and translation
 - start from the input string, rather than 1-best tree
 - essentially considering all trees (search space too big)
- our approach: *packed forest* (*poly. encoding of exp. space*)
 - almost as fast as 1-best, almost as good as combined

Outline

- Overview: Tree-based Translation
- Forest-based Translation
 - Packed Forest
 - Translation on a Forest
 - Experiments
- Forest-based Rule Extraction
 - Large-scale Experiments

Packed Forest

- a compact representation of many many parses
- by sharing common sub-derivations
- polynomial-space encoding of exponentially large set



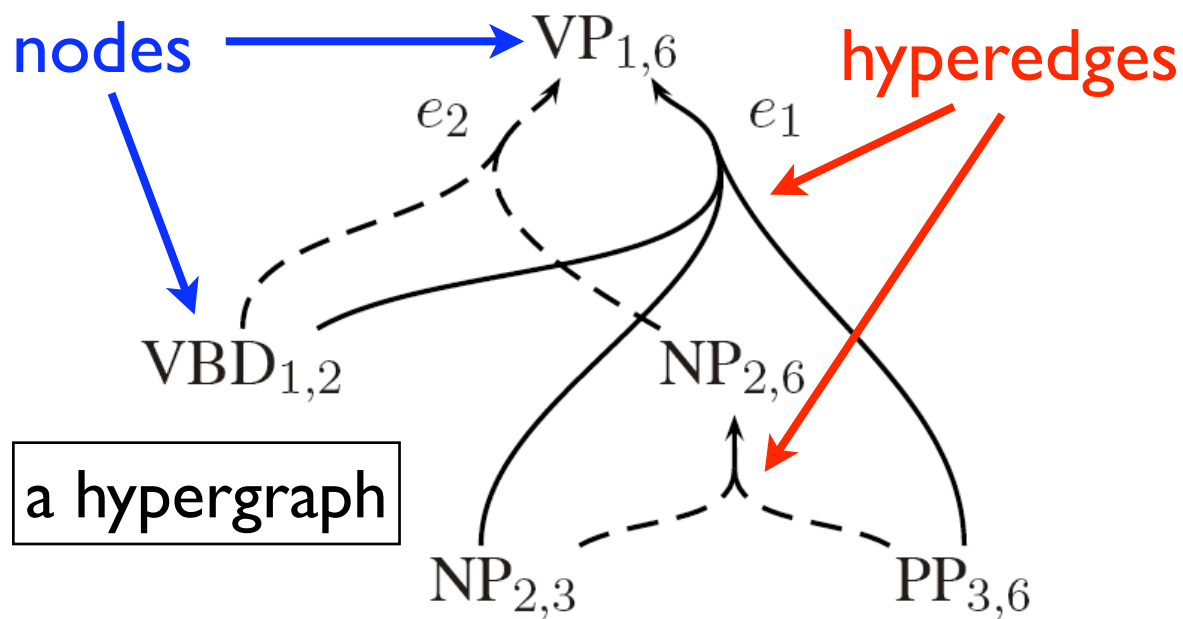
$$e_1 \frac{\text{VBD}_{1,2} \quad \text{NP}_{2,3} \quad \text{PP}_{3,6}}{\text{VP}_{1,6}}$$

0 I 1 saw 2 him 3 with 4 a 5 mirror 6

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

Packed Forest

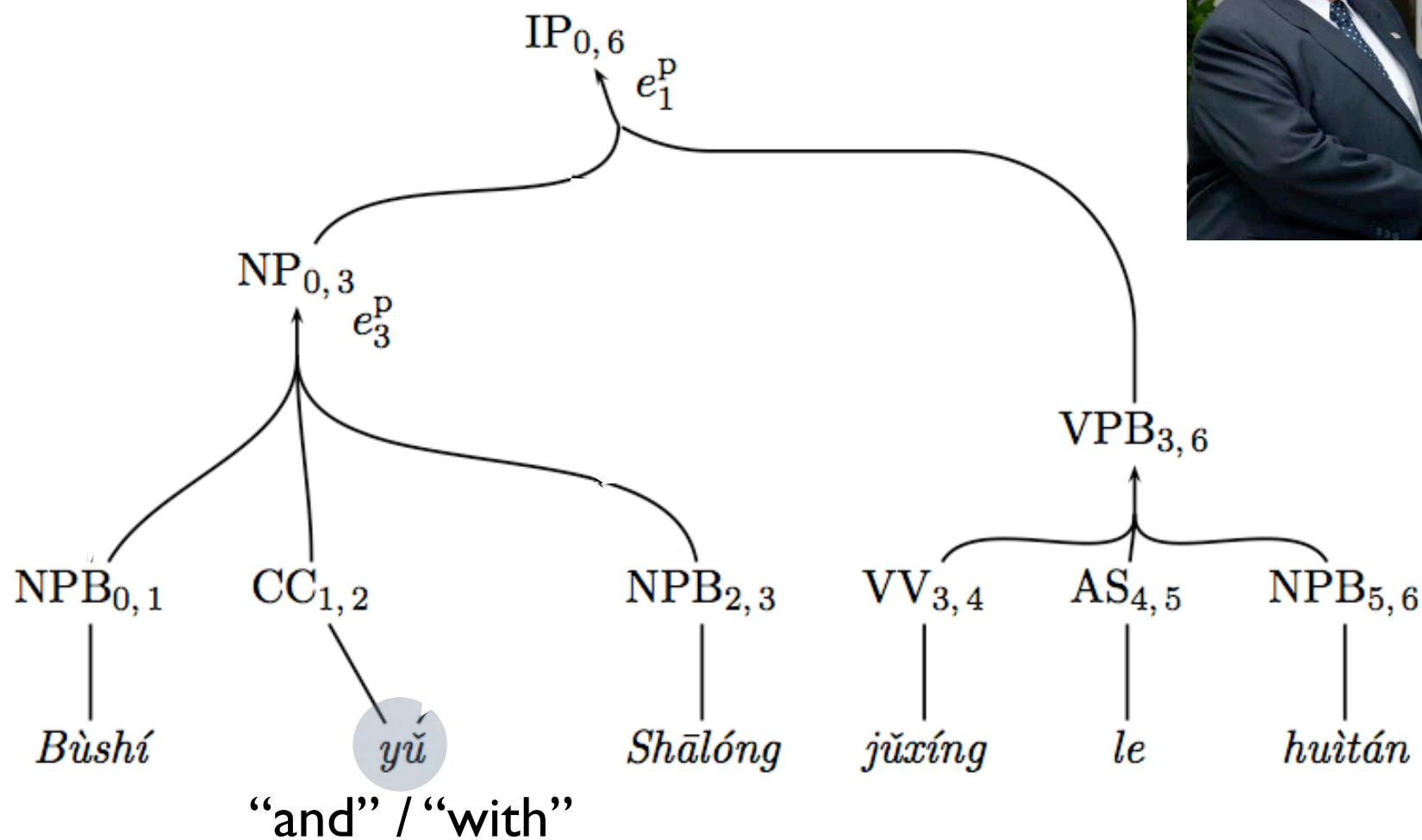
- a compact representation of many many parses
 - by sharing common sub-derivations
 - polynomial-space encoding of exponentially large set



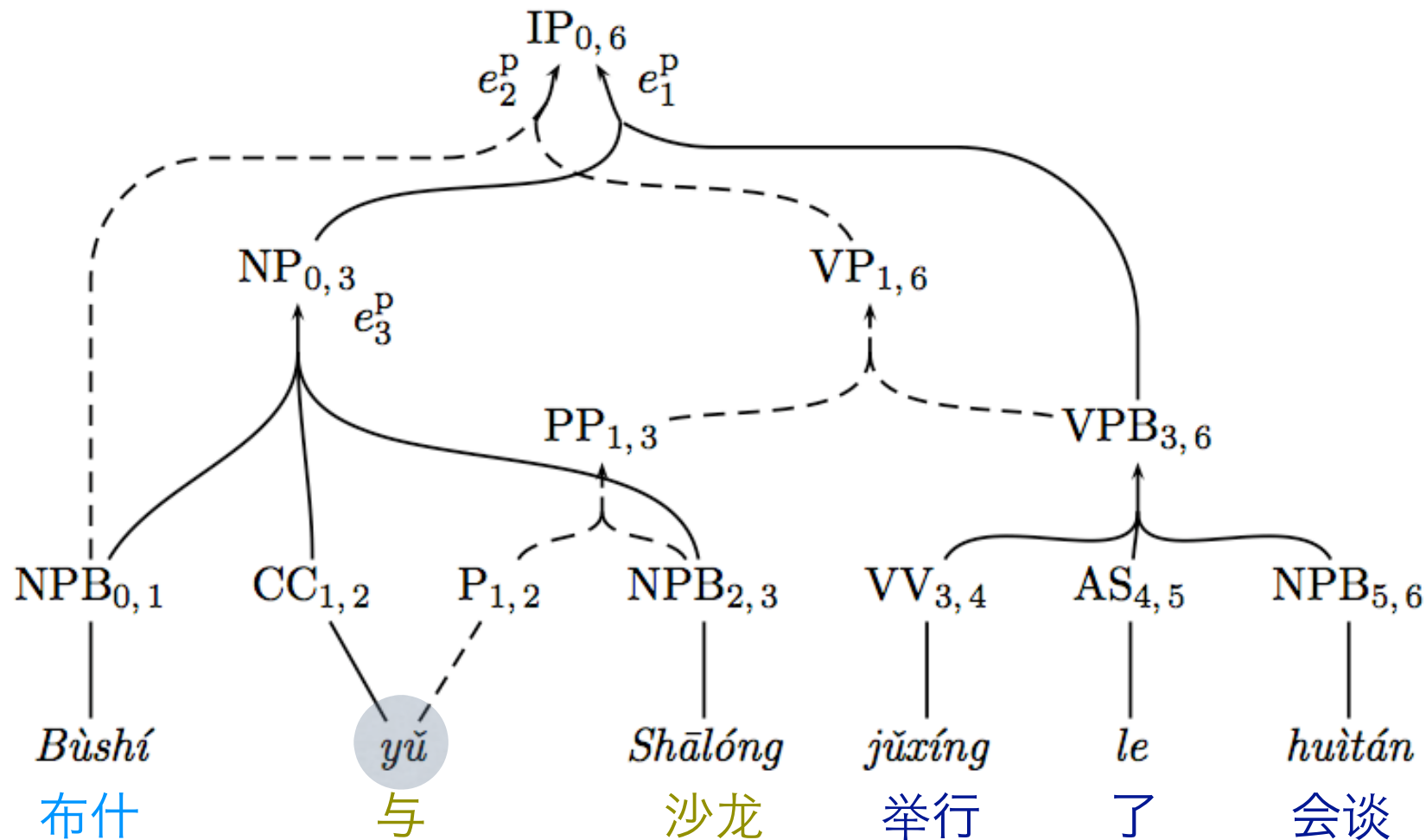
$$e_1 \frac{VBD_{1,2} \quad NP_{2,3} \quad PP_{3,6}}{VP_{1,6}}$$

0 I 1 saw 2 him 3 with 4 a 5 mirror 6

Forest-based Translation



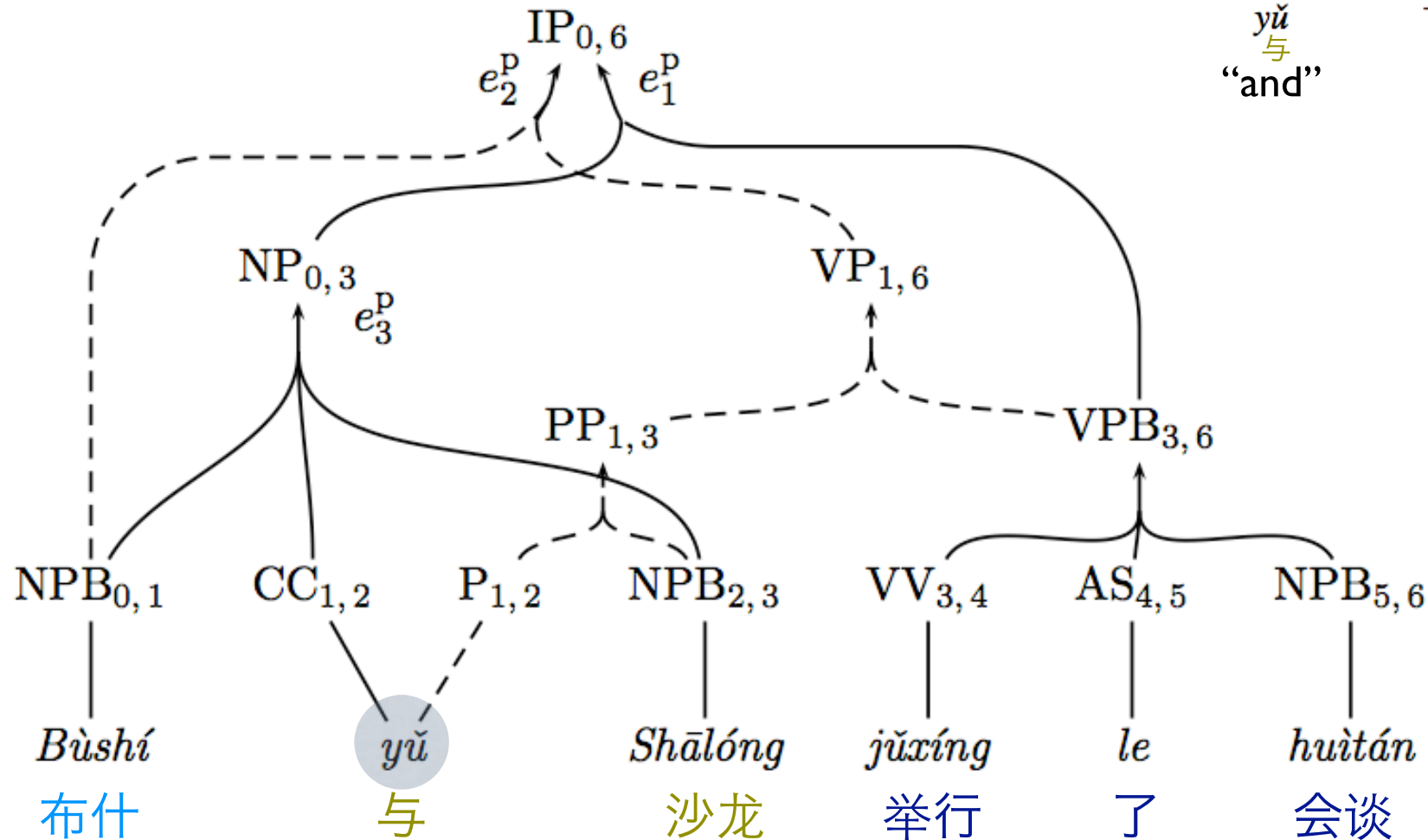
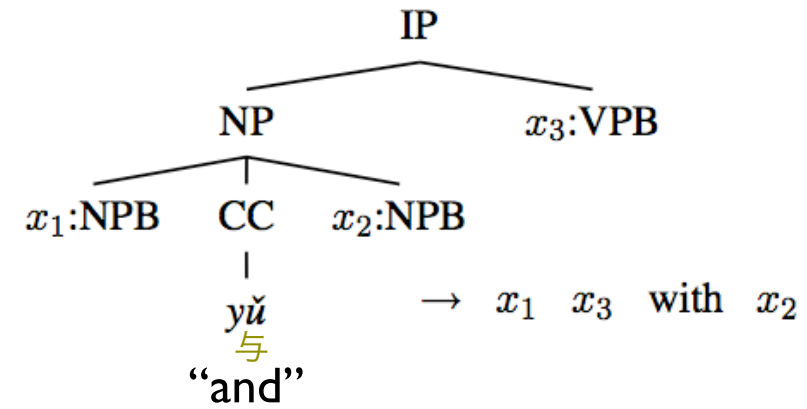
Forest-based Translation



“and” / “with”

Forest-based Translation

pattern-matching
on forest
(linear-time in forest size)

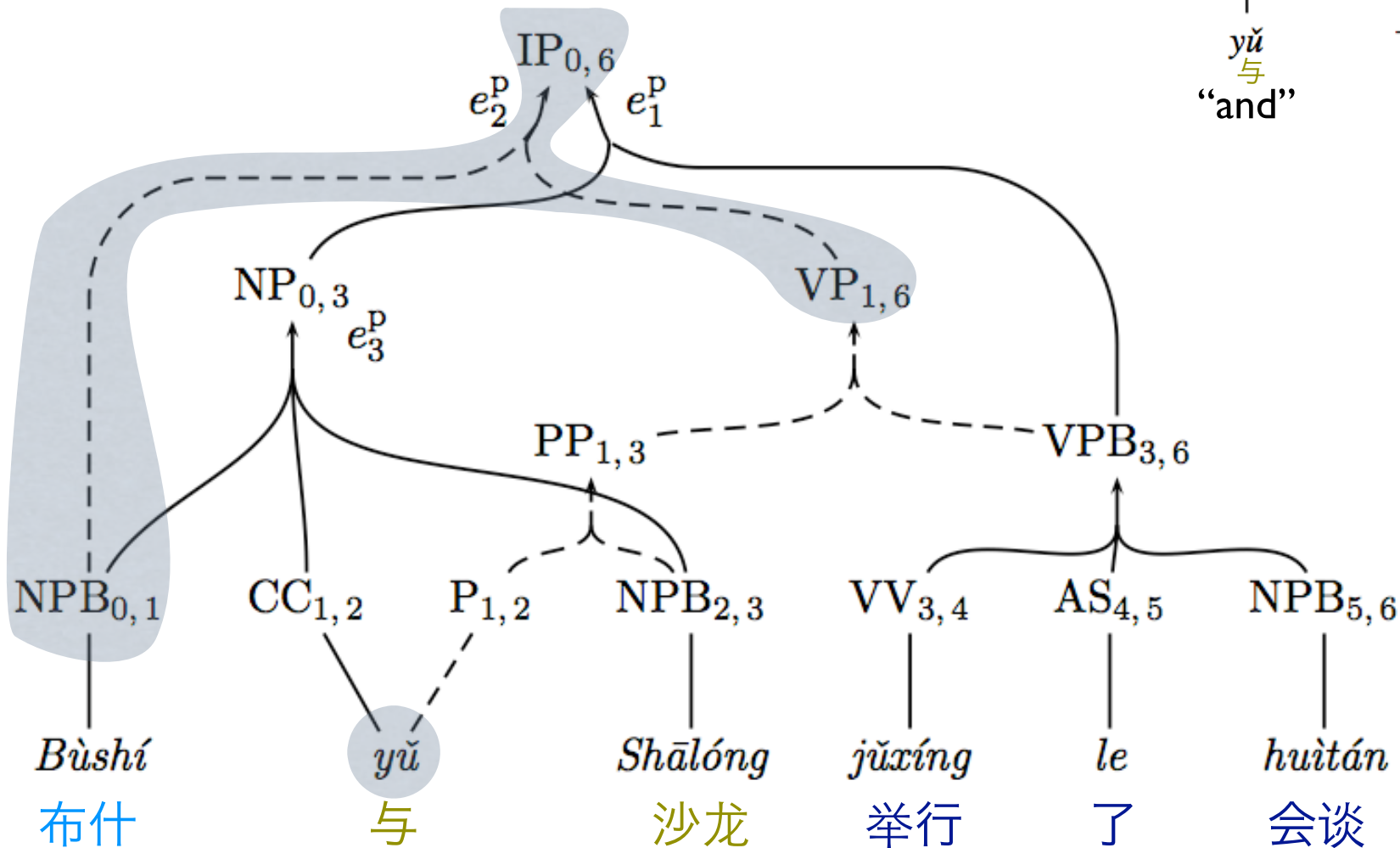
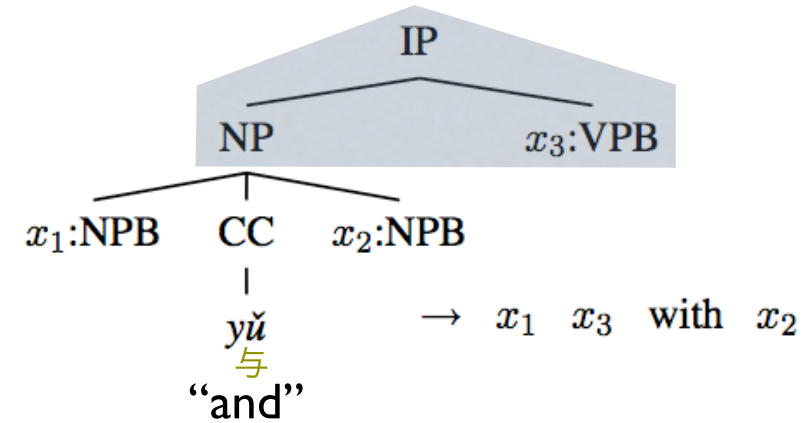


“and” / “with”

Forest-based Translation

pattern-matching
on forest

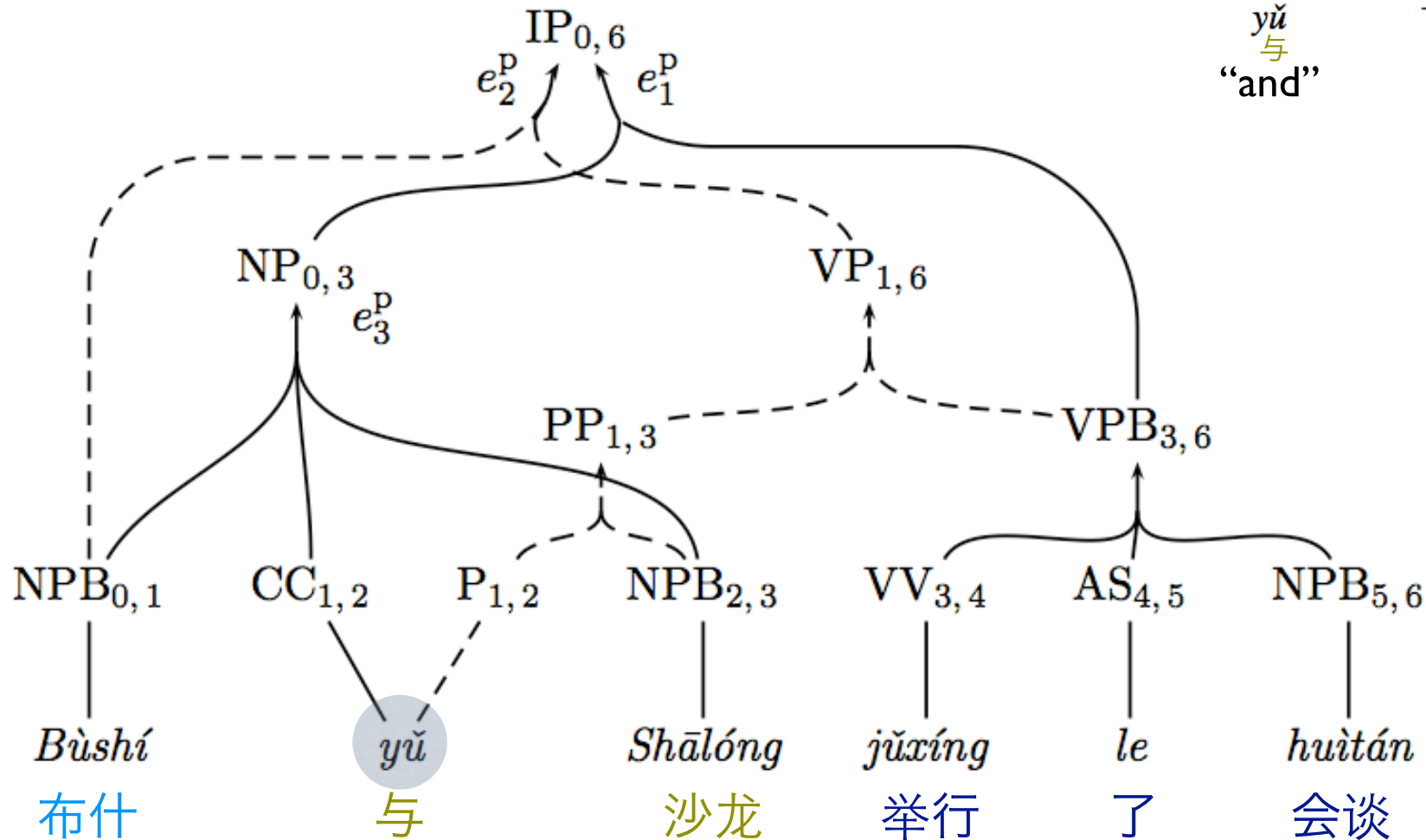
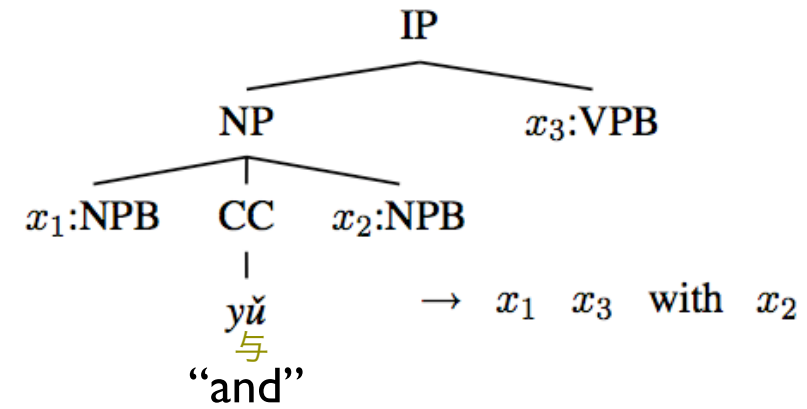
(linear-time in forest size)



“and” / “with”

Forest-based Translation

pattern-matching
on forest
(linear-time in forest size)

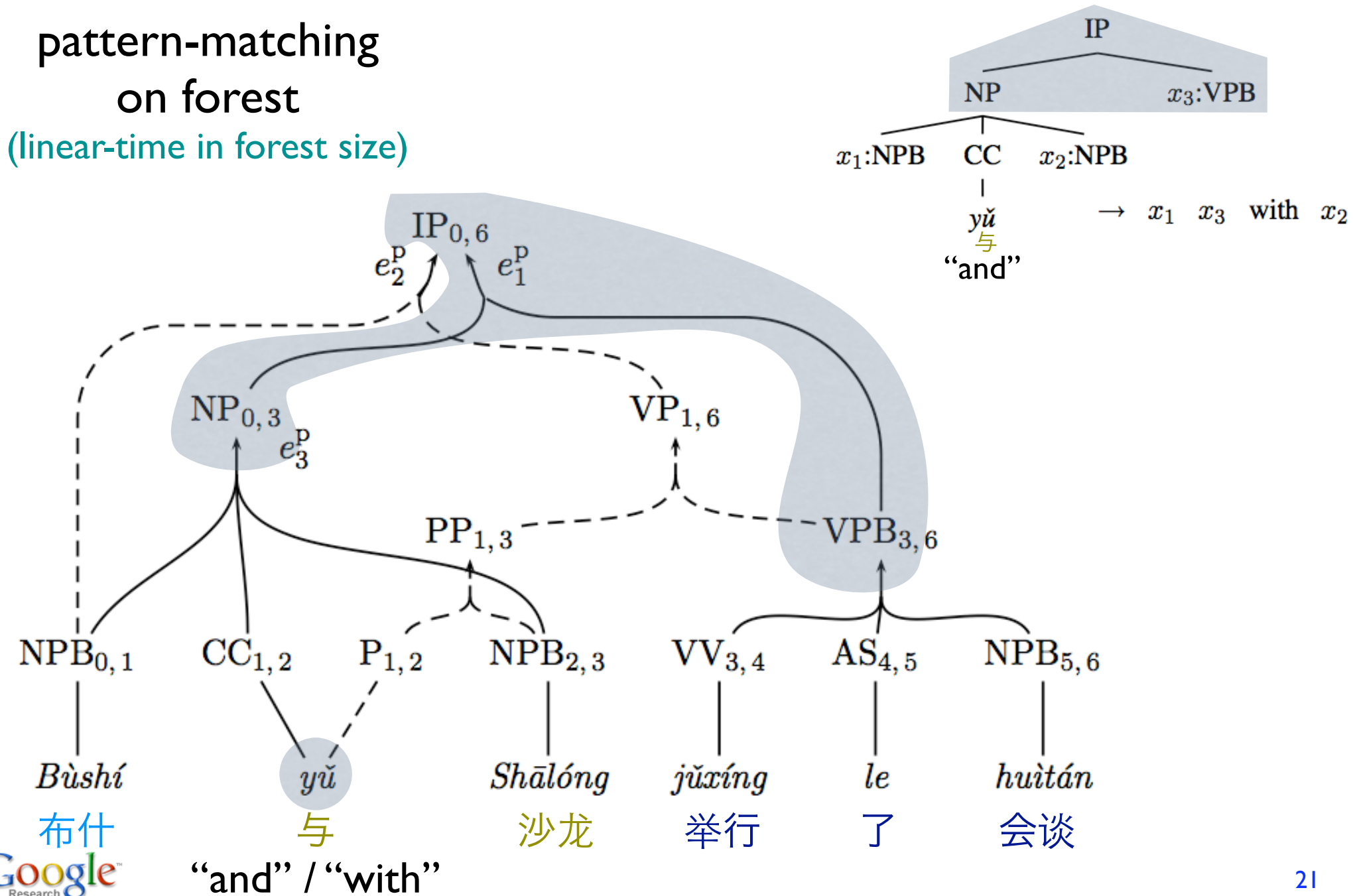


“and” / “with”

Forest-based Translation

pattern-matching
on forest

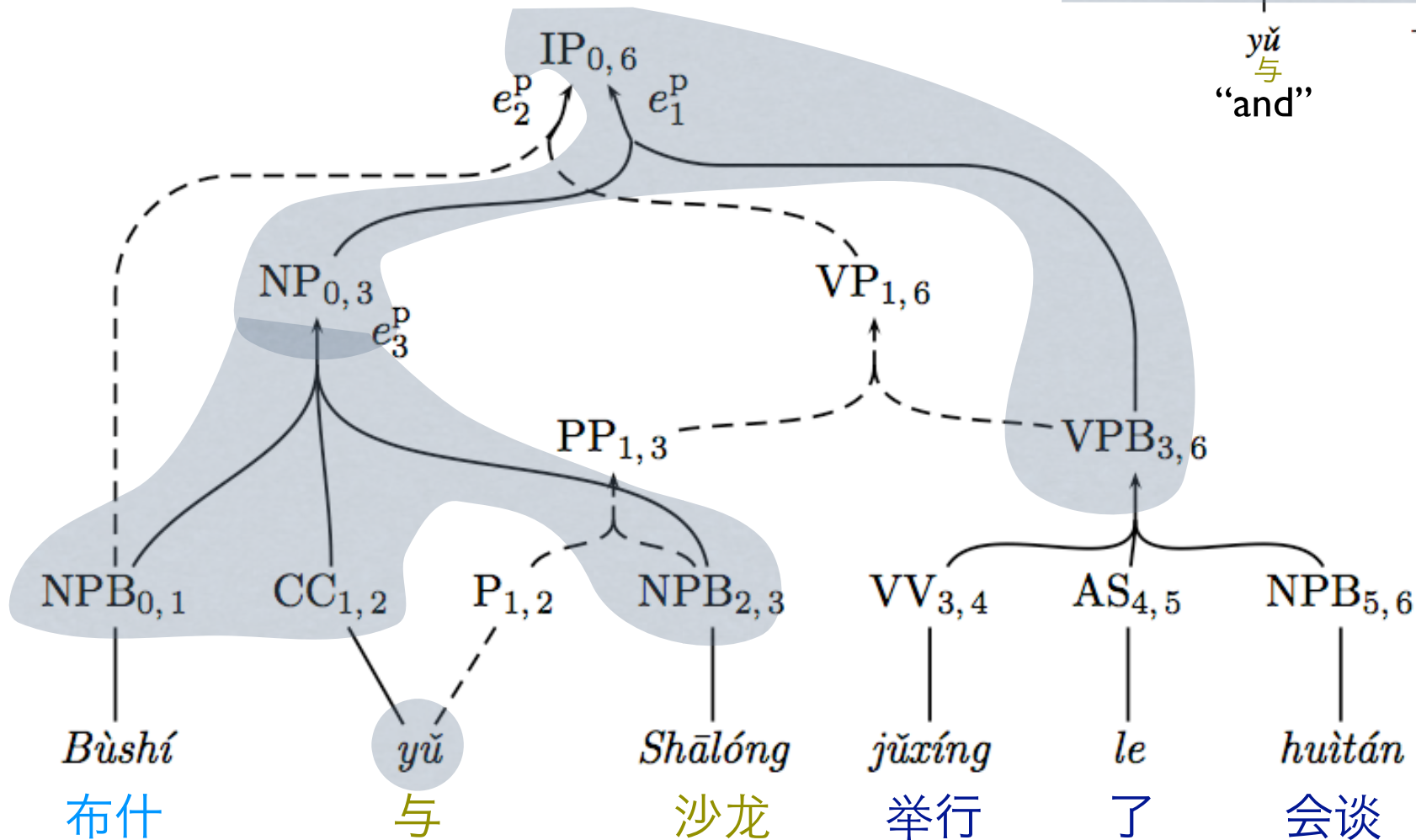
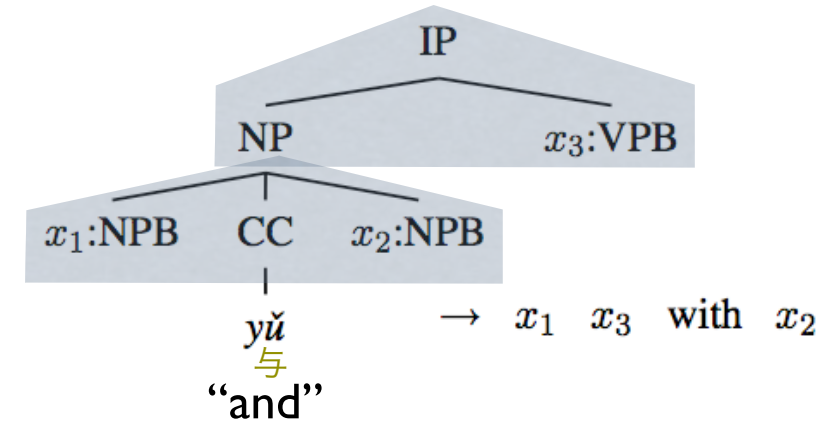
(linear-time in forest size)



Forest-based Translation

pattern-matching
on forest

(linear-time in forest size)

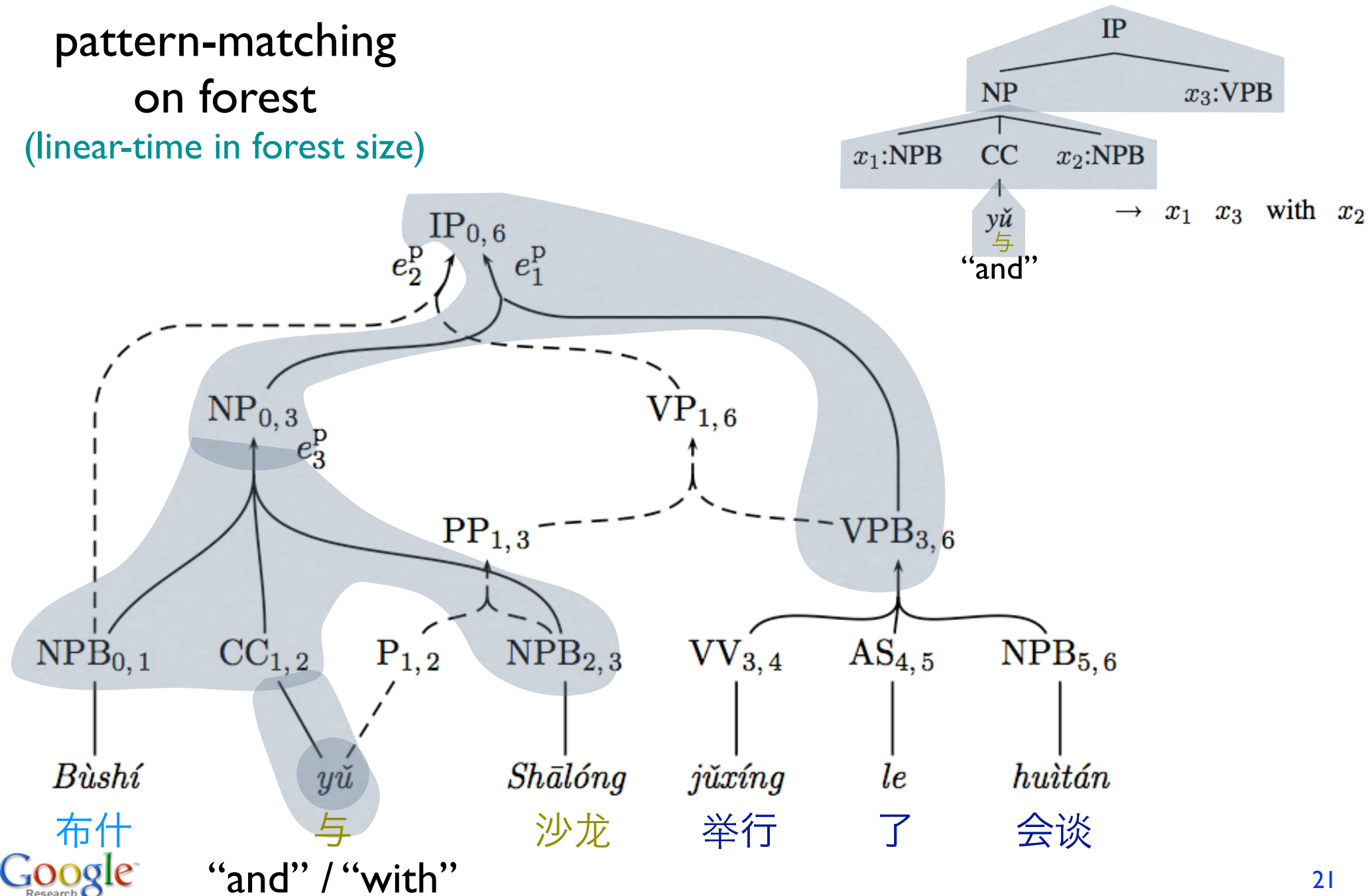


“and” / “with”

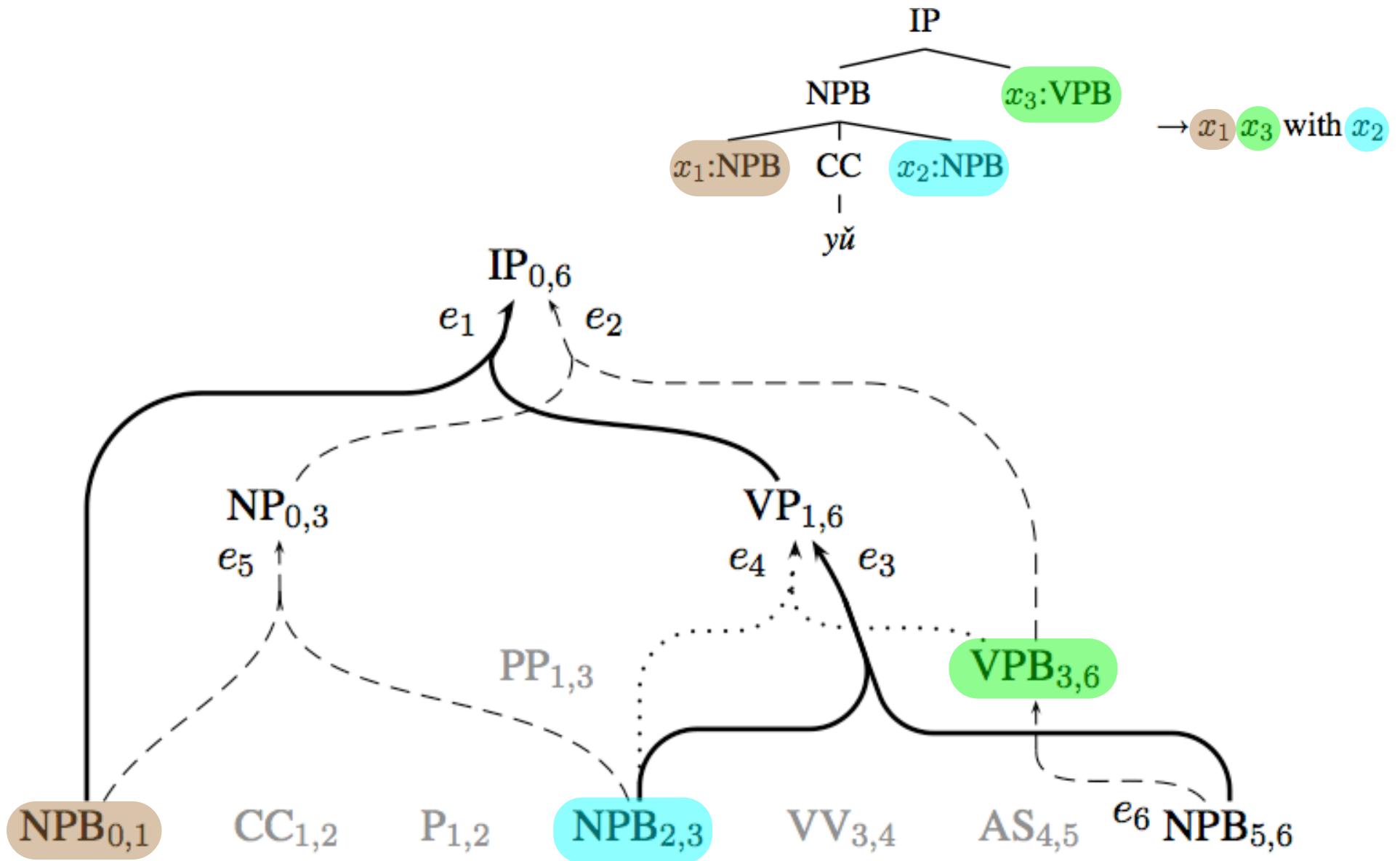
Forest-based Translation

pattern-matching
on forest

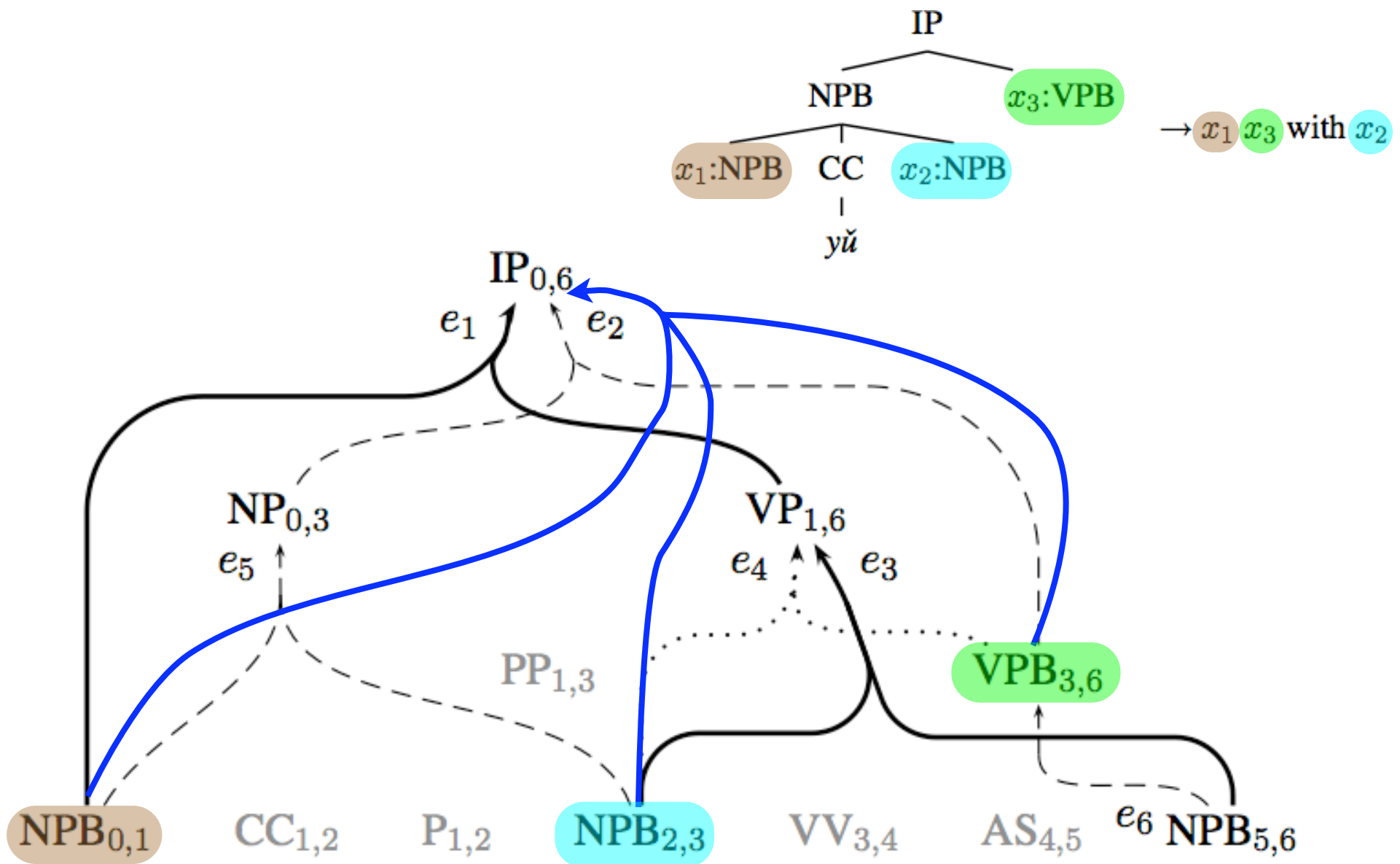
(linear-time in forest size)



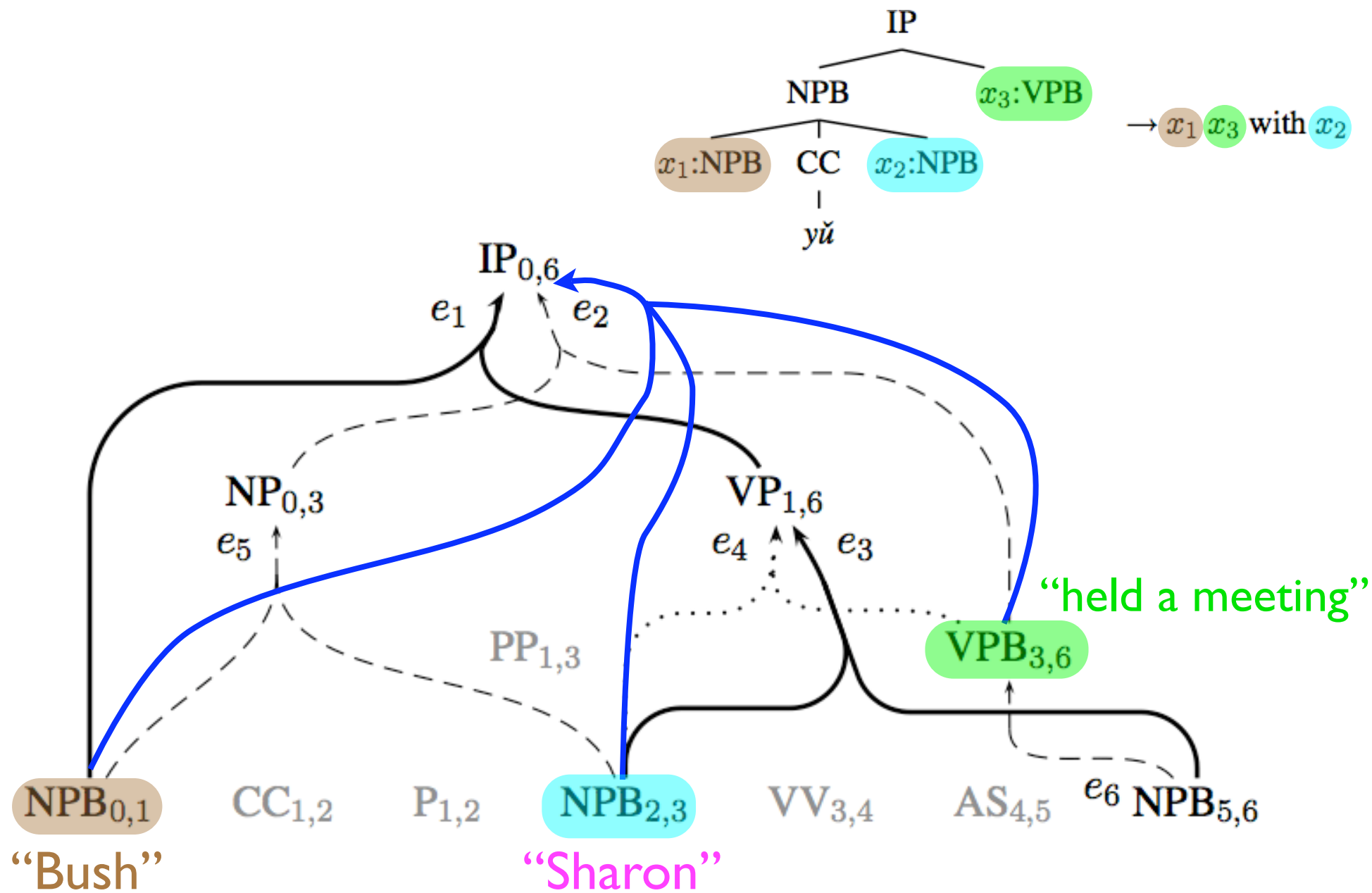
Translation Forest



Translation Forest

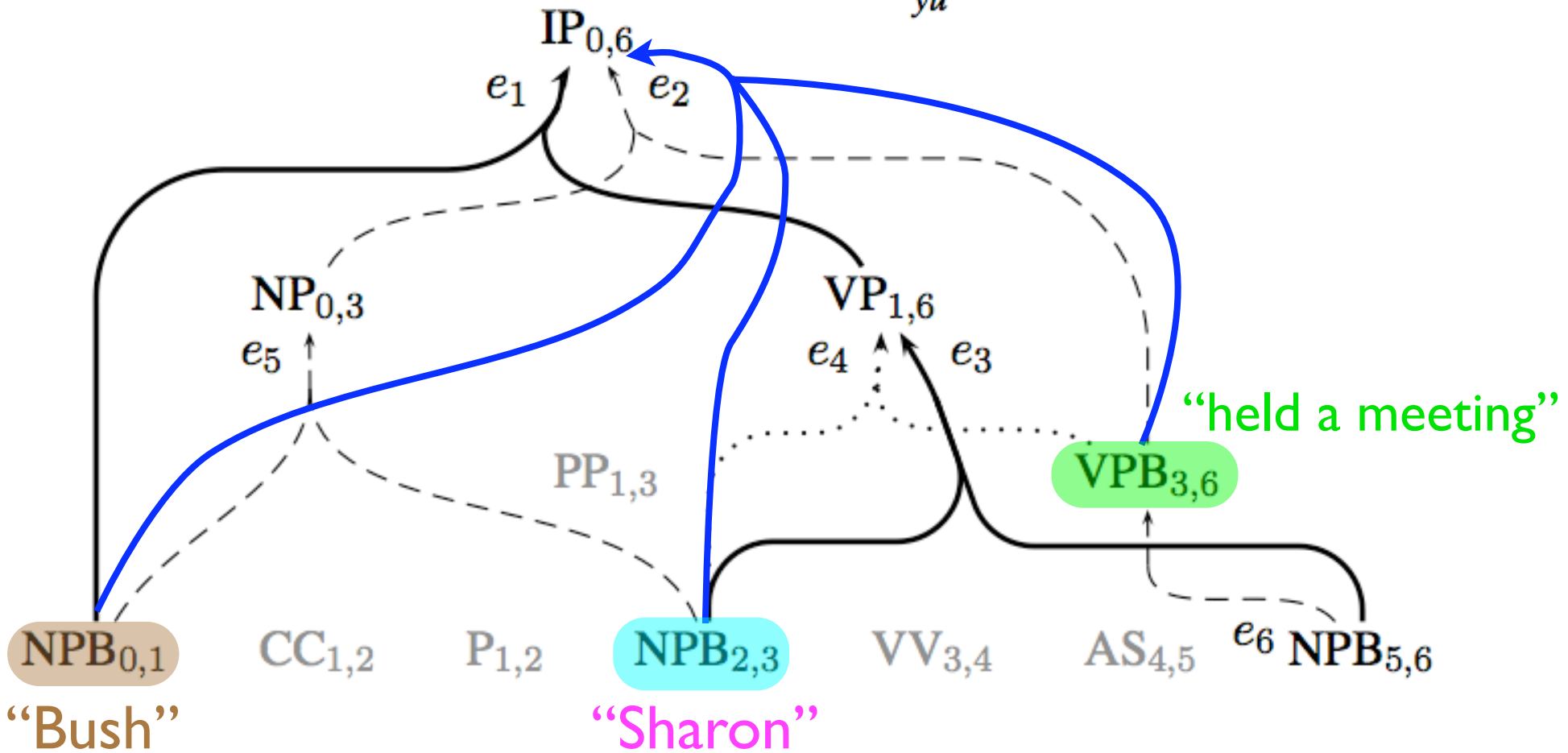
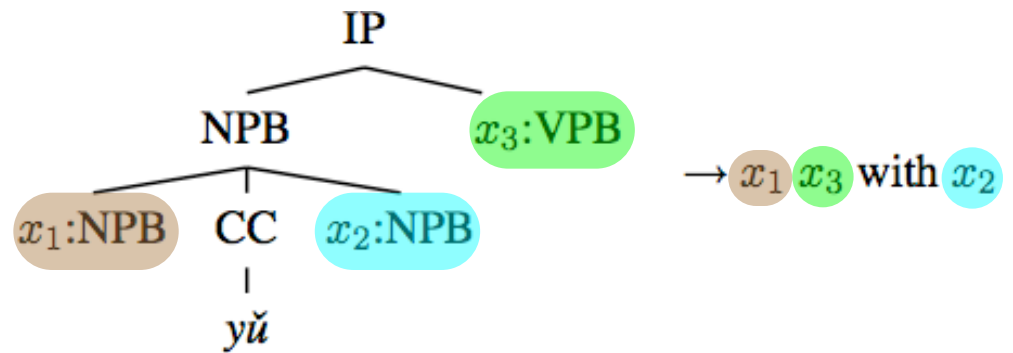


Translation Forest

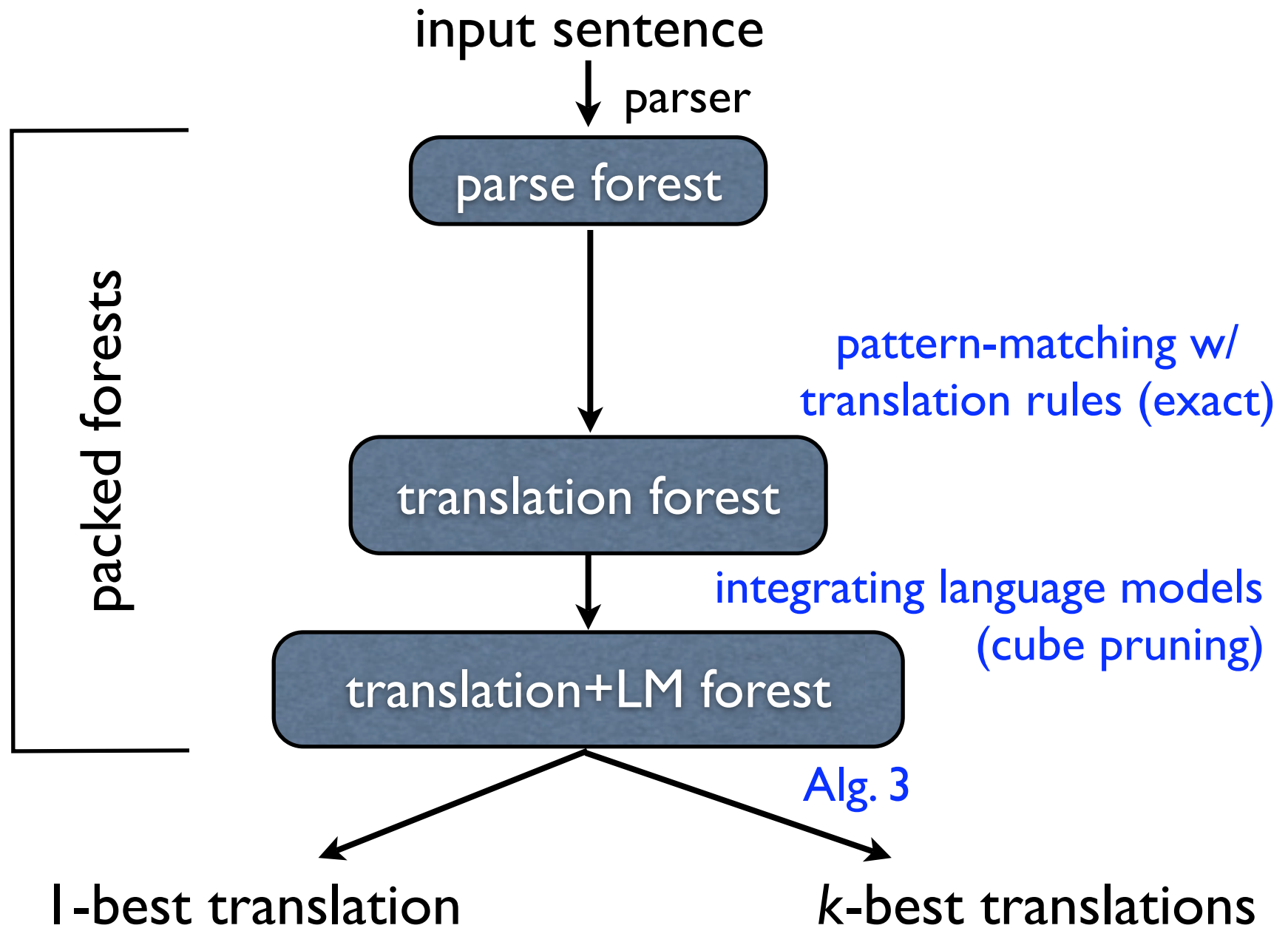


Translation Forest

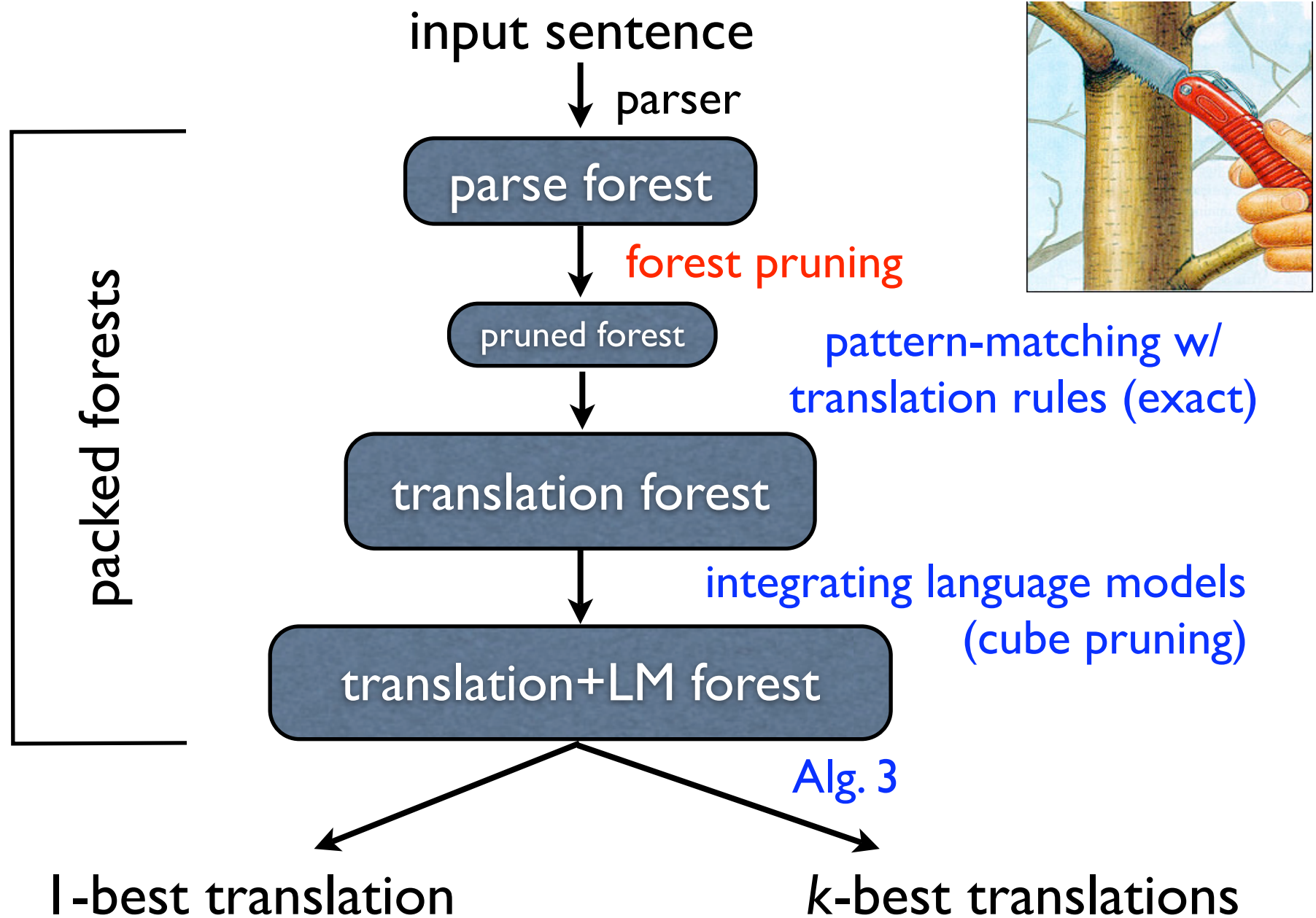
“Bush held a meeting with Sharon”



The Whole Pipeline

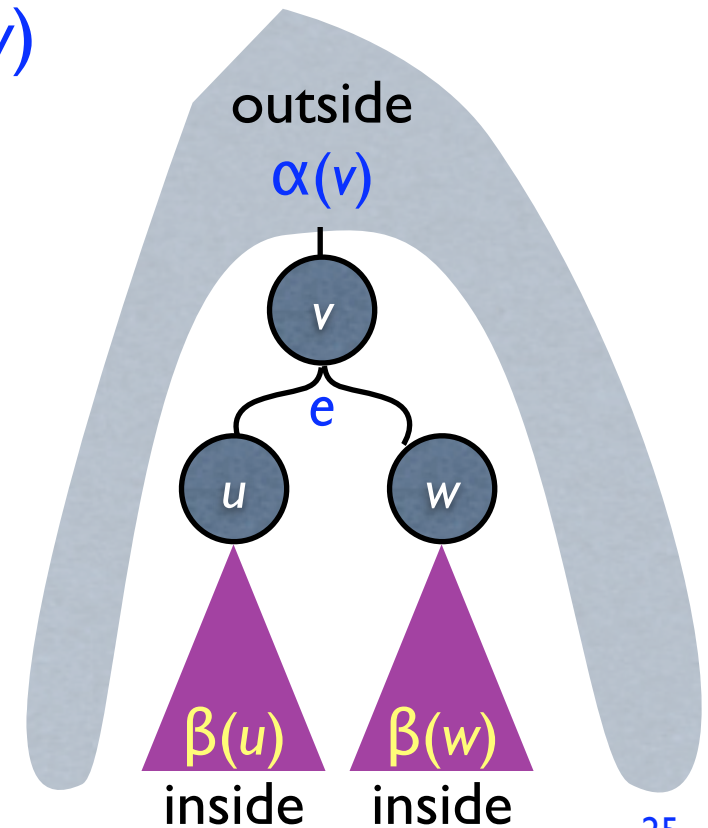


The Whole Pipeline



Parse Forest Pruning

- prune *unpromising* hyperedges
- principled way: inside-outside
 - first compute Viterbi inside β , outside α
- then $\alpha\beta(e) = \alpha(v) + c(e) + \beta(u) + \beta(w)$
 - cost of best deriv that traverses e
 - similar to “expected count” in EM
- prune away hyperedges that have
$$\alpha\beta(e) - \alpha\beta(\text{TOP}) > p$$
for some threshold p

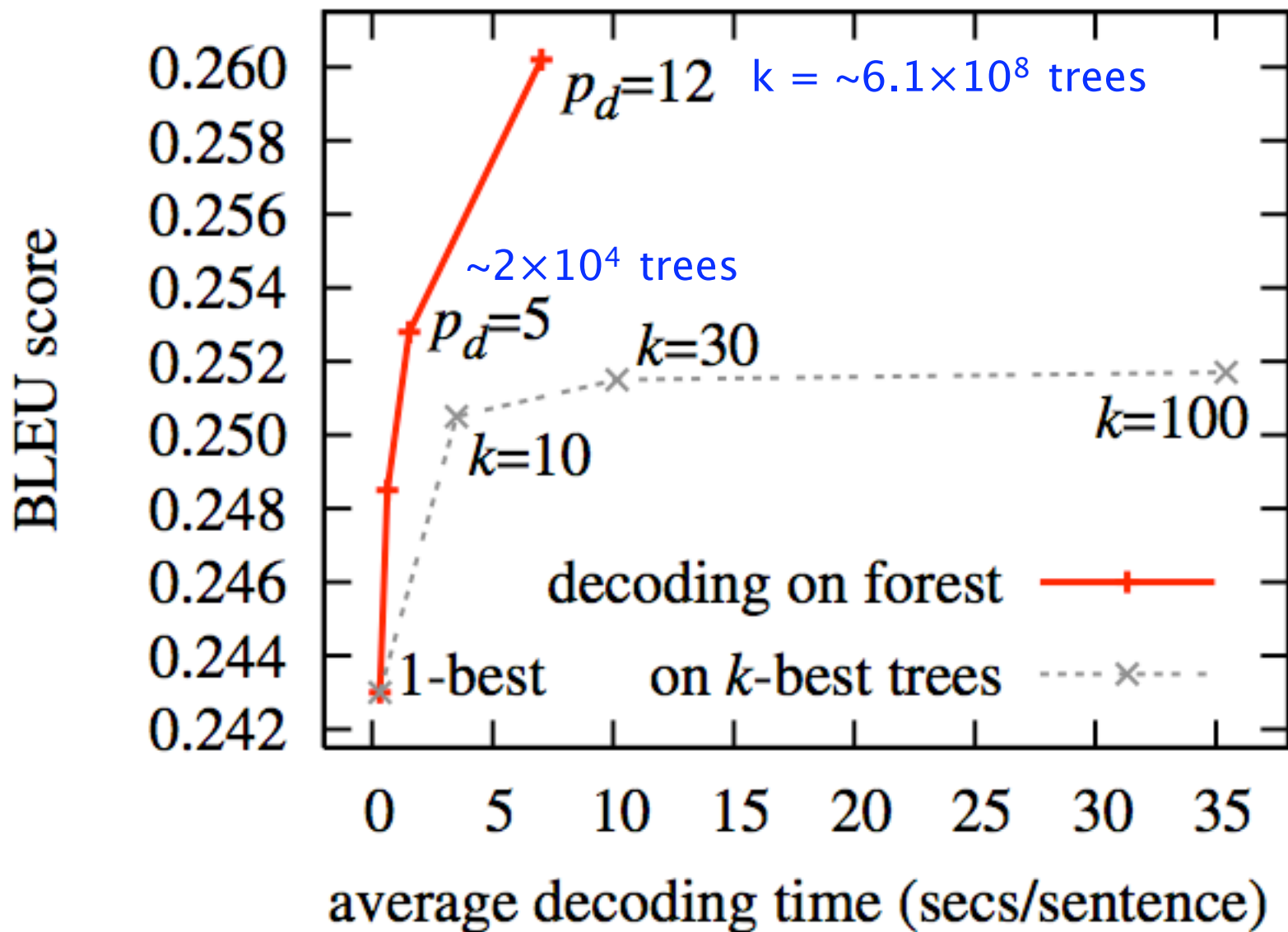


Small-Scale Experiments

- Chinese-to-English translation
 - on a tree-to-string system similar to (Liu et al, 2006)
- 31k sentences pairs (0.8M Chinese & 0.9M English words)
- GIZA++ aligned
- trigram language model trained on the English side
- dev: NIST 2002 (878 sent.); test: NIST 2005 (1082 sent.)
- Chinese-side parsed by the parser of Xiong et al. (2005)
 - modified to output a forest for each sentence (Huang 2008)
- BLEU score: I-best baseline: 0.2430 vs. Pharaoh: 0.2297

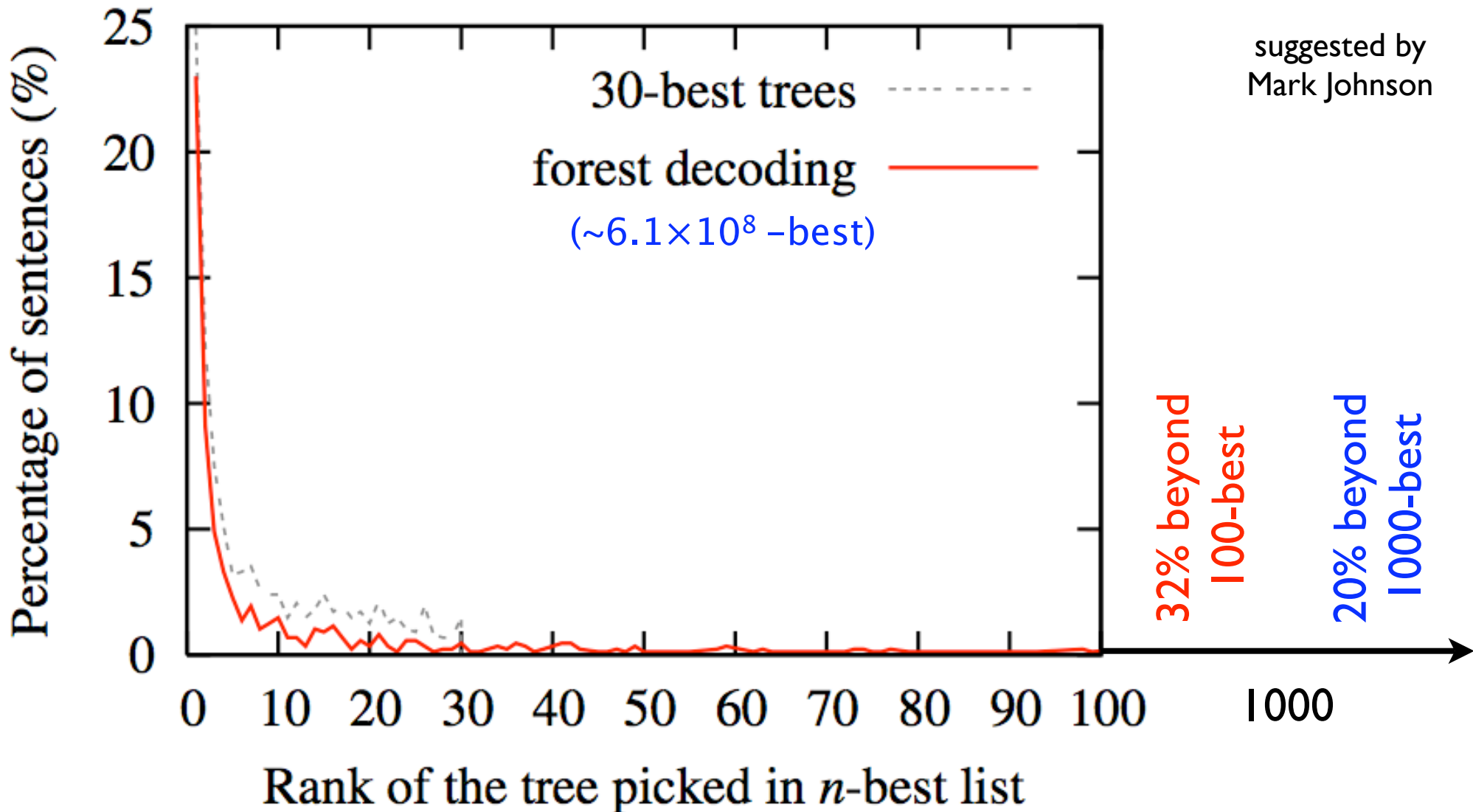
k-best trees vs. forest-based

1.7 Bleu improvement over 1-best,
0.8 over 30-best, and even faster!



forest as virtual ∞ -best list

- how often is the i th-best tree picked by the decoder?



wait a sec... where are the rules from?



wait a sec... where are the rules from?

xiǎoxīn

小心 X \Leftrightarrow be careful not to X



wait a sec... where are the rules from?

xiǎoxīn gǒu
小心 狗 \Leftrightarrow be aware of dog



xiǎoxīn X
小心 X \Leftrightarrow be careful not to X



wait a sec... where are the rules from?

小心 VP \Leftrightarrow be careful **not to** VP

小心 NP \Leftrightarrow be careful **of** NP

...

xiǎoxīn gǒu

小心 狗 \Leftrightarrow be aware of **dog**

xiǎoxīn

小心 X \Leftrightarrow be careful not to X

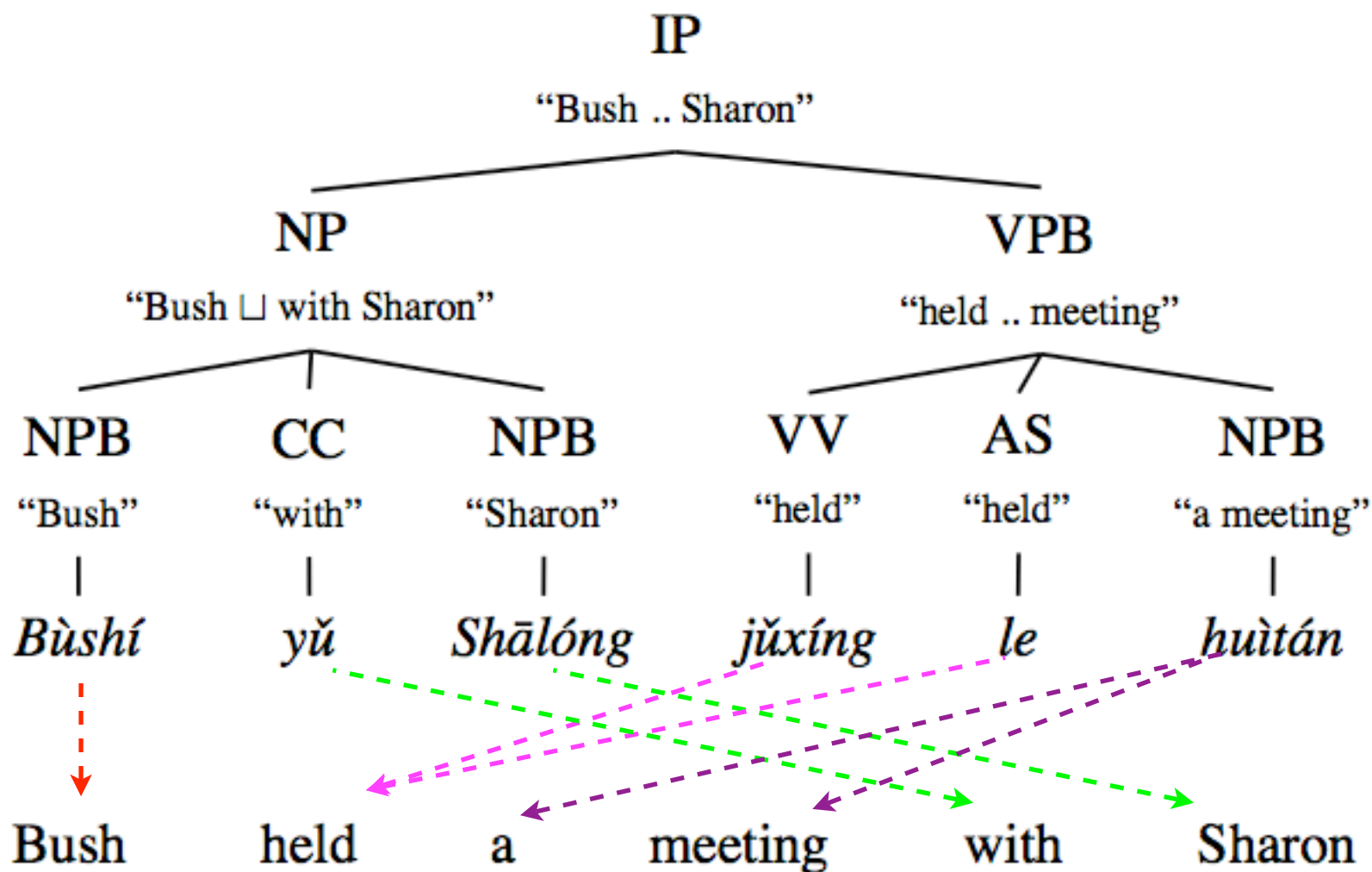


Outline

- Overview: Tree-based Translation
- Forest-based Translation
- **Forest-based Rule Extraction**
 - background: tree-based rule extraction (Galley et al., 2004)
 - extension to forest-based
 - large-scale experiments

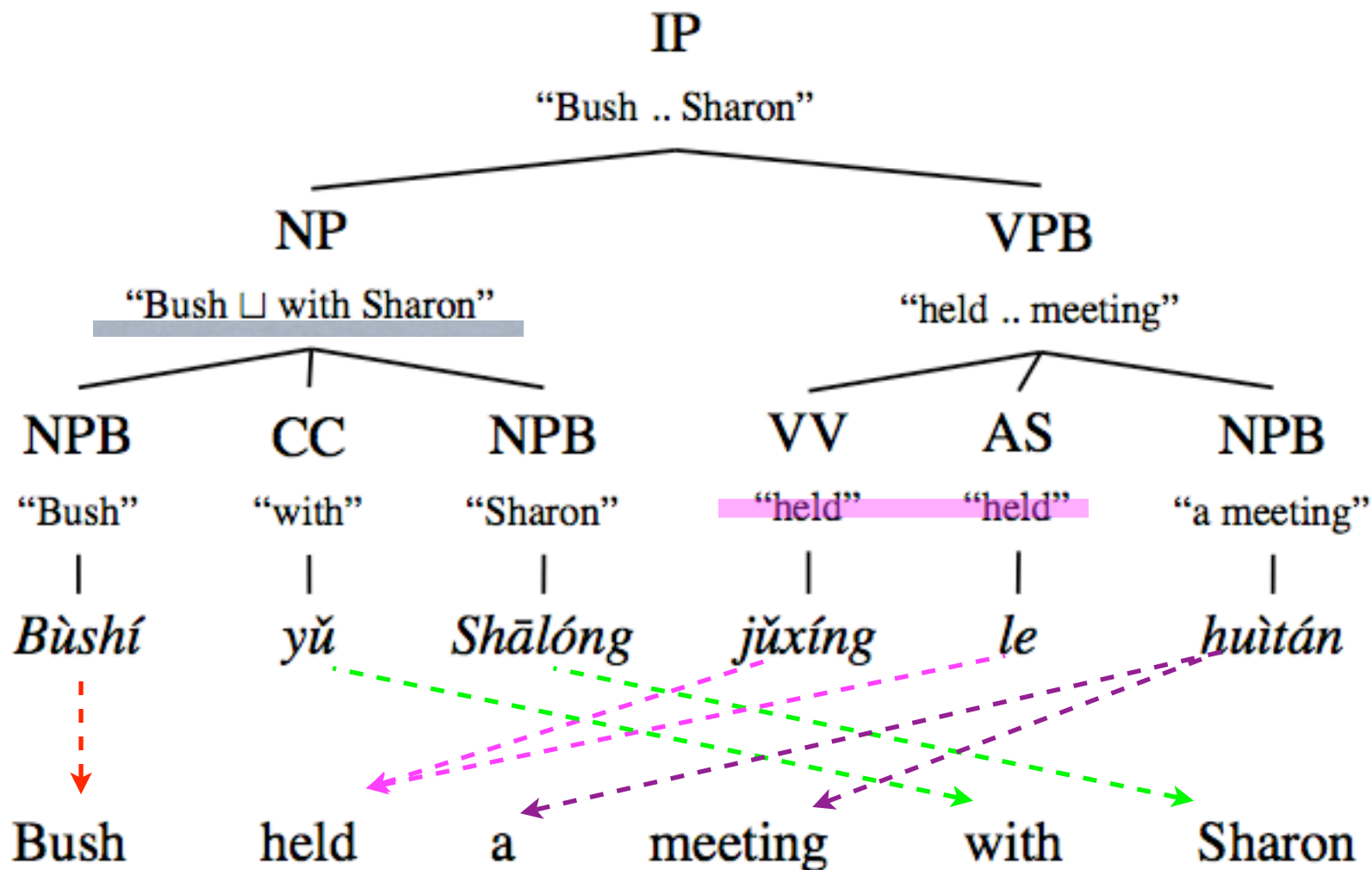
Where are the rules from?

- source parse tree, target sentence, and alignment
- compute target spans



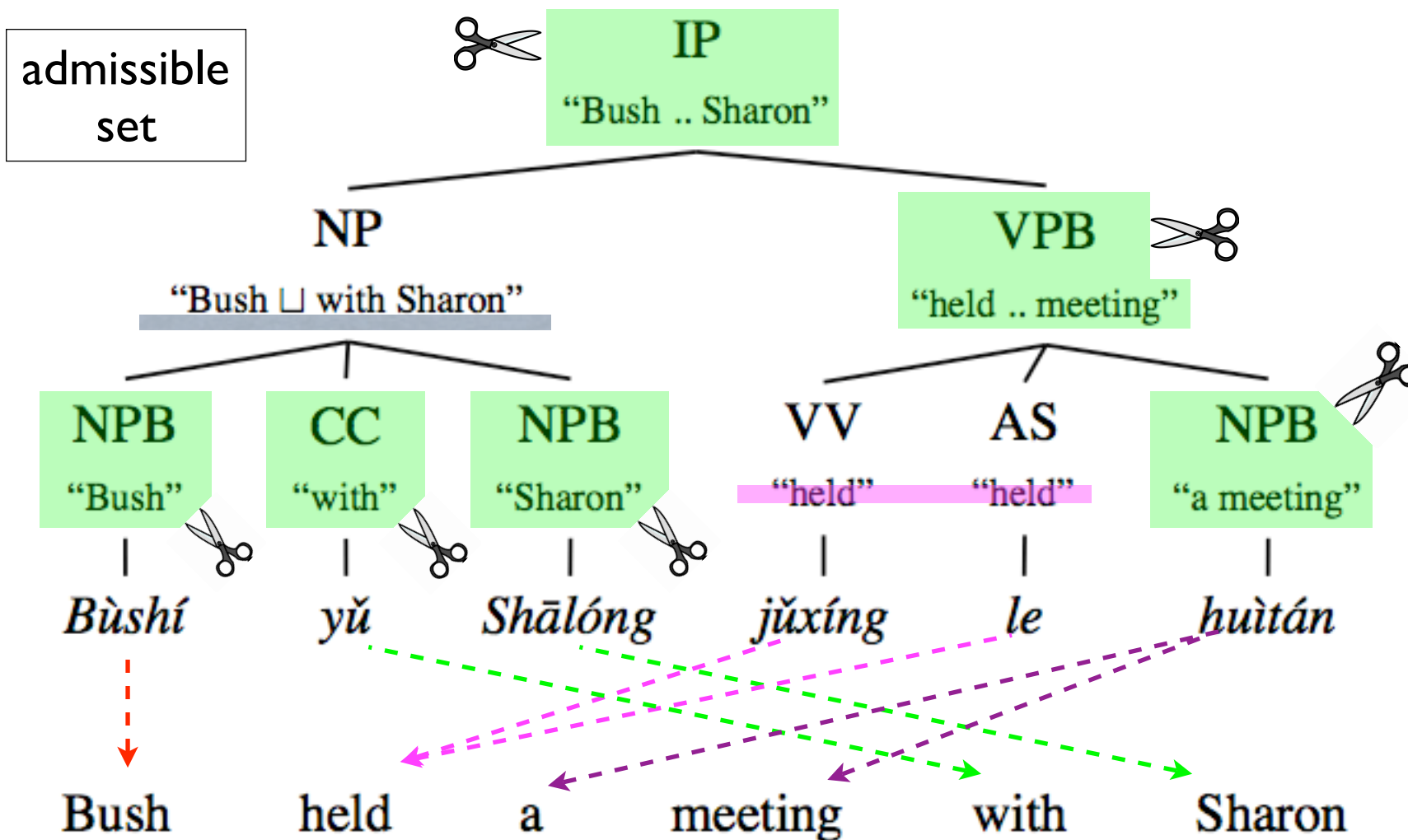
Where are the rules from?

- source parse tree, target sentence, and alignment
- well-formed fragment: contiguous and faithful t-span



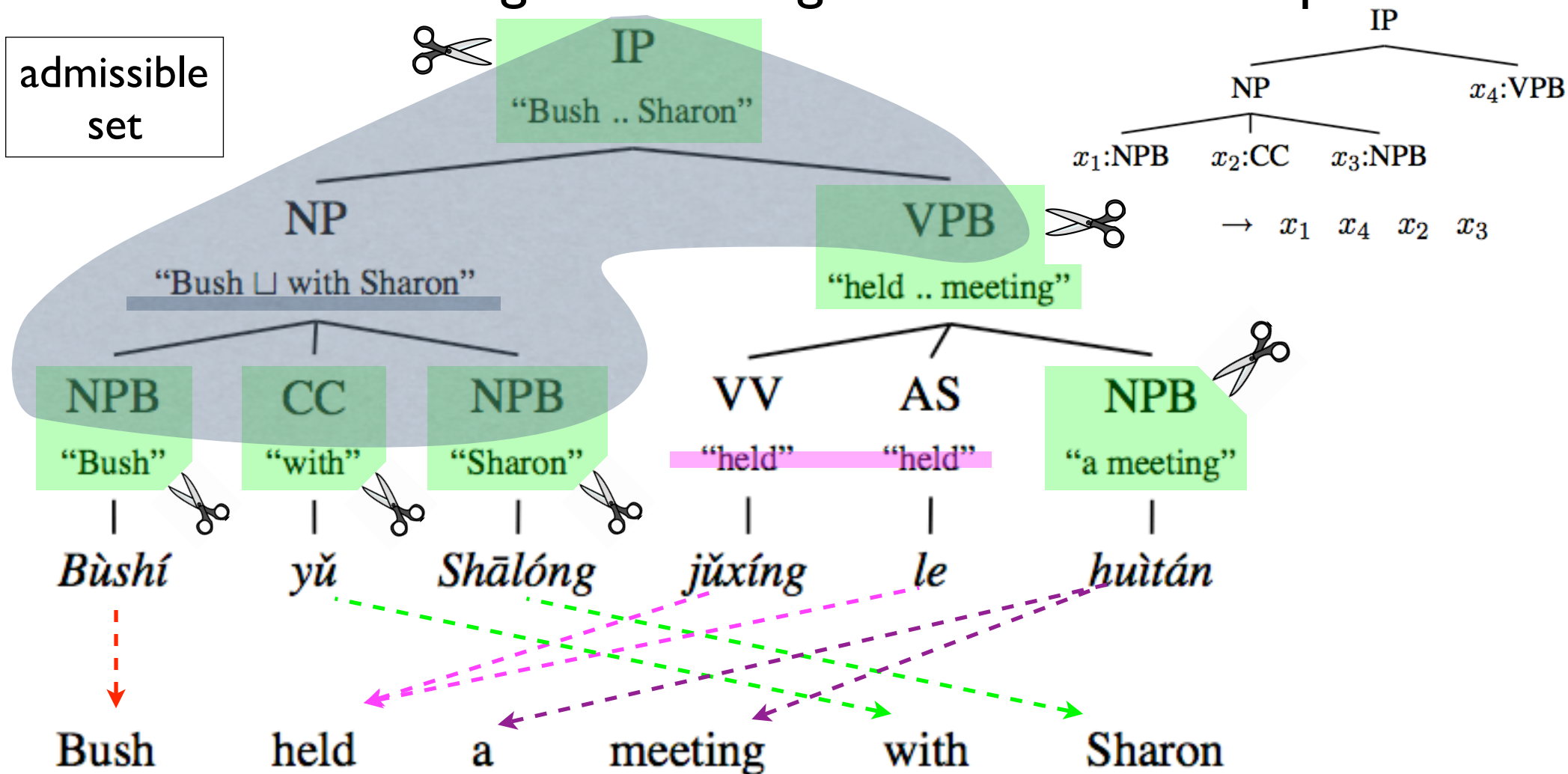
Where are the rules from?

- source parse tree, target sentence, and alignment
- well-formed fragment: contiguous and faithful t-span



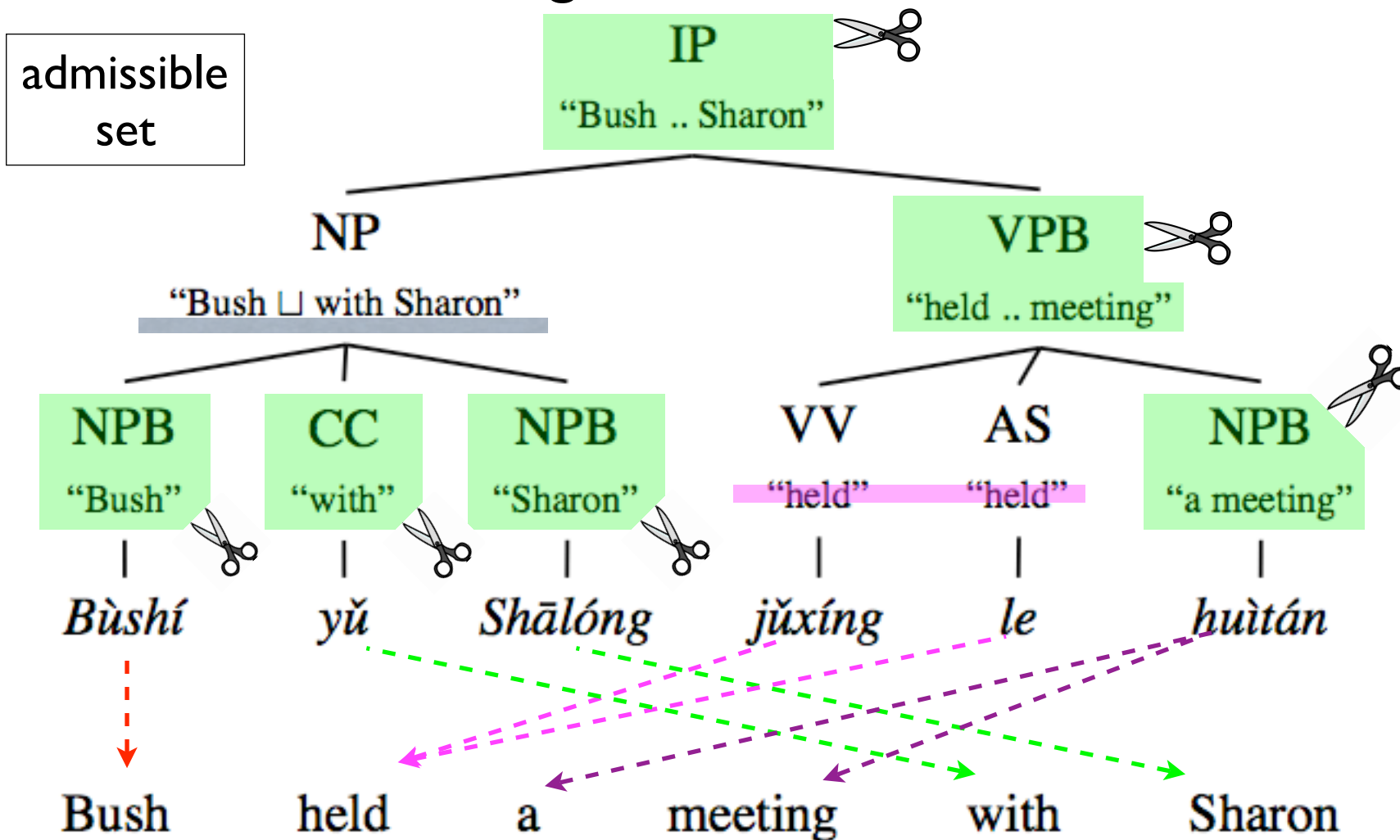
Where are the rules from?

- source parse tree, target sentence, and alignment
- well-formed fragment: contiguous and faithful t-span



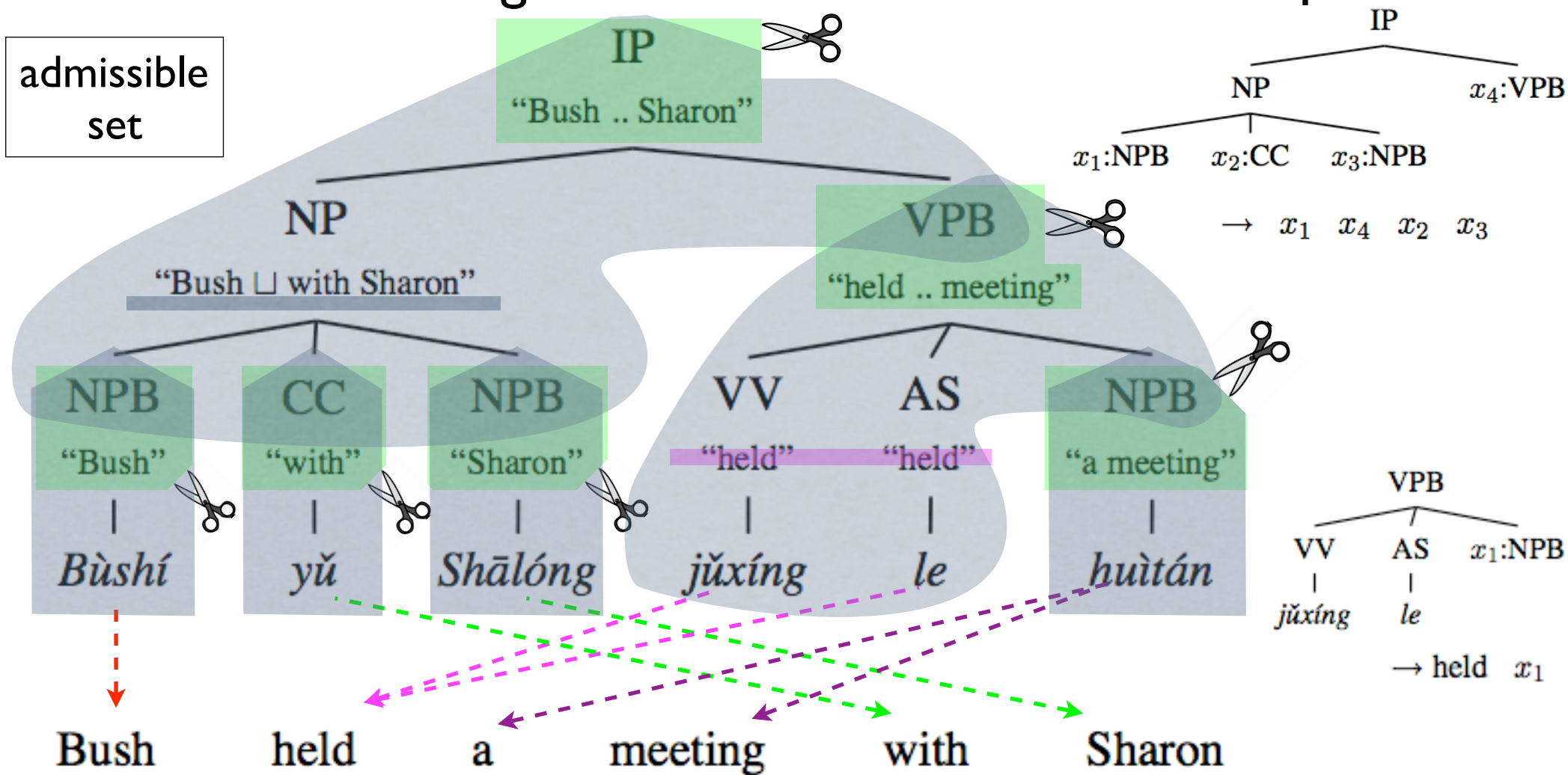
Where are the rules from?

- source parse tree, target sentence, and alignment
- well-formed fragment: contiguous and faithful t-span



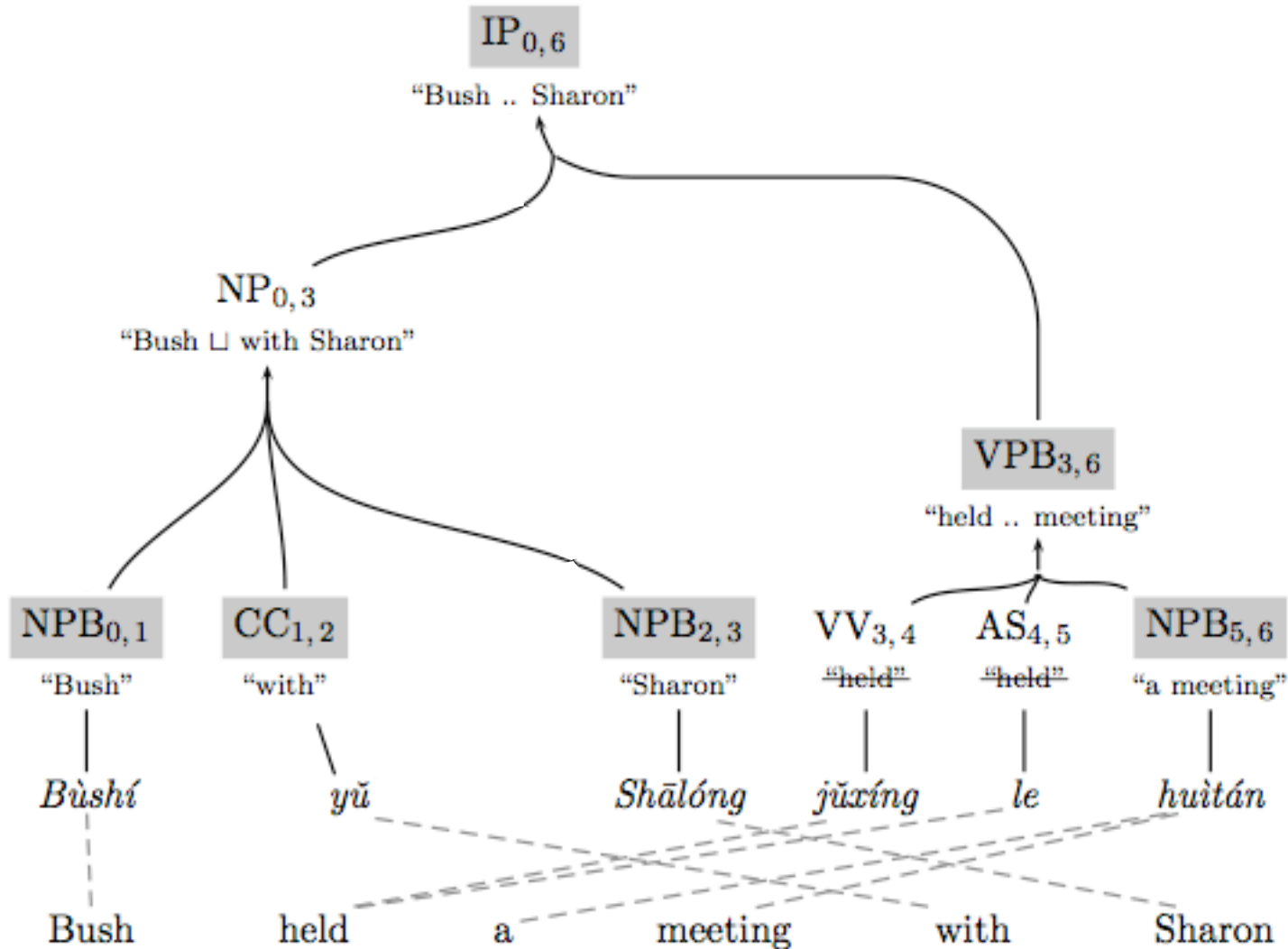
Where are the rules from?

- source parse tree, target sentence, and alignment
- well-formed fragment: contiguous and faithful t-span



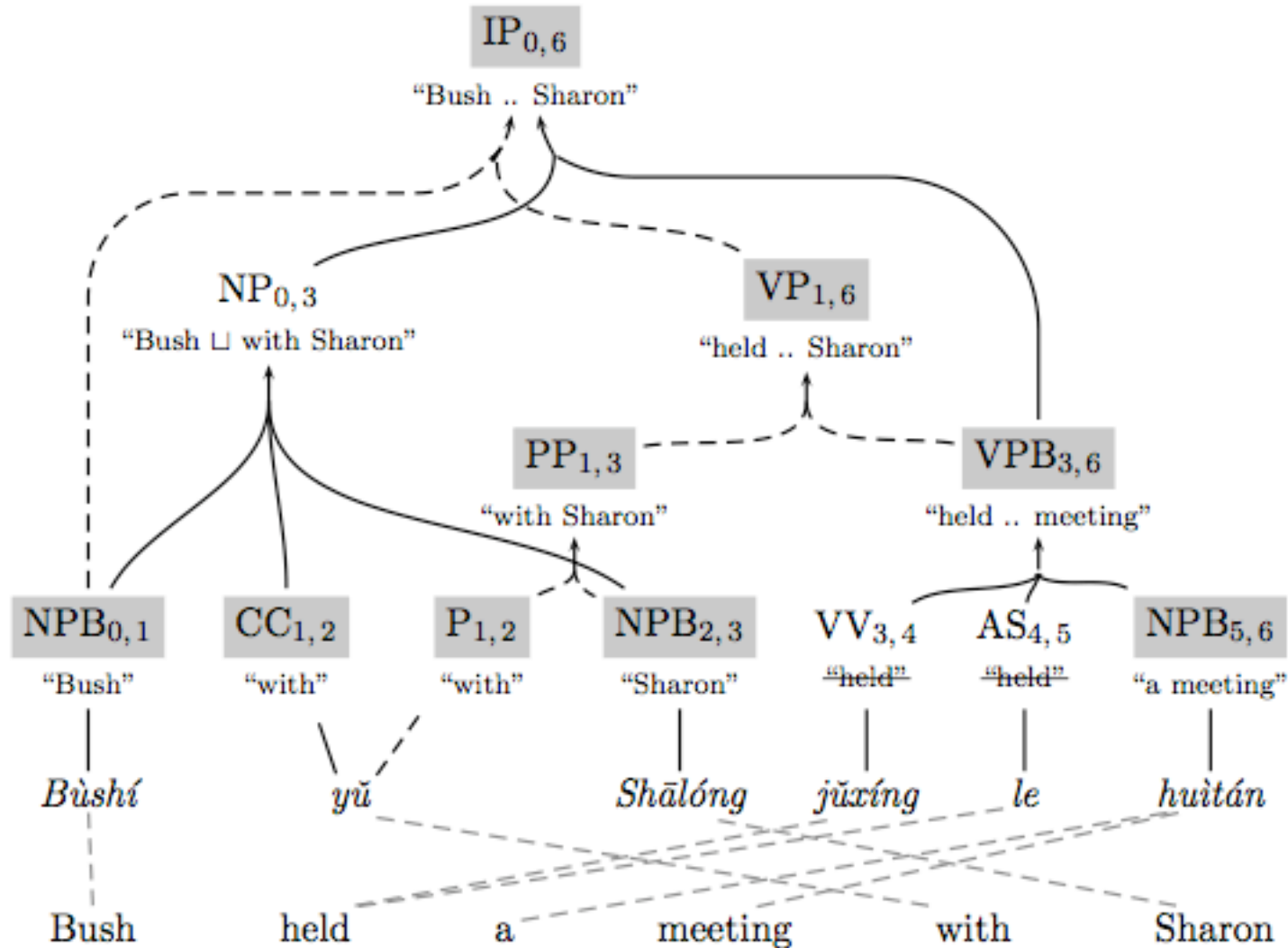
Forest-based Rule Extraction

- same cut set computation; different fragmentation



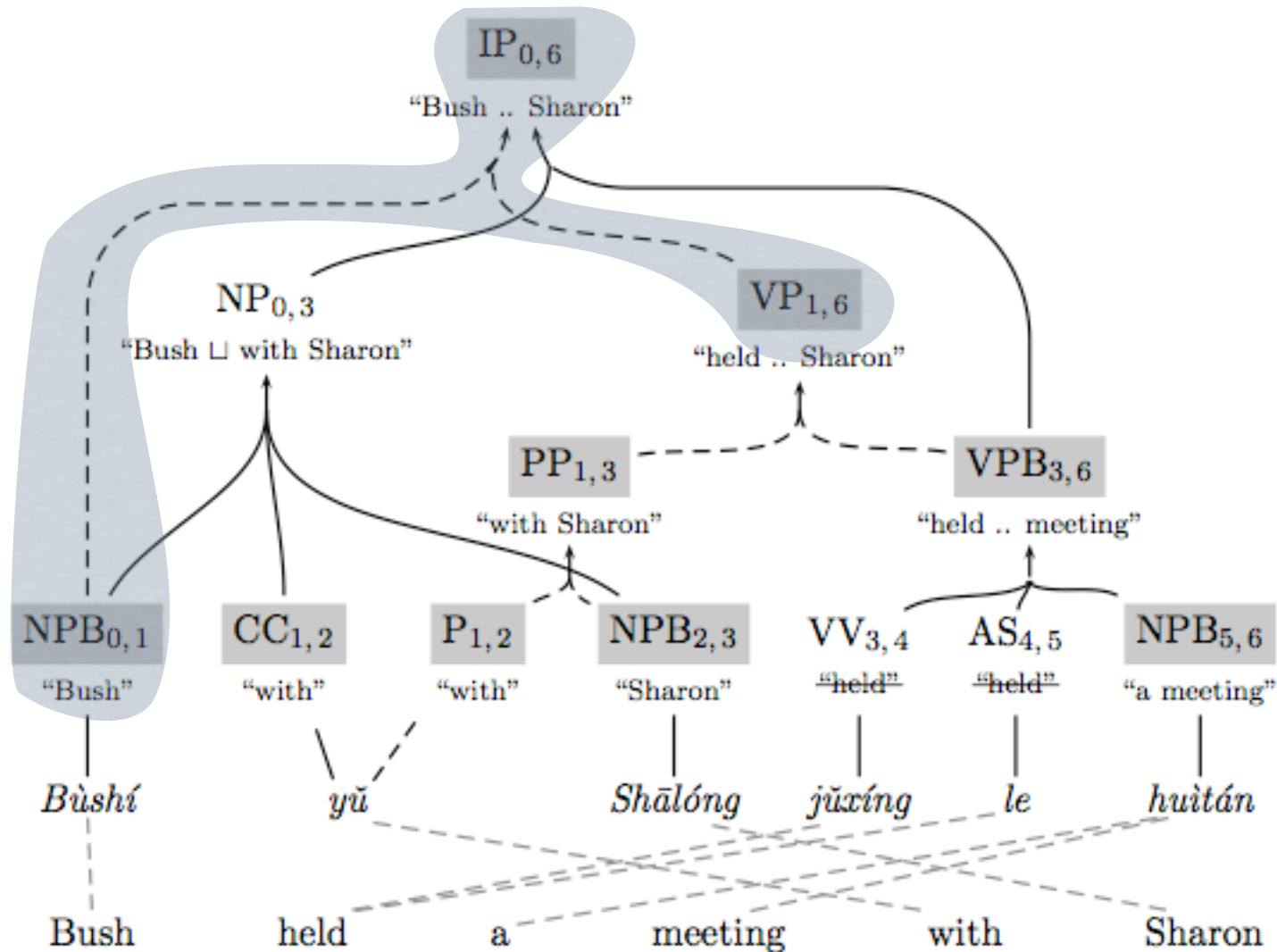
Forest-based Rule Extraction

- same cut set computation; different fragmentation



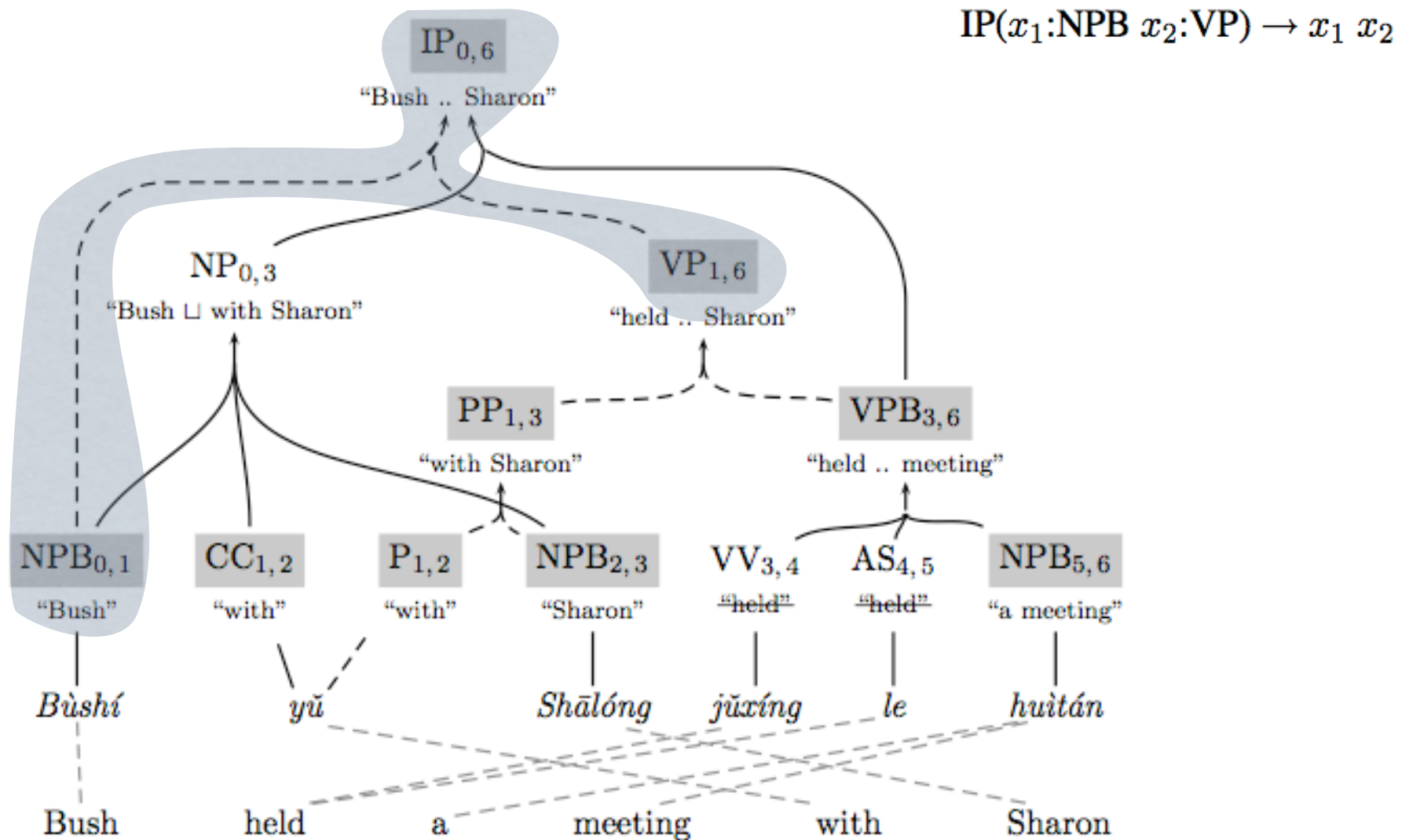
Forest-based Rule Extraction

- same cut set computation; different fragmentation



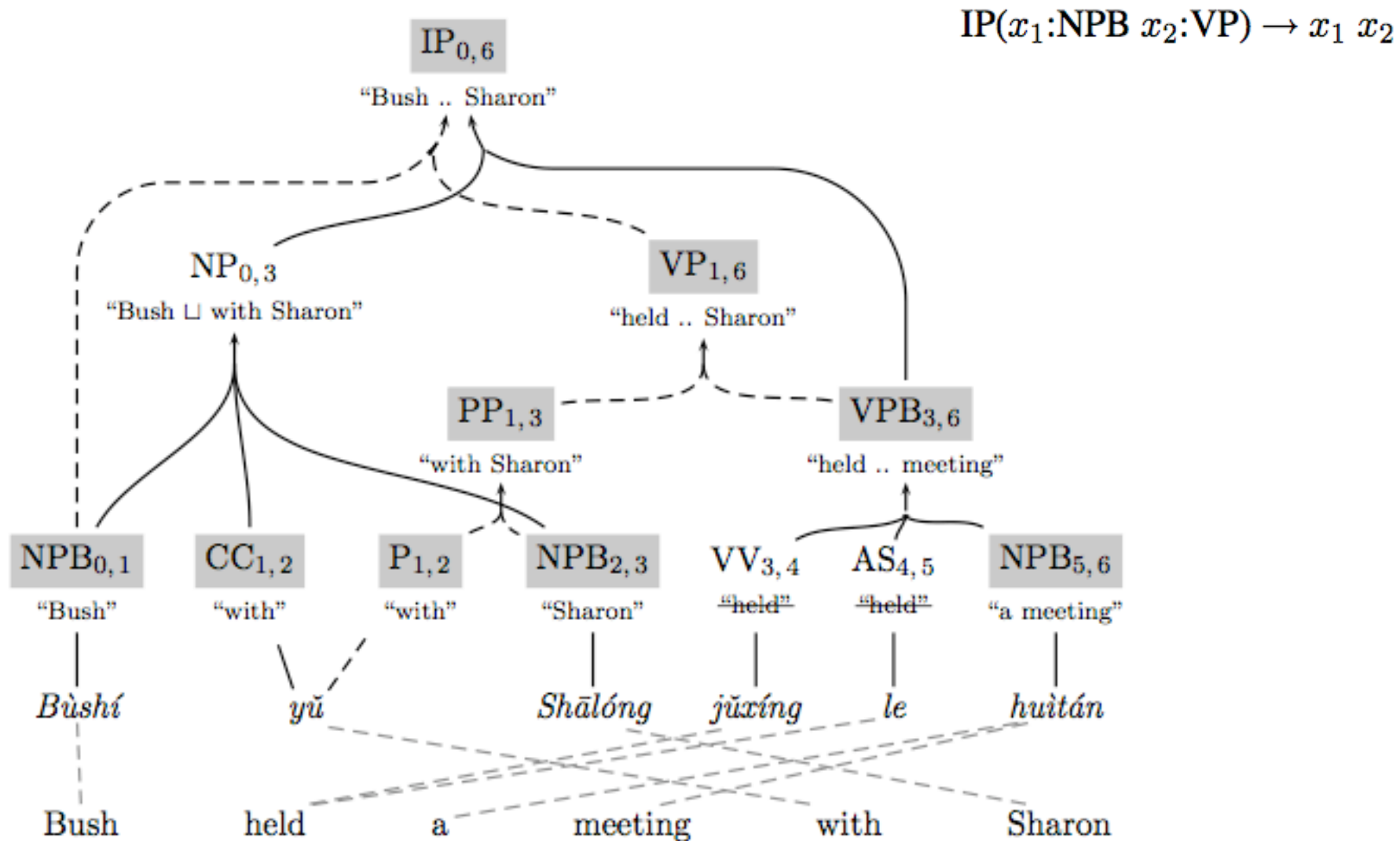
Forest-based Rule Extraction

- same cut set computation; different fragmentation



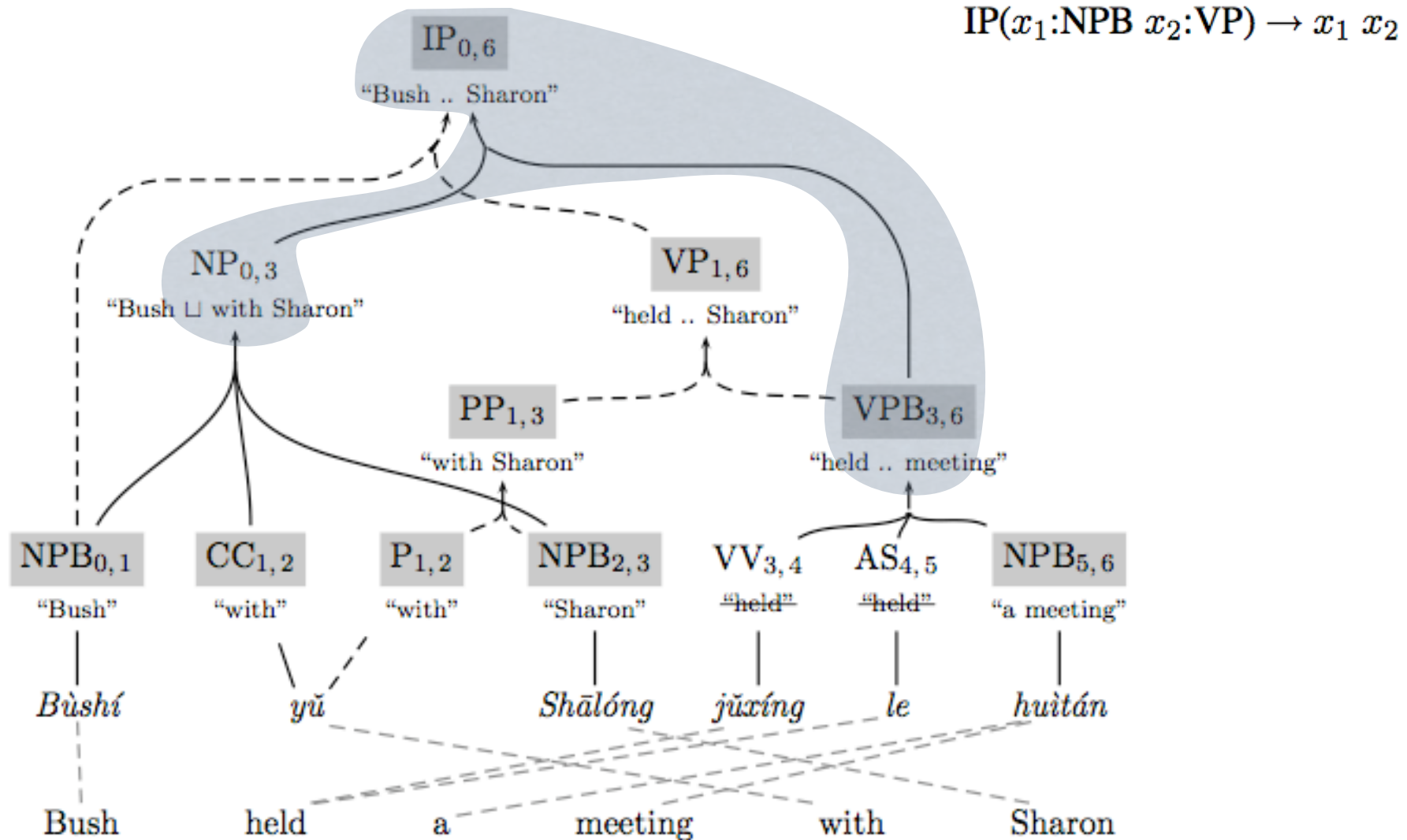
Forest-based Rule Extraction

- same admissible set definition; different fragmentation



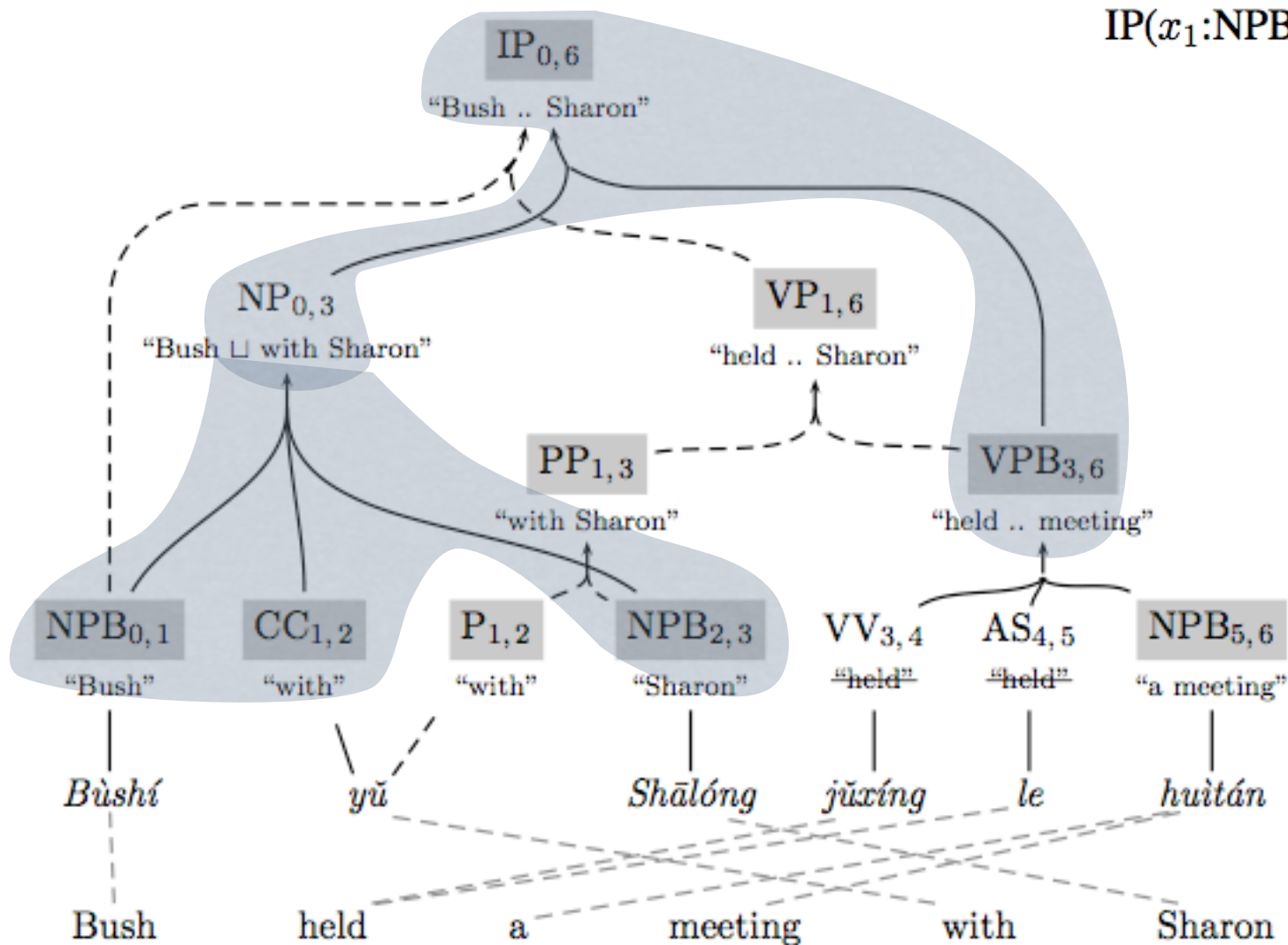
Forest-based Rule Extraction

- same admissible set definition; different fragmentation



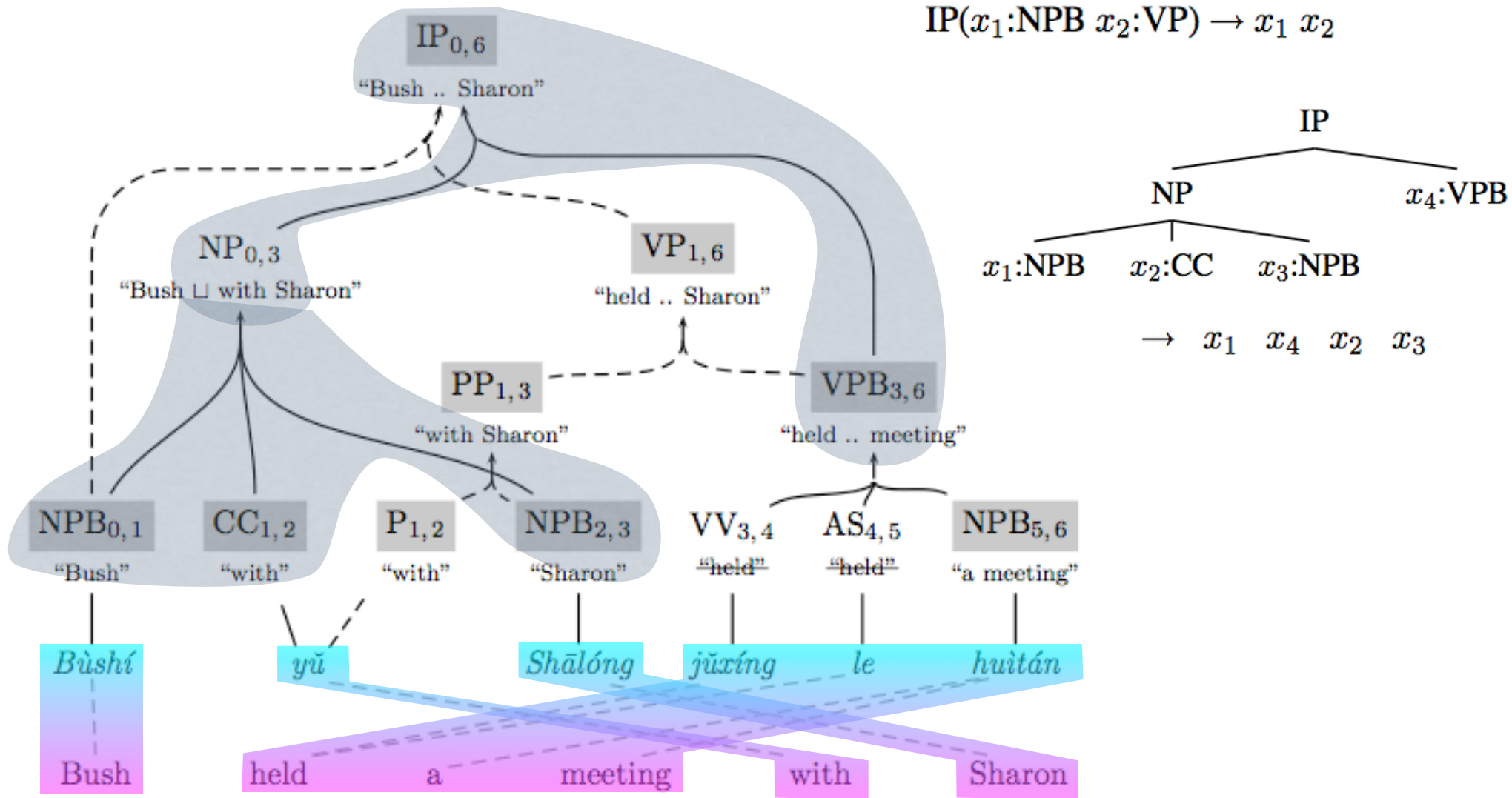
Forest-based Rule Extraction

- same admissible set definition; different fragmentation



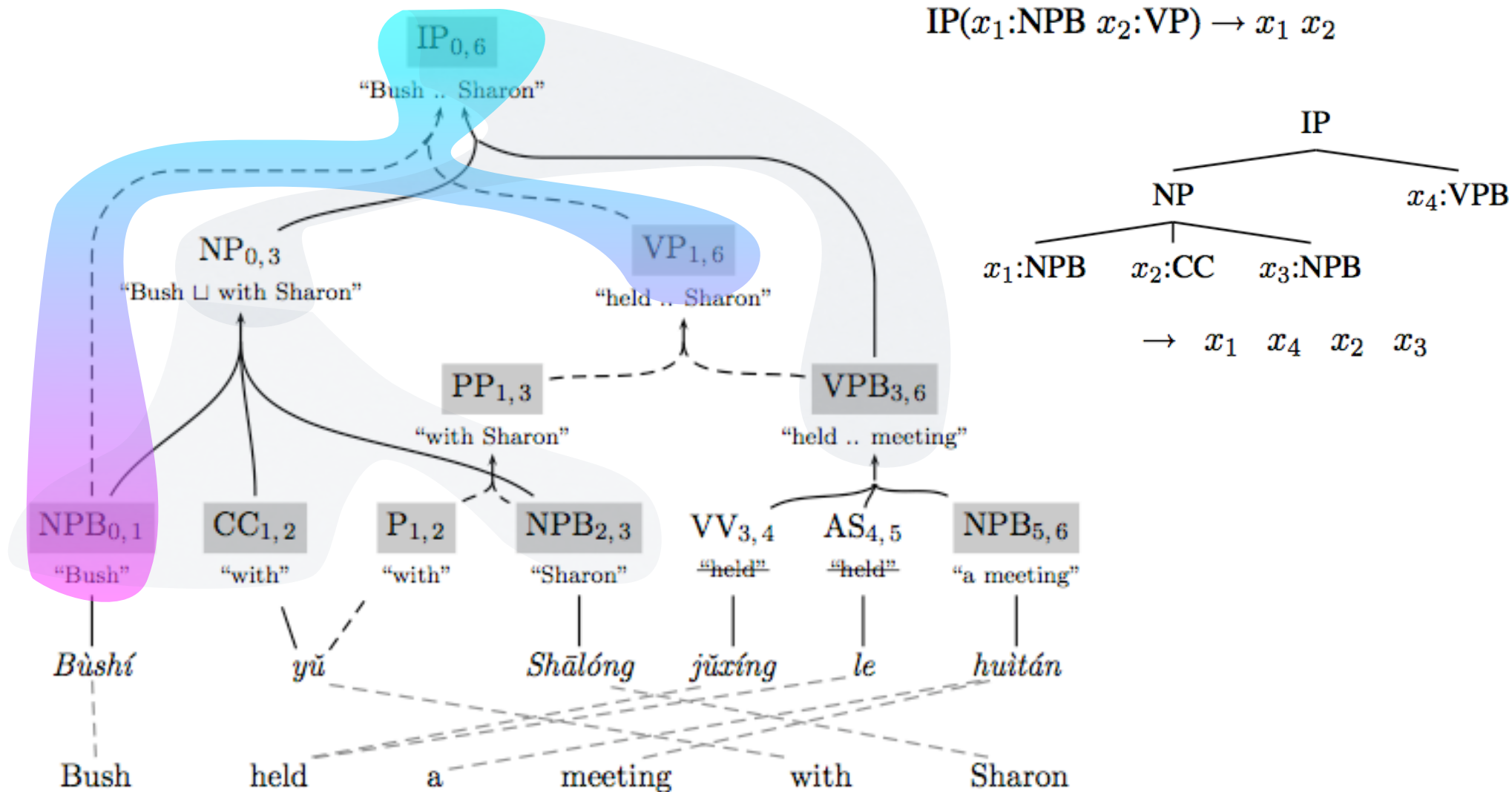
Forest-based Rule Extraction

- same admissible set definition; different fragmentation



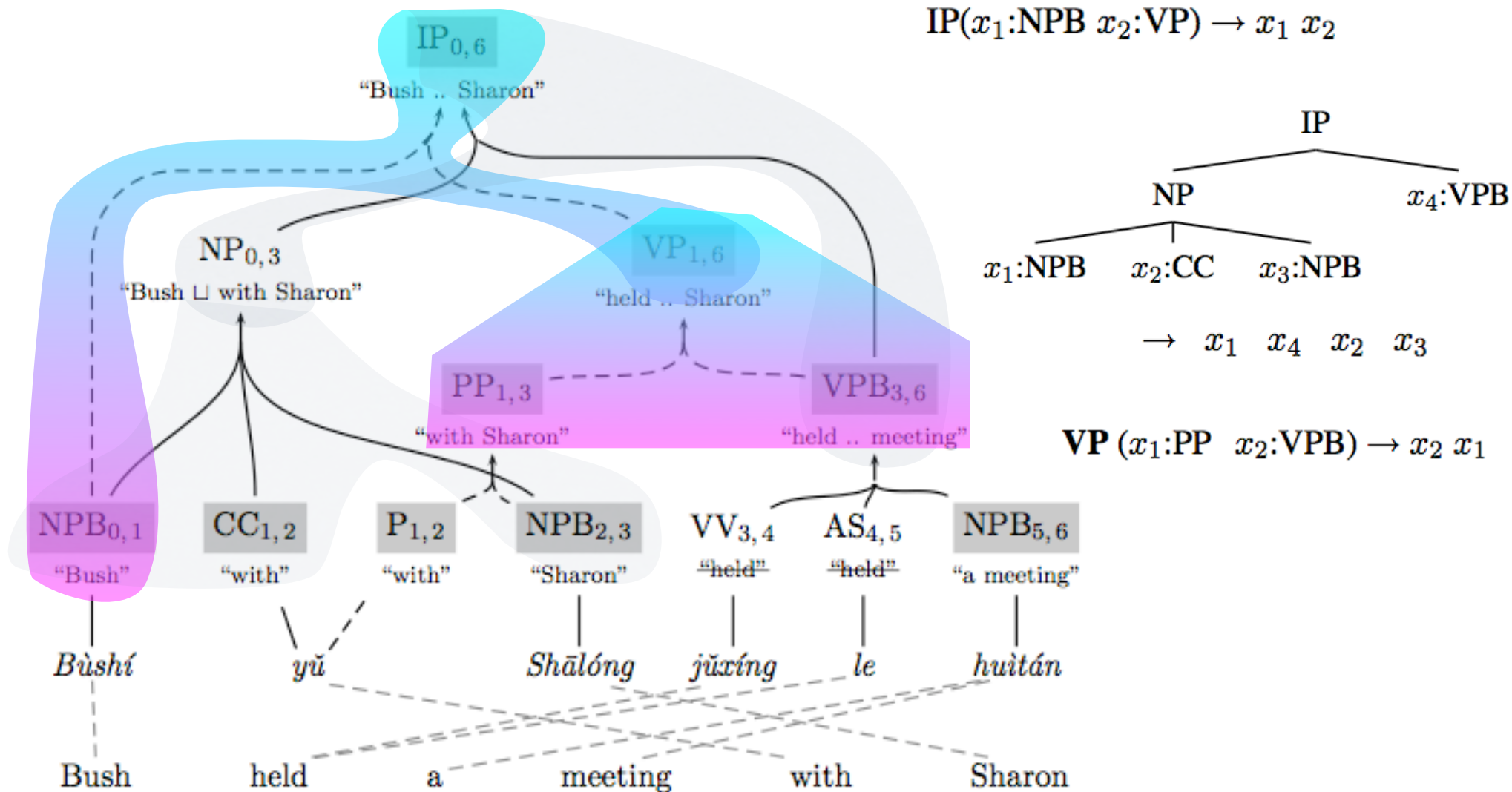
Forest-based Rule Extraction

- forest can extract smaller chunks of rules



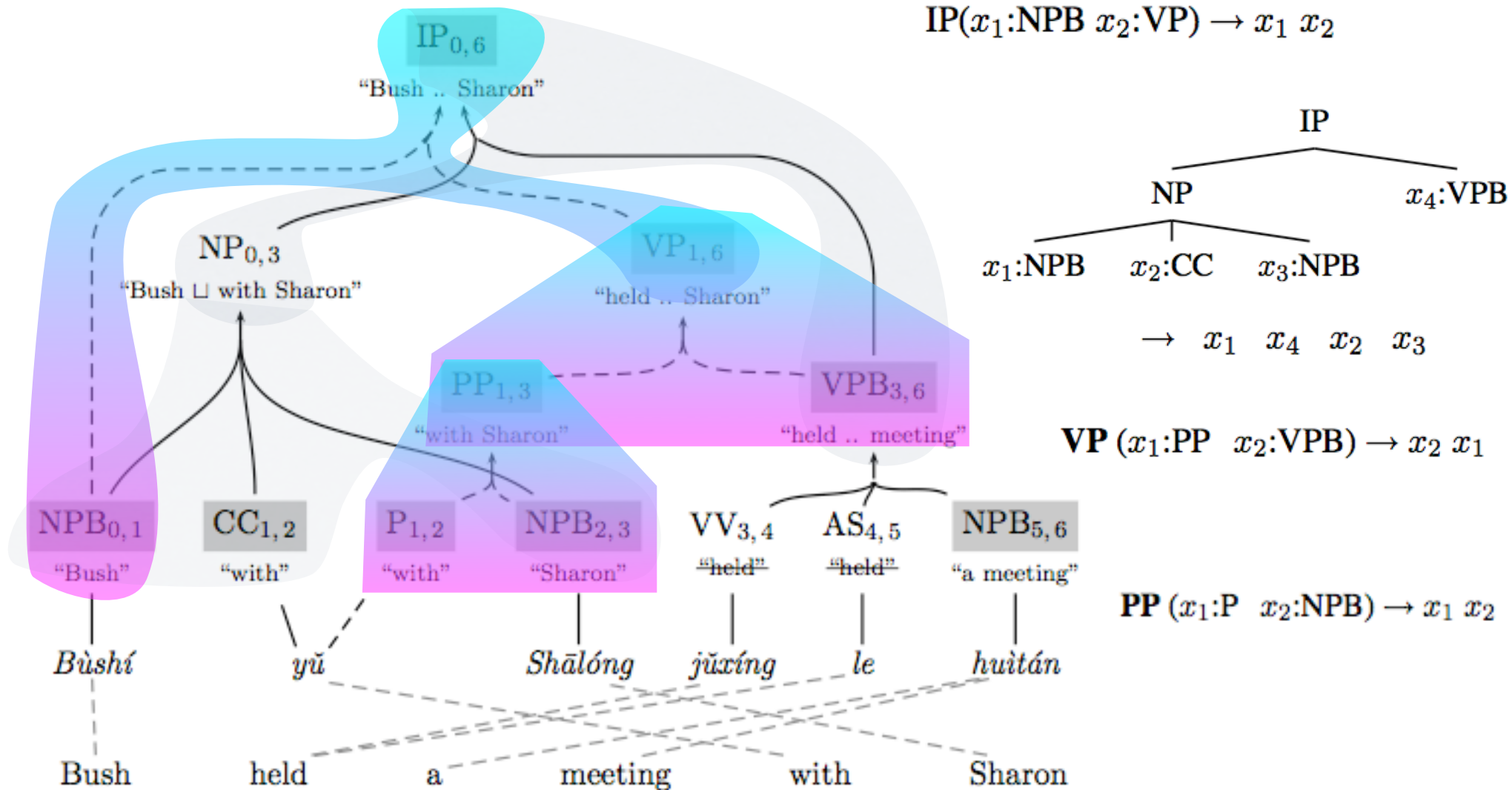
Forest-based Rule Extraction

- forest can extract smaller chunks of rules

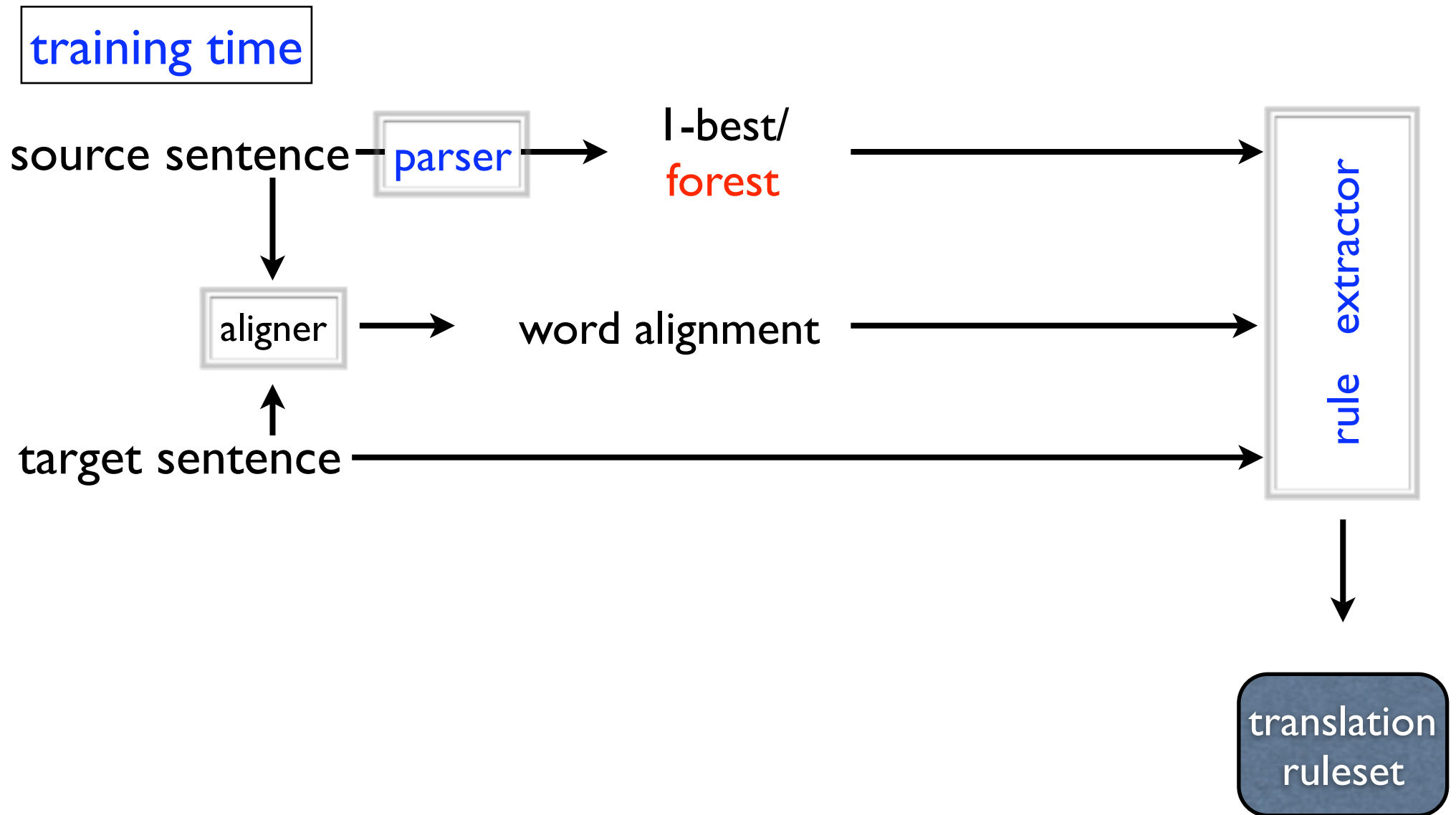


Forest-based Rule Extraction

- forest can extract smaller chunks of rules

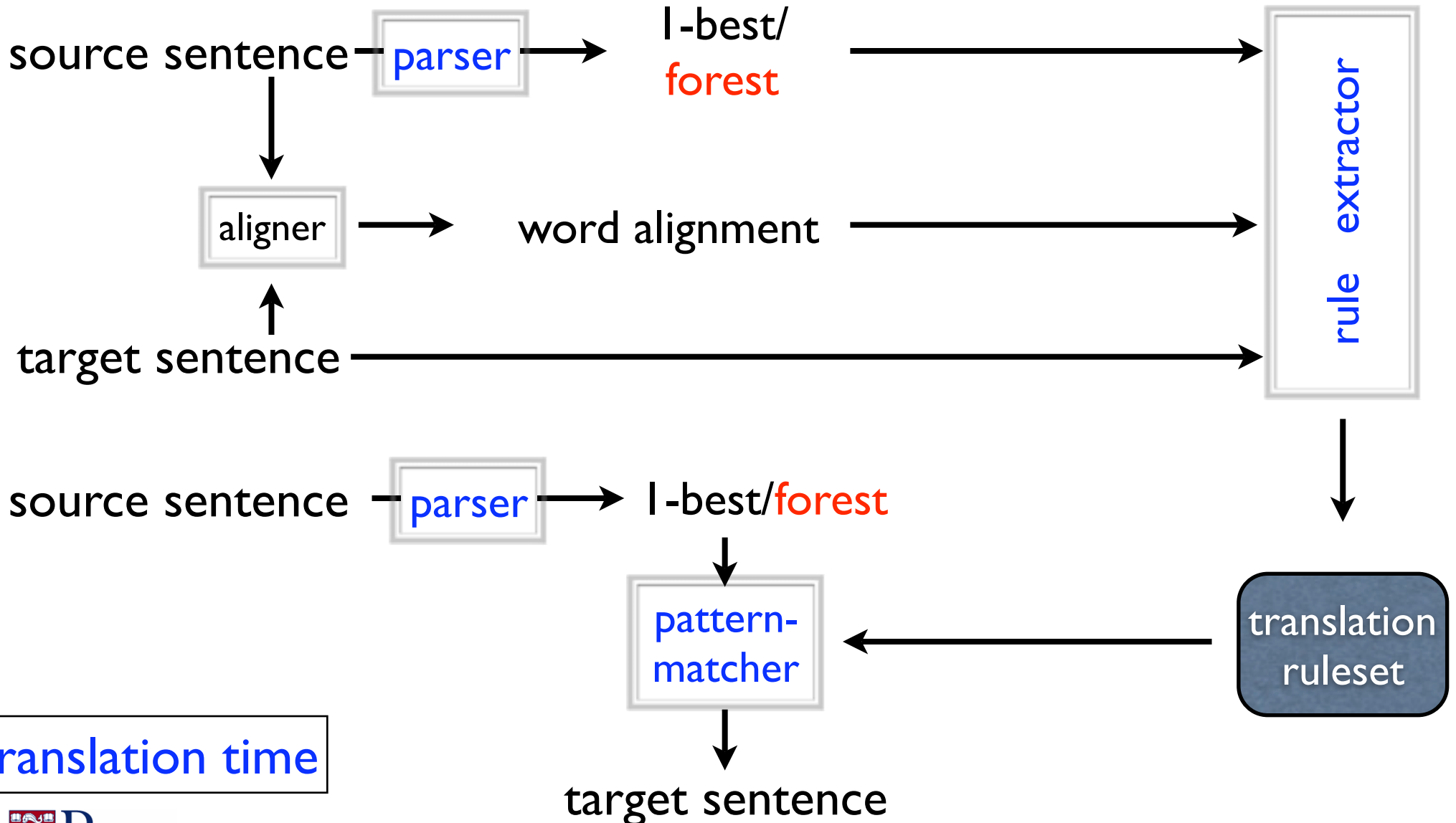


The Forest² Pipeline



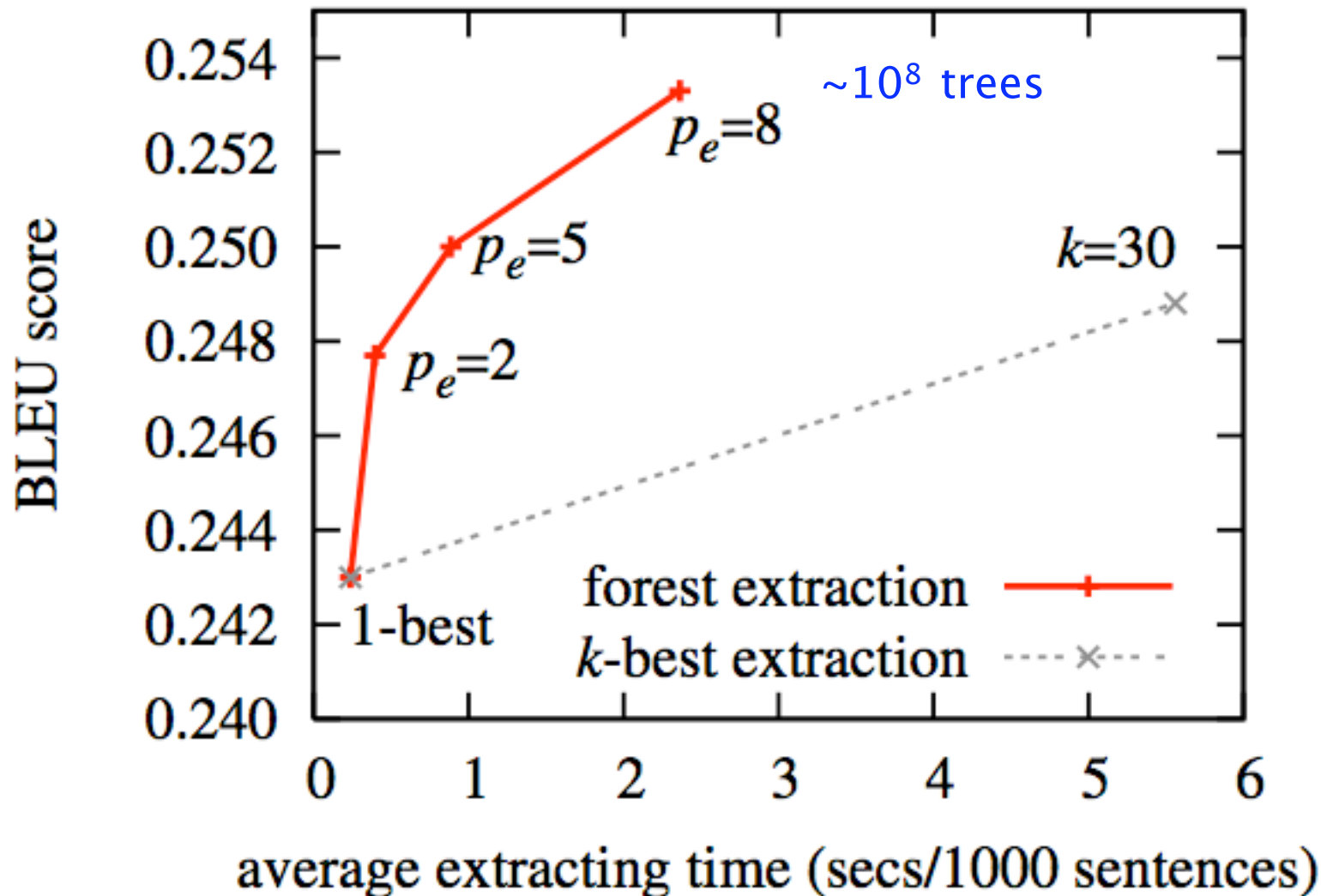
The Forest² Pipeline

training time



Forest vs. k -best Extraction

1.0 Bleu improvement over 1-best,
twice as fast as 30-best extraction



Forest²

- FBIS: 239k sentence pairs (7M/9M Chinese/English words)
- forest in both extraction and decoding
- forest² results is 2.5 points better than I-best²
- and outperforms Hiero (Chiang 2007) by quite a bit

translating on ... →

rules from ... ↓

	I-best tree	forest
I-best tree	0.2560	0.2674
30-best trees	0.2634	0.2767
forest	0.2679	0.2816
Hiero	0.2738	

Translation Examples



- **src** 鲍威尔 说 与 阿拉法特 会谈 很 重要
Bàowēir shuō yǔ Alāfǎtè huìtán hěn zhòngyào
Powell say with Arafat talk very important

- **I-best²** Powell said the very important talks with Arafat

- **forest²** Powell said his meeting with Arafat is very important

- **hiero** Powell said very important talks with Arafat

Conclusions

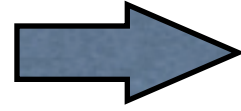
- main theme: efficient syntax-directed translation
- forest-based translation
 - forest = “underspecified syntax”: polynomial vs. exponential
 - still fast (with pruning), yet does not commit to 1-best tree
 - translating millions of trees is faster than just on top- k trees
- forest-based rule extraction: improving rule set quality
- very simple idea, but works well in practice
 - significant improvement over 1-best syntax-directed
 - final result outperforms hiero by quite a bit

Forest is your friend in machine translation.



help save the forest.

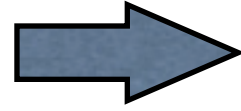
More “forest-based” algorithms in my thesis (this talk is about Chap. 6).



self-service terminals

carefully slide

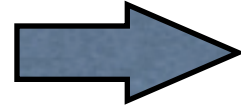
<http://translate.google.com>



self-service terminals



carefully slide



self-service terminals



carefully



<http://translate.google.com>

Larger Decoding Experiments (ACL)

- 2.2M sentence pairs (57M Chinese and 62M English words)
- larger trigram models (1/3 of Xinhua Gigaword)
- also use **bilingual phrases** (BP) as flat translation rules
 - phrases that are consistent with syntactic constituents
- forest enables larger improvement with BP

	T2S	T2S+BP
1-best tree	0.2666	0.2939
30-best trees	0.2755	0.3084
forest	0.2839	0.3149
improvement	1.7	2.1