

Rule Markov Models for Fast Tree-to-String Translation

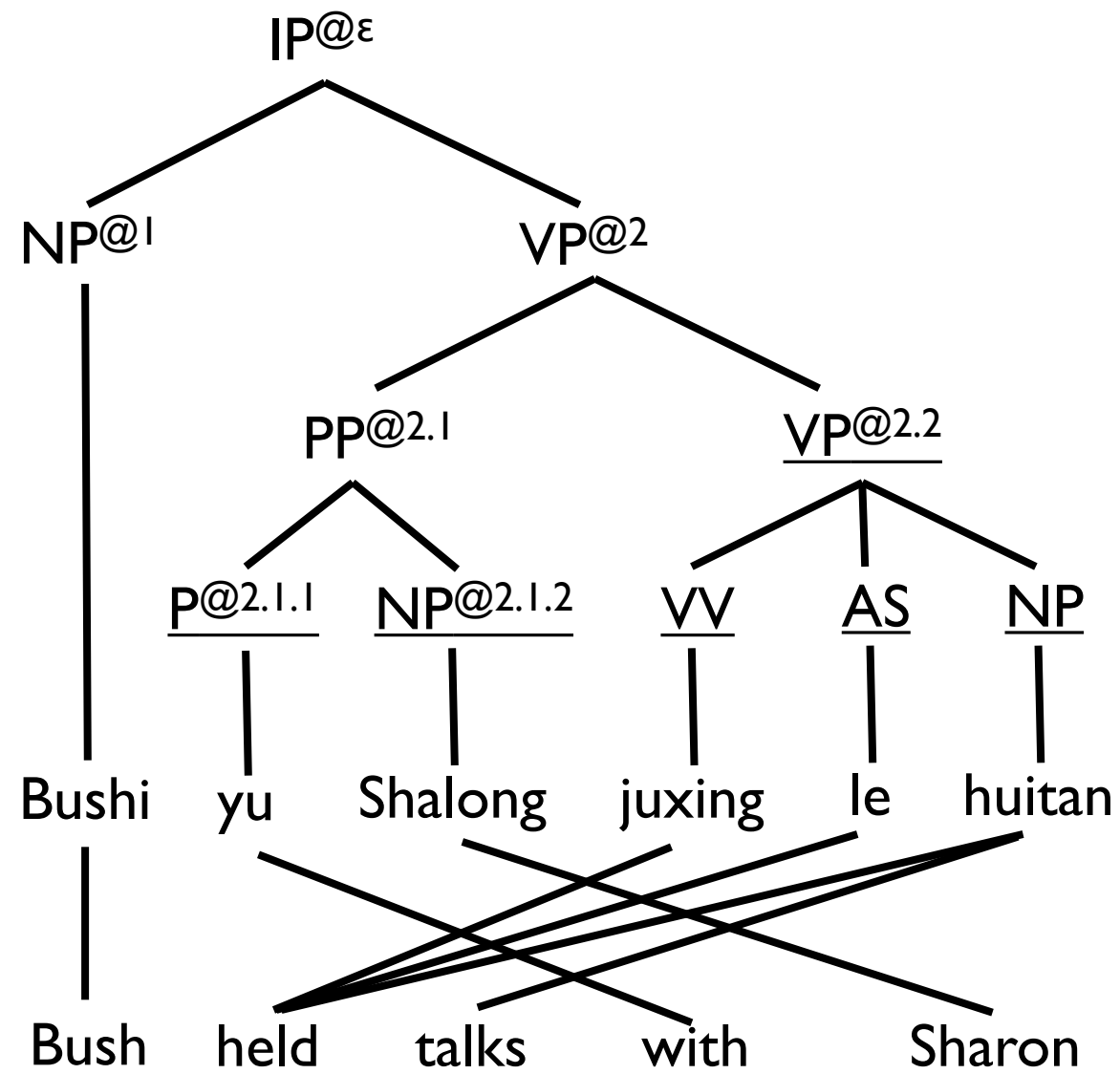
Ashish Vaswani
USC / ISI

Haitao Mi
CAS / ICT

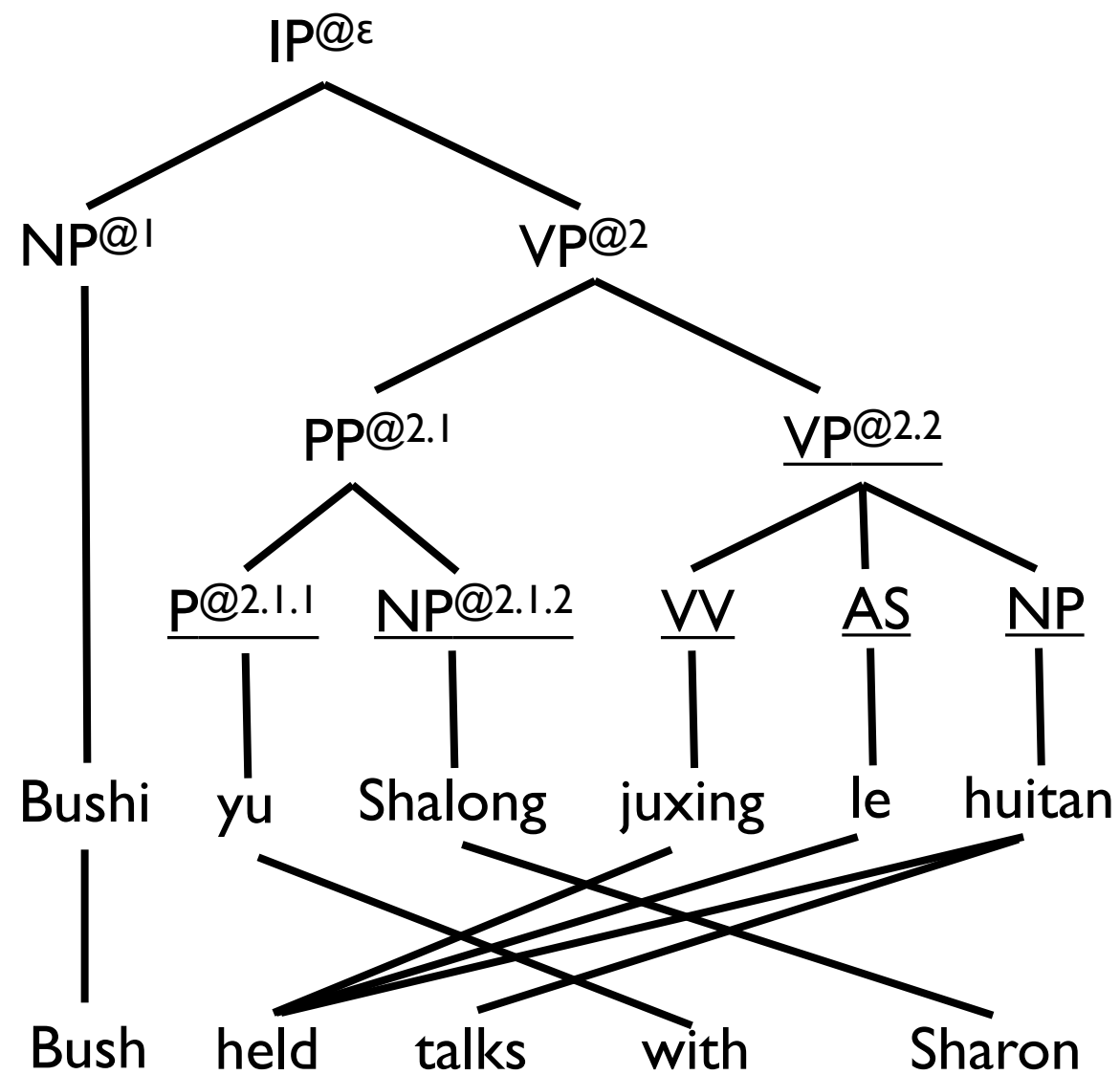
Liang Huang
USC / ISI

David Chiang
USC / ISI

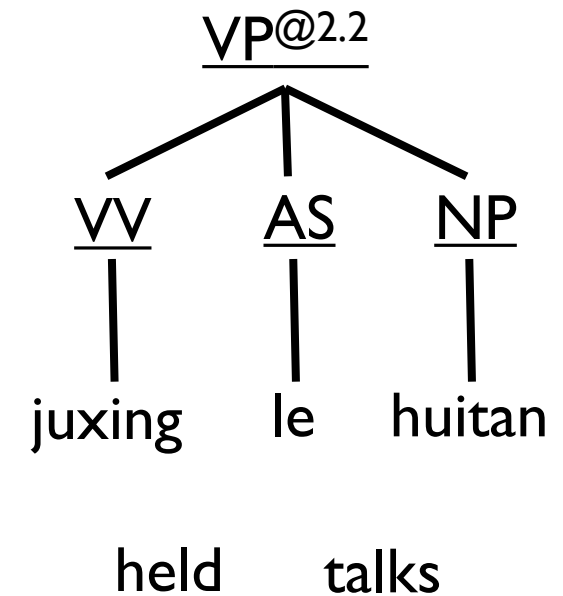
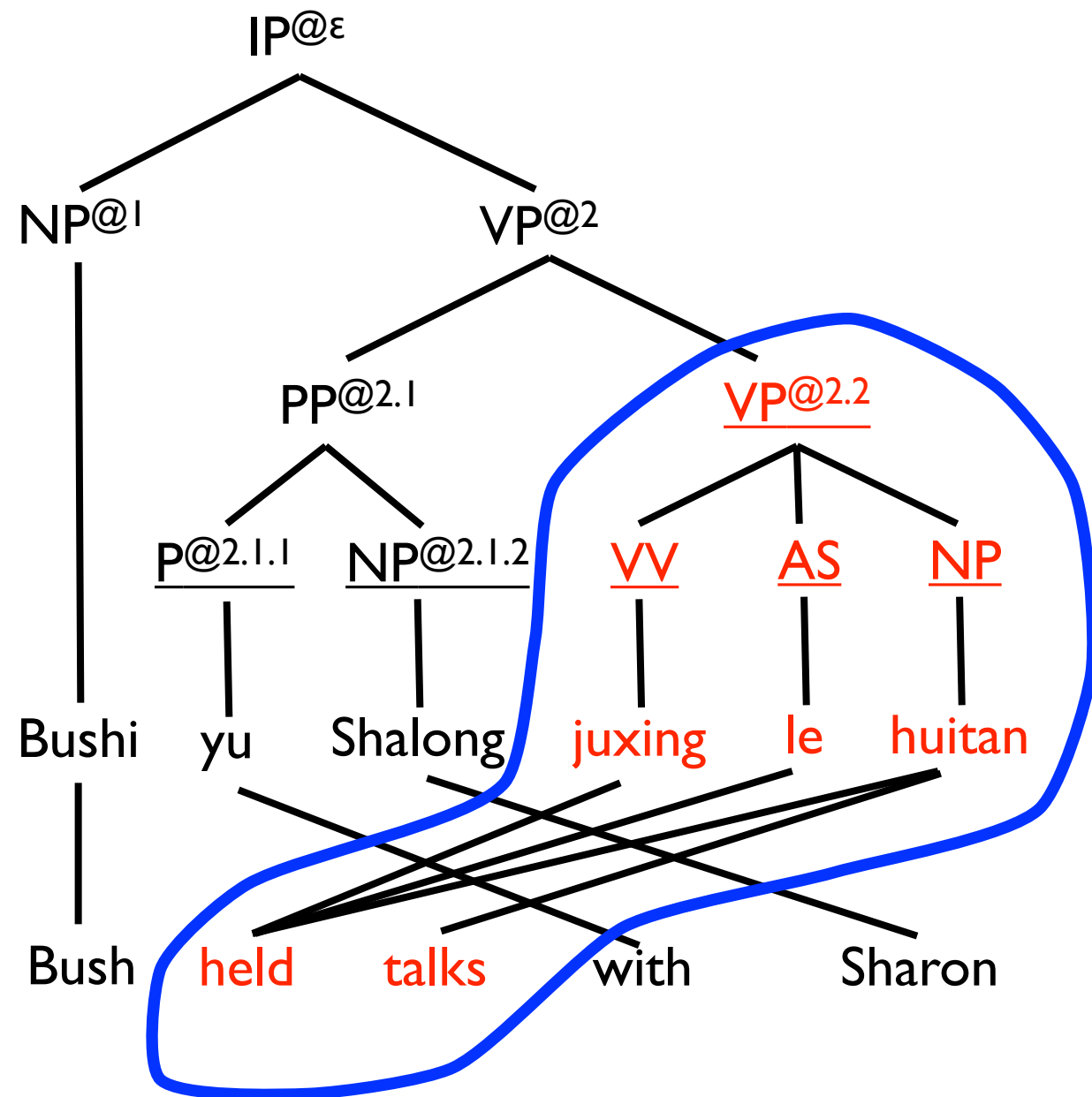
Syntax Based Machine Translation (SBMT)



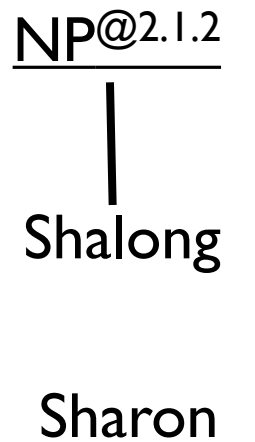
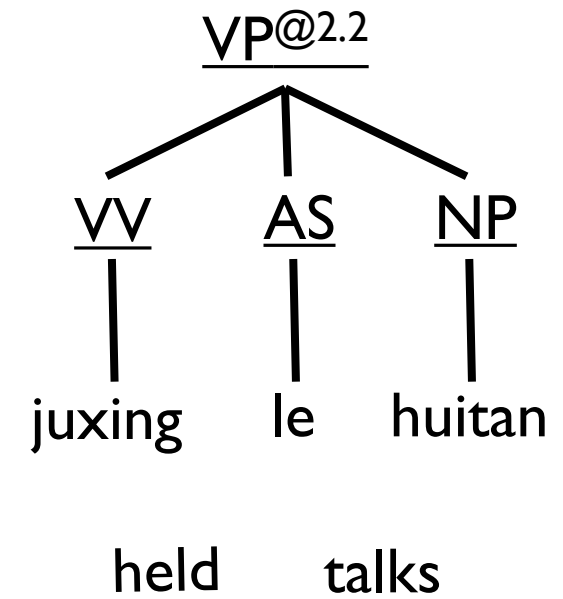
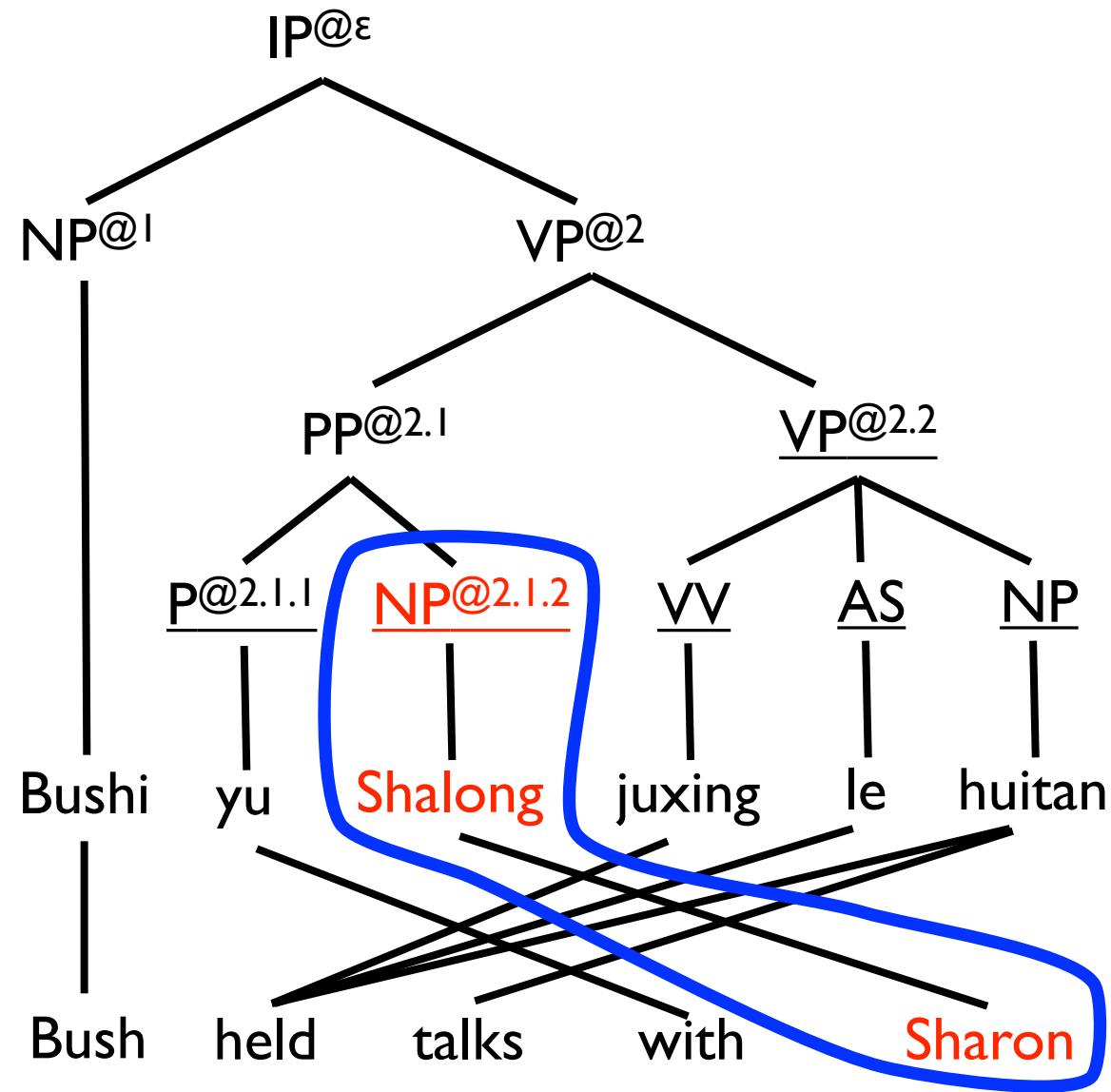
Syntax Based Machine Translation (SBMT)



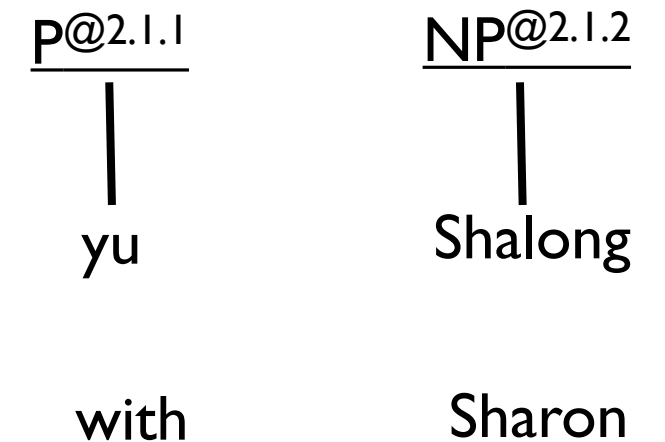
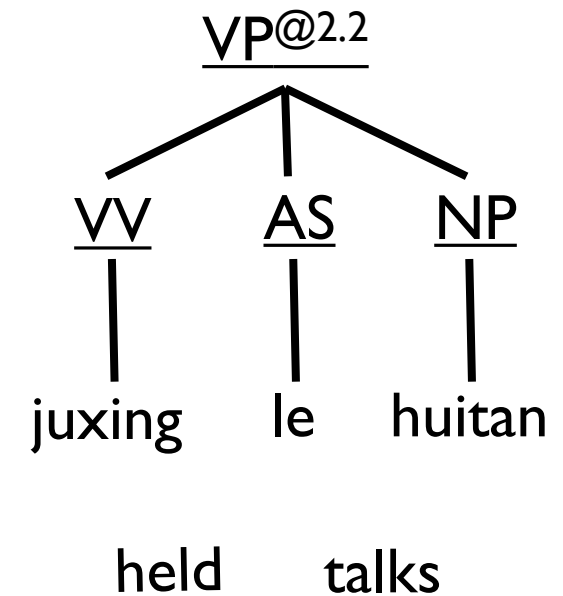
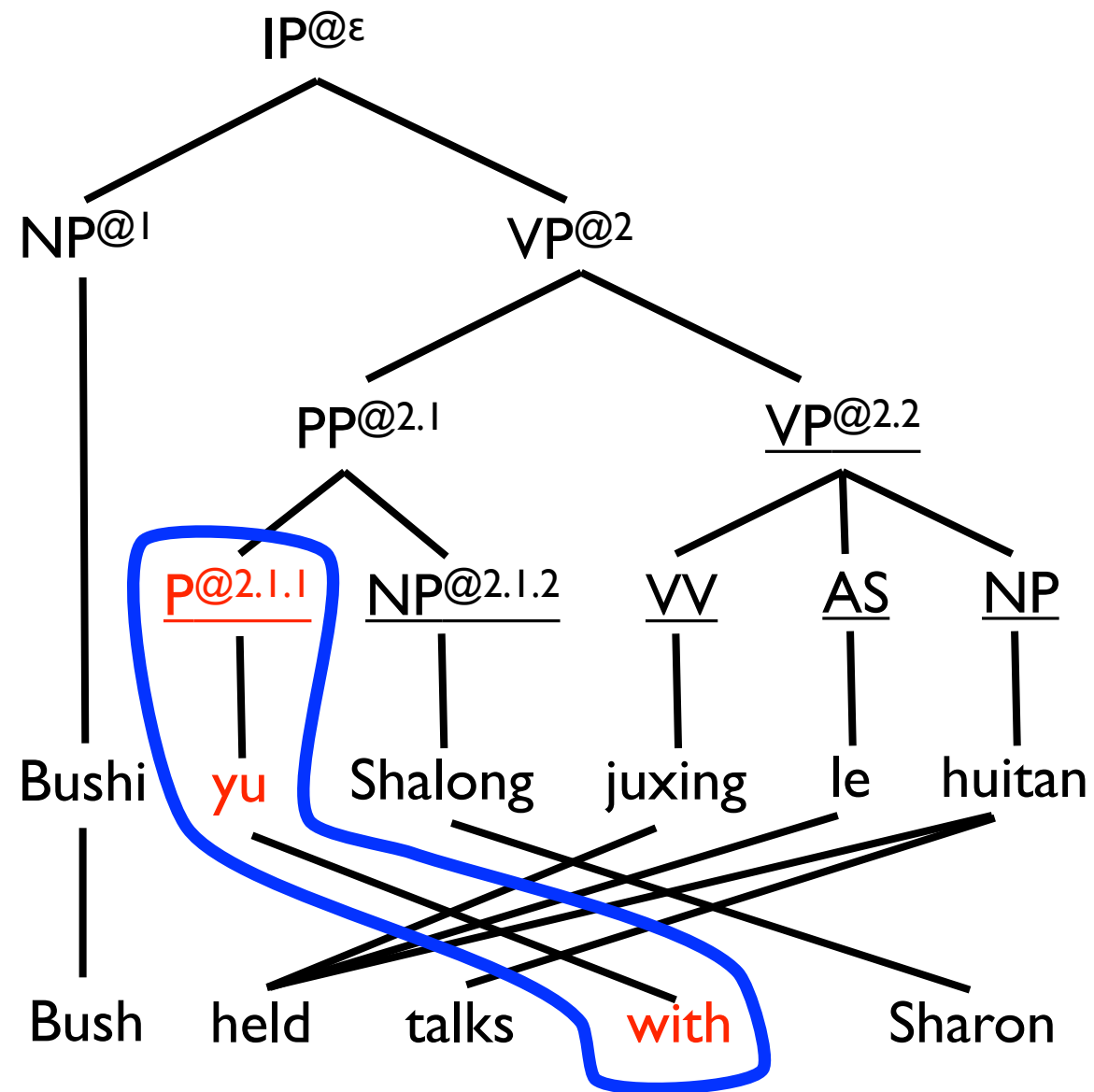
Syntax Based Machine Translation



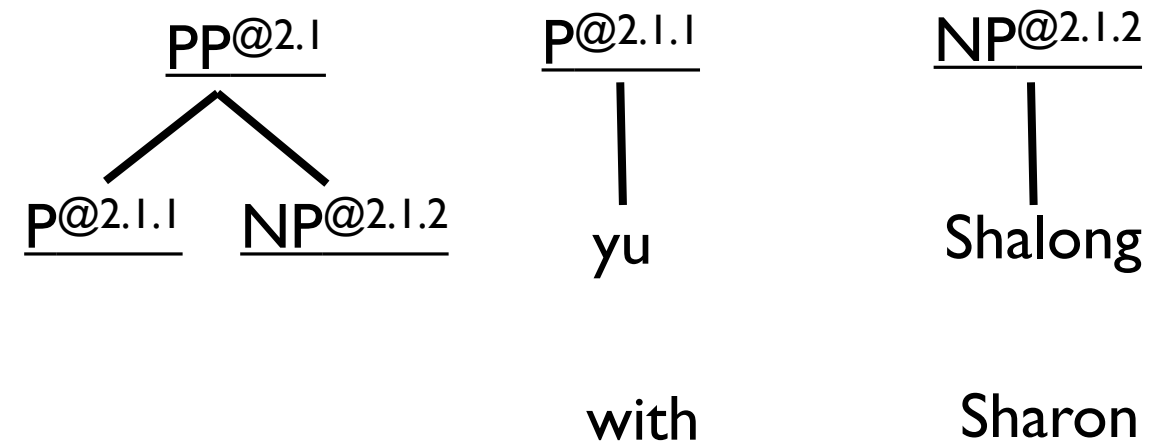
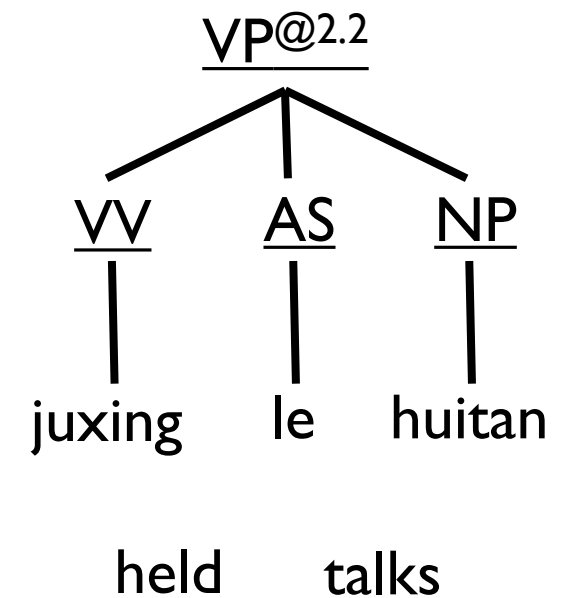
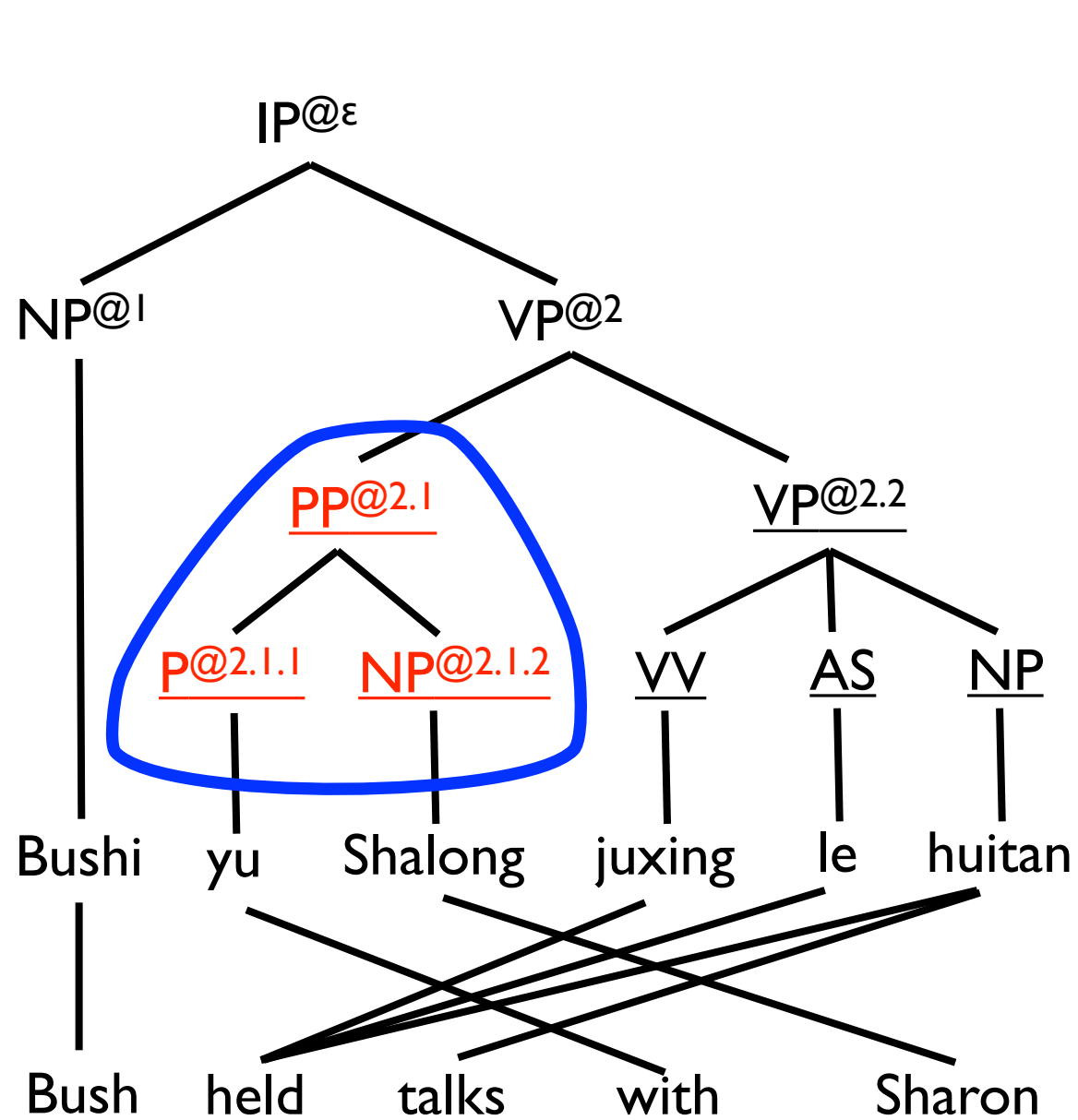
Syntax Based Machine Translation



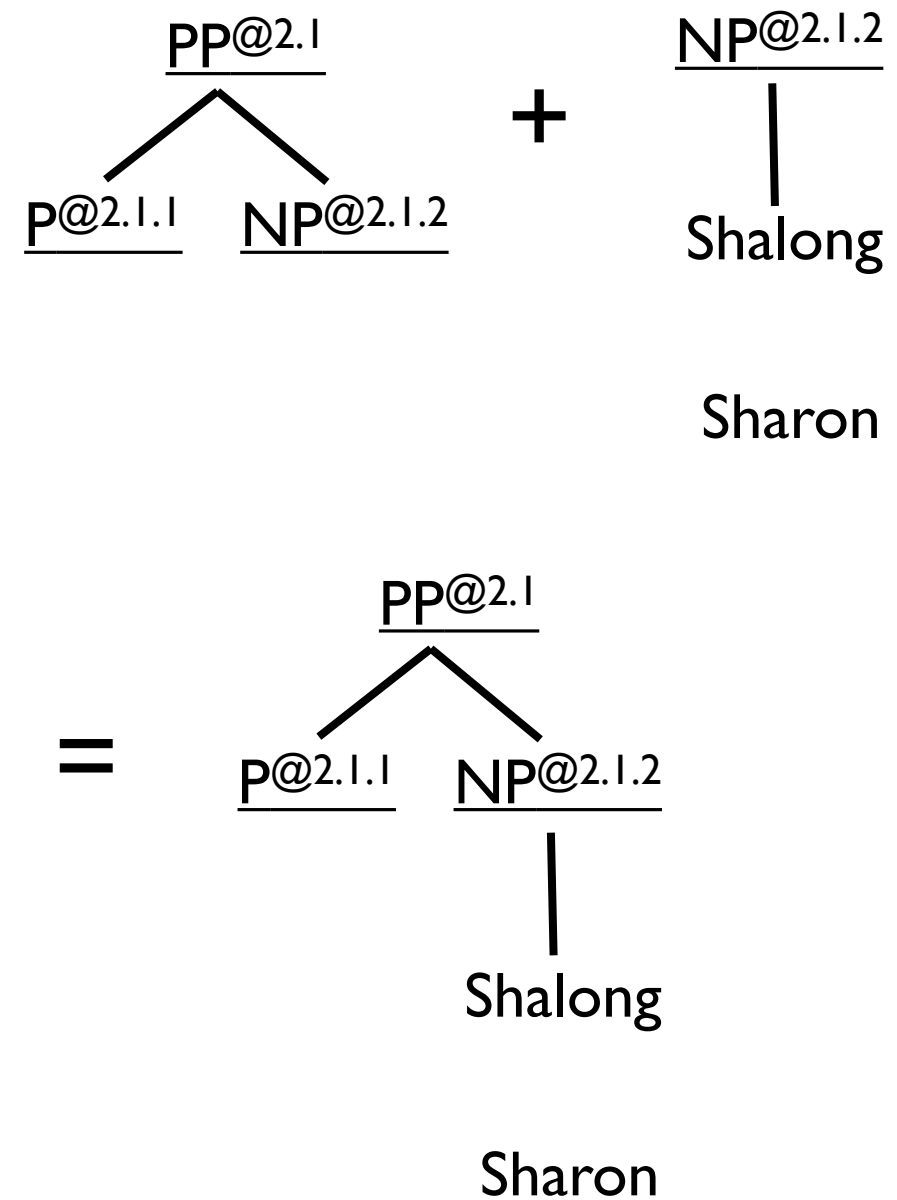
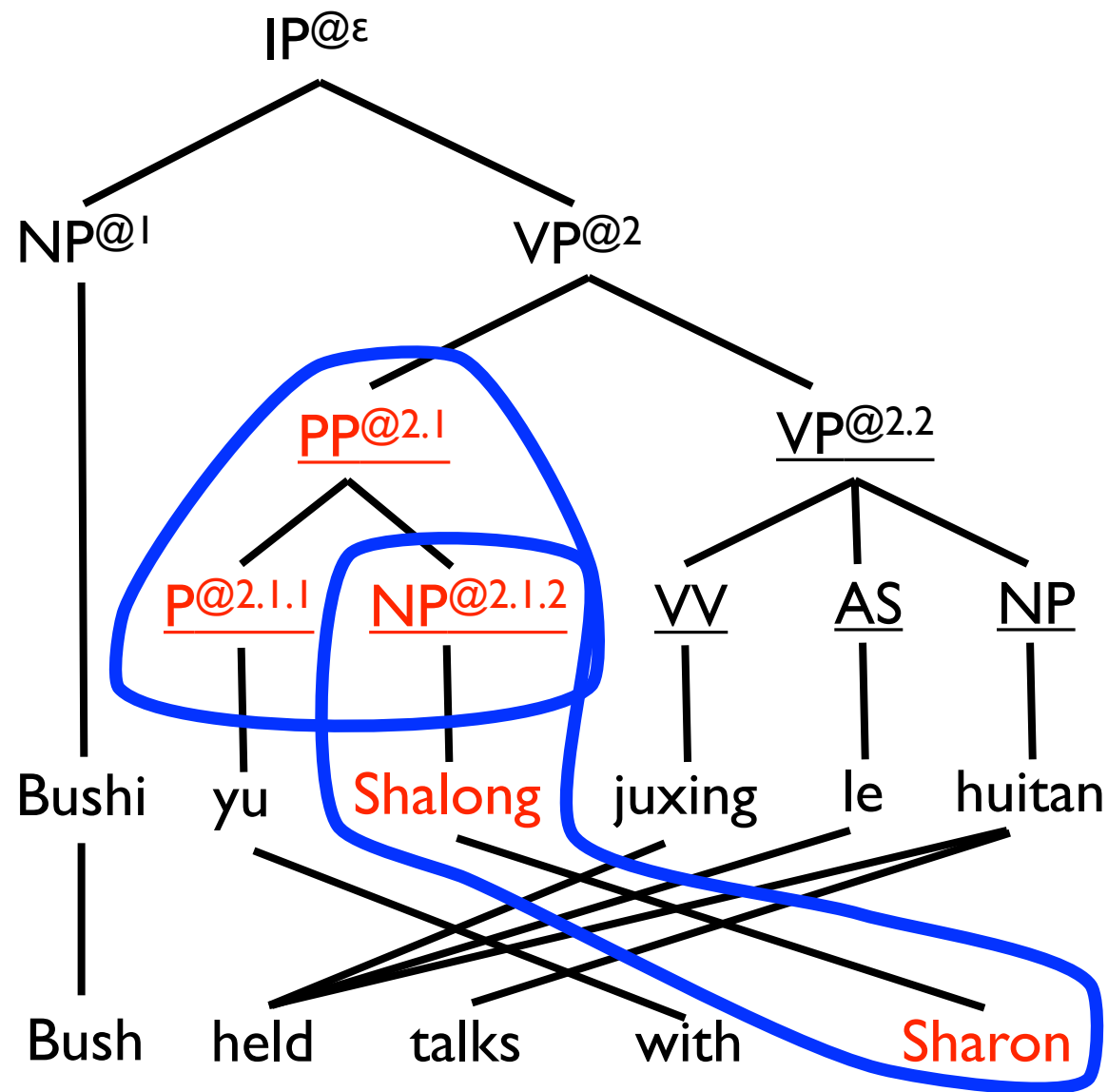
Syntax Based Machine Translation



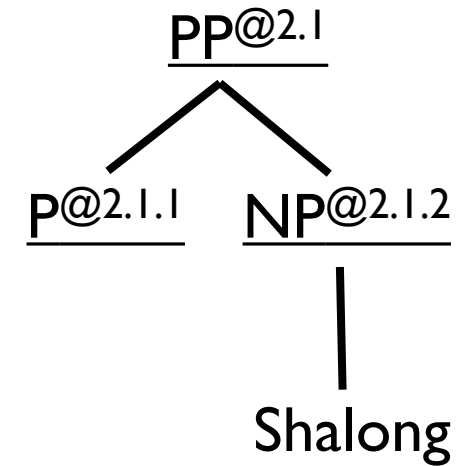
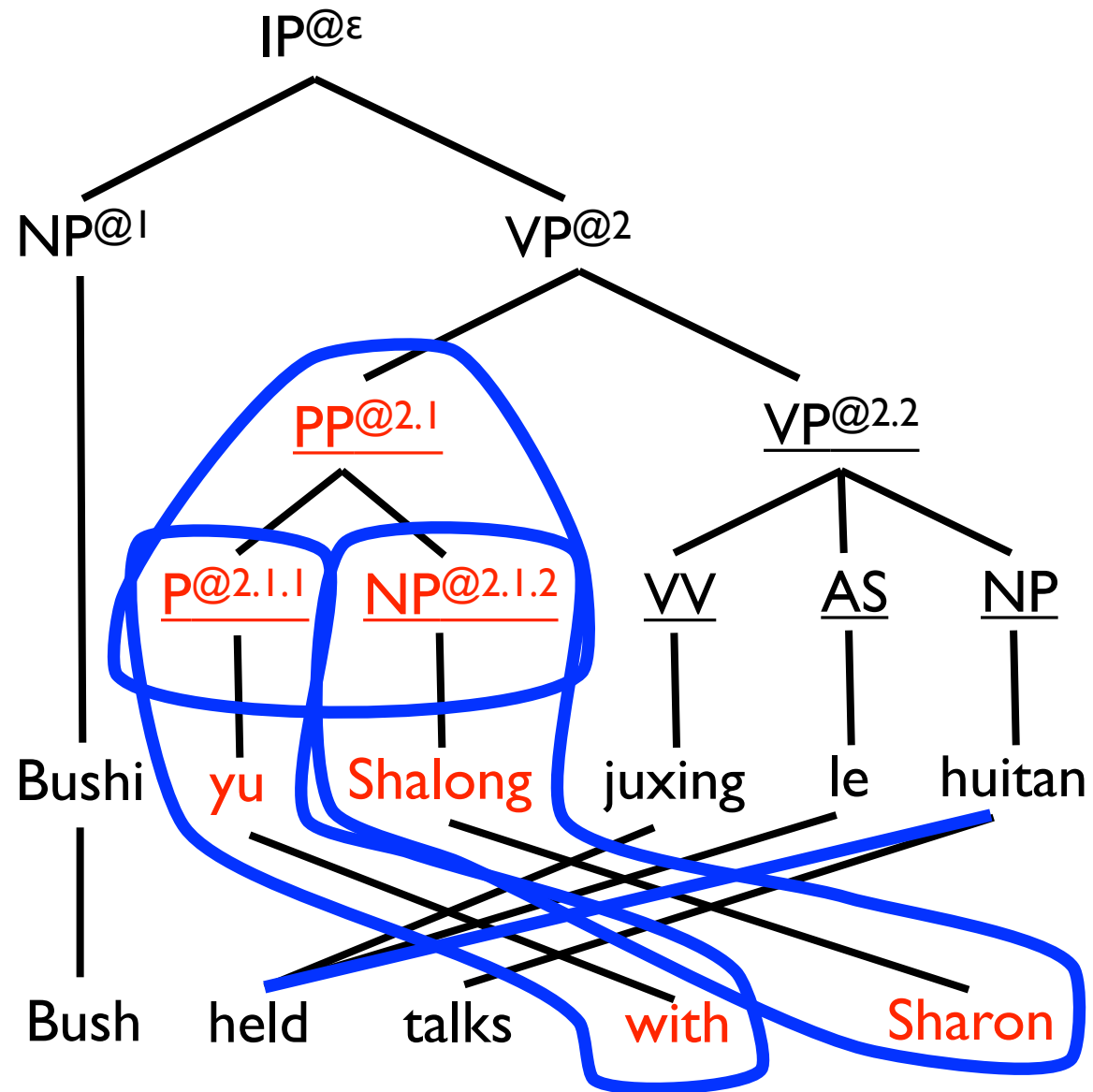
Syntax Based Machine Translation



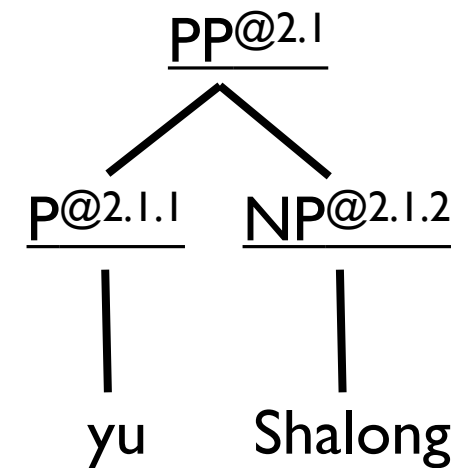
Syntax Based Machine Translation



Syntax Based Machine Translation

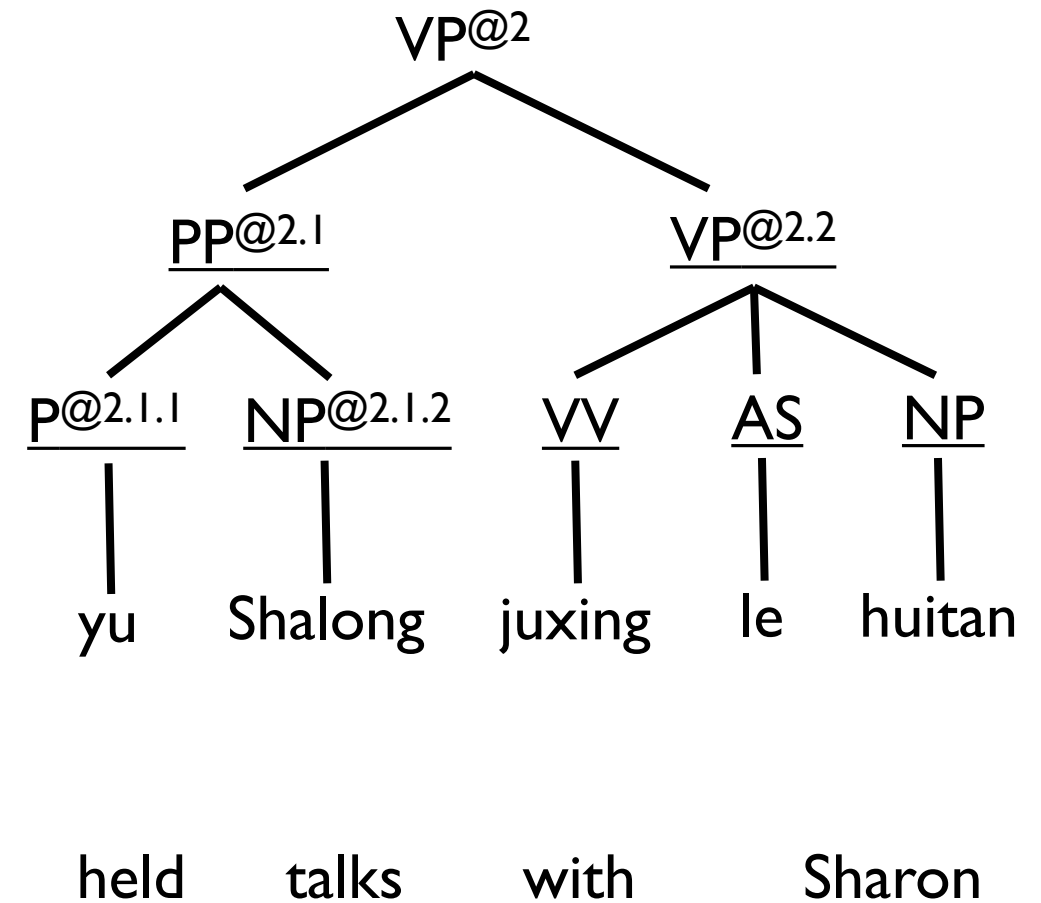
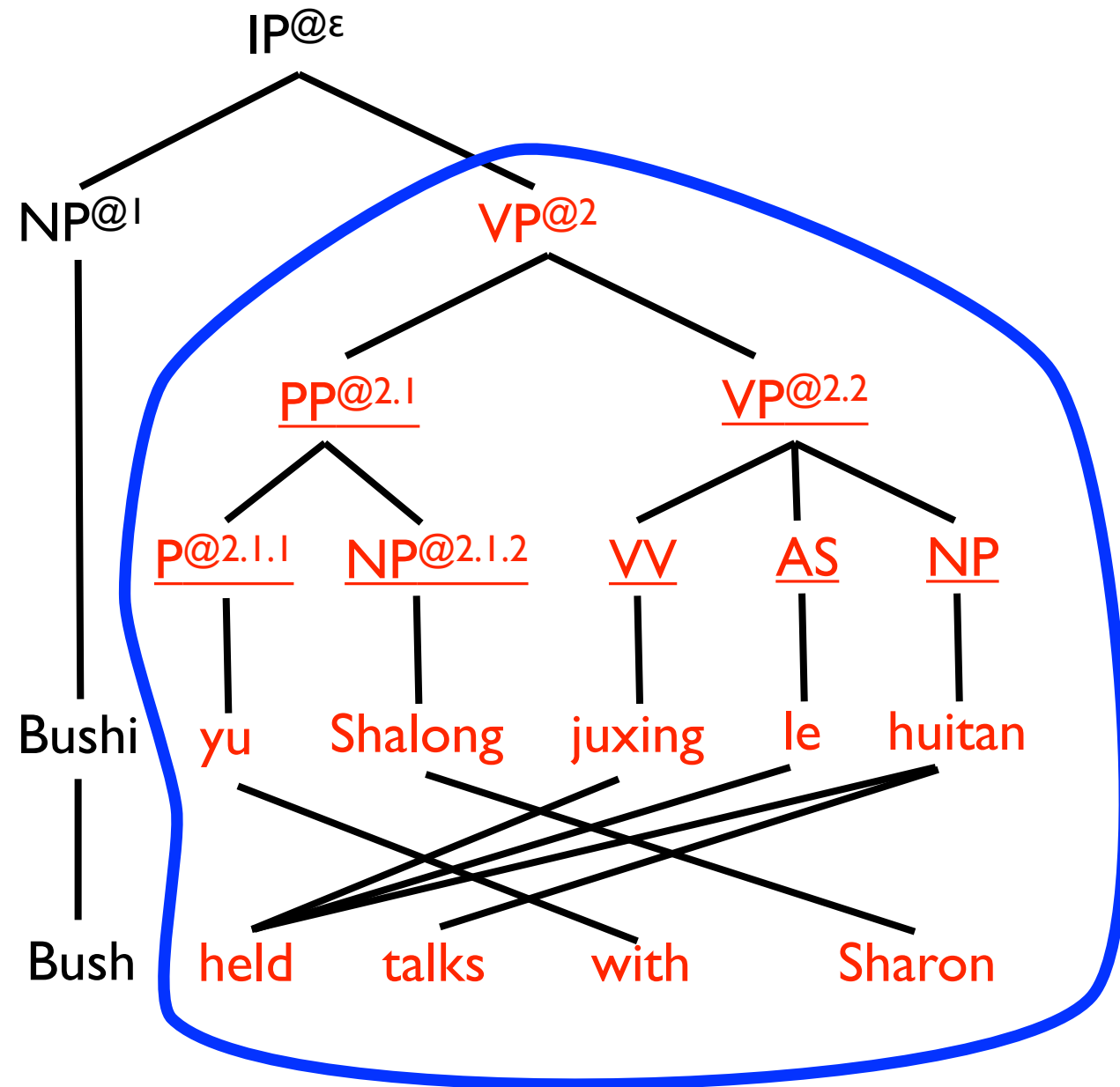


Sharon









with Sharon

Syntax Based Machine Translation



Goal

	BLEU	Decoding time	Grammar size
composed			
minimal			
??			

Our Approach: Rule Markov Models

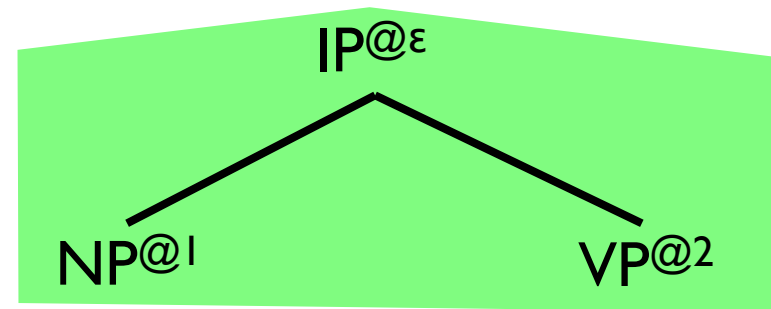
- composed rules alleviate independence b/w rules
- but we can **explicitly** model this dependence by conditioning on the history (grammar annotation)
- Learn the ancestral dependencies between minimal rules with rule Markov models (RMM)
- Use rule Markov models with top down incremental decoder (Huang and Mi, 2010)

Generative Story

IP@ ϵ

Generative Story

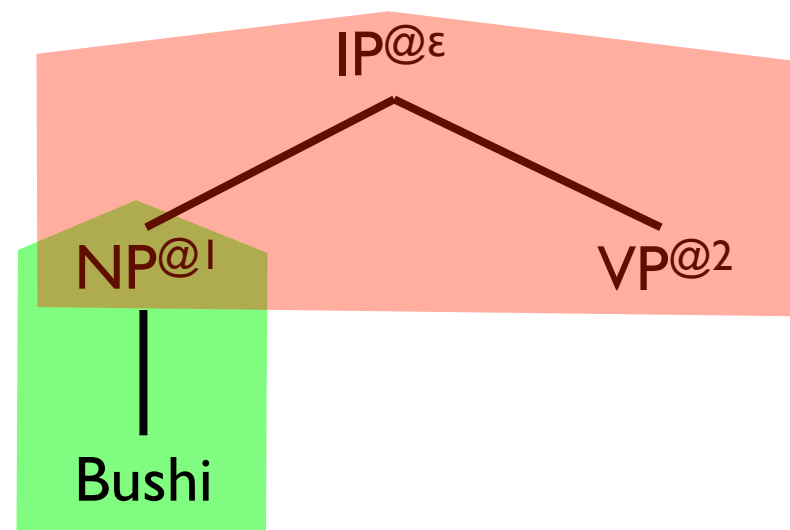
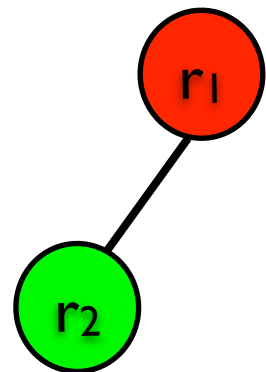
r_1



$P(r_1|\epsilon)$

$NP@1 \quad VP@2$

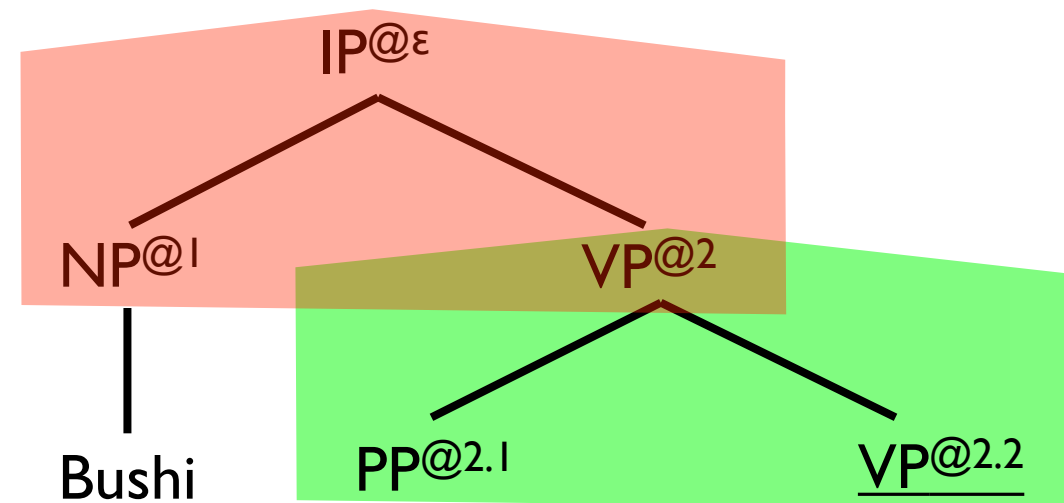
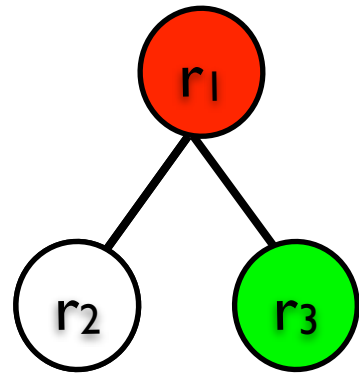
Generative Story



$P(r_2|r_1)$

Bush $VP@2$

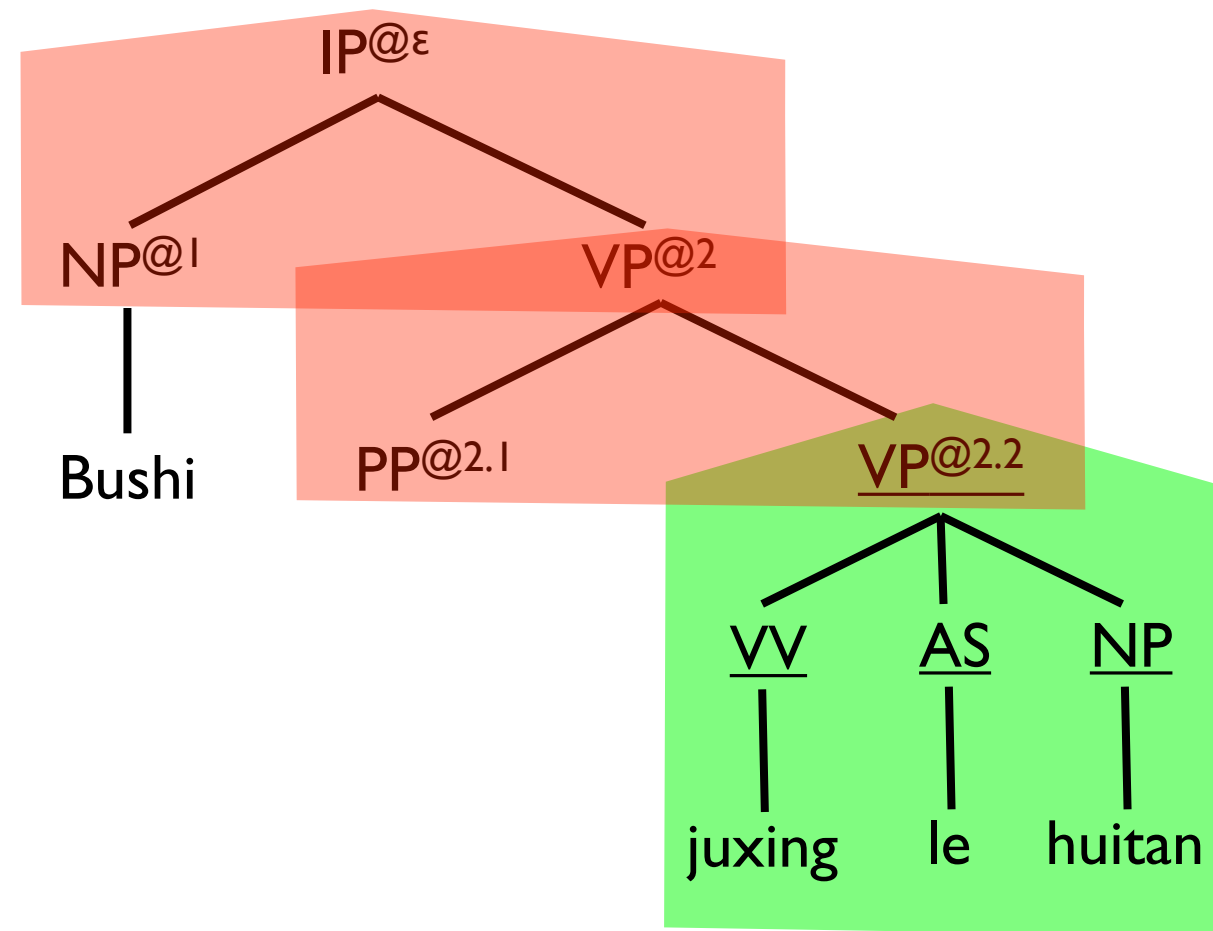
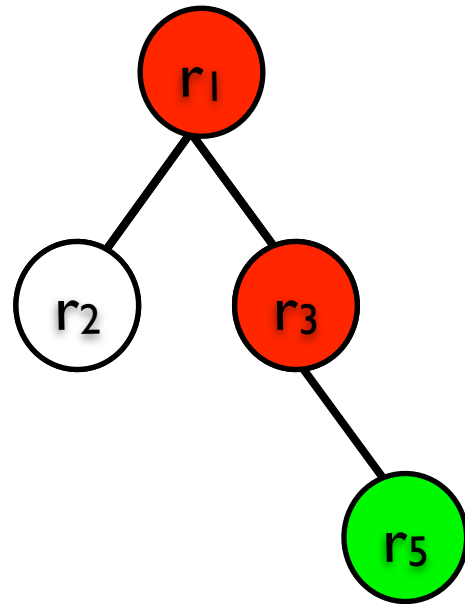
Generative Story



$P(r_3|r_1)$

Bush VP@2.2 PP@2.1

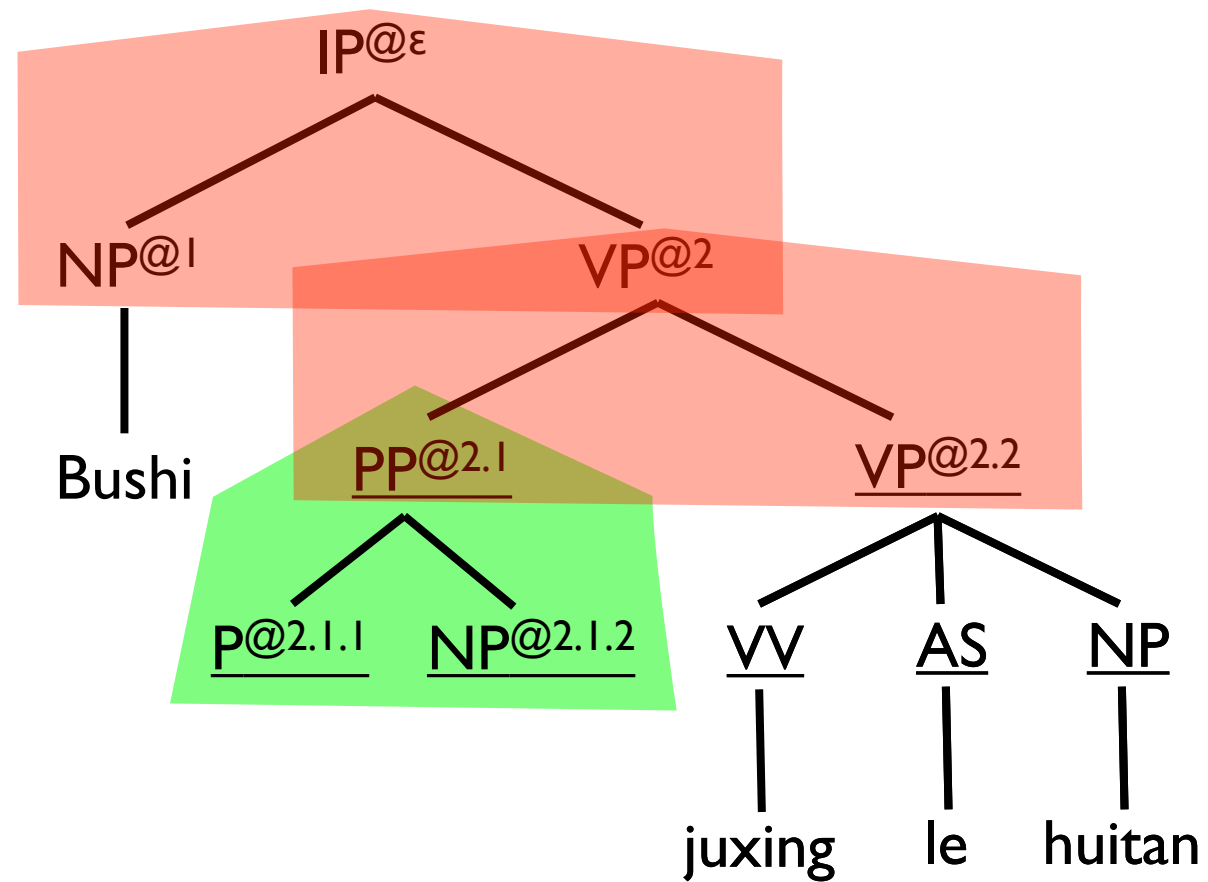
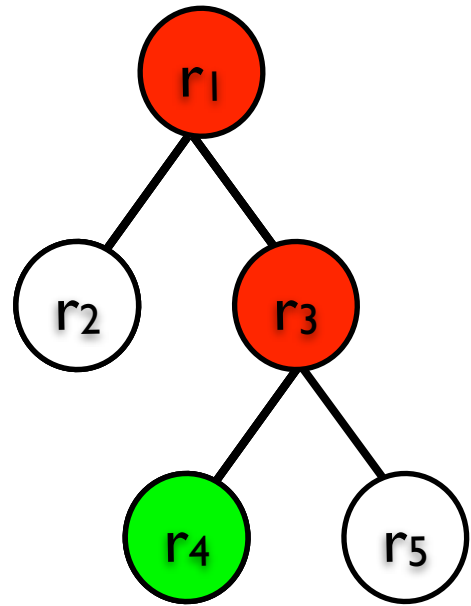
Generative Story



$P(r_5|r_1, r_3)$

Bush **held talks** PP@2.1

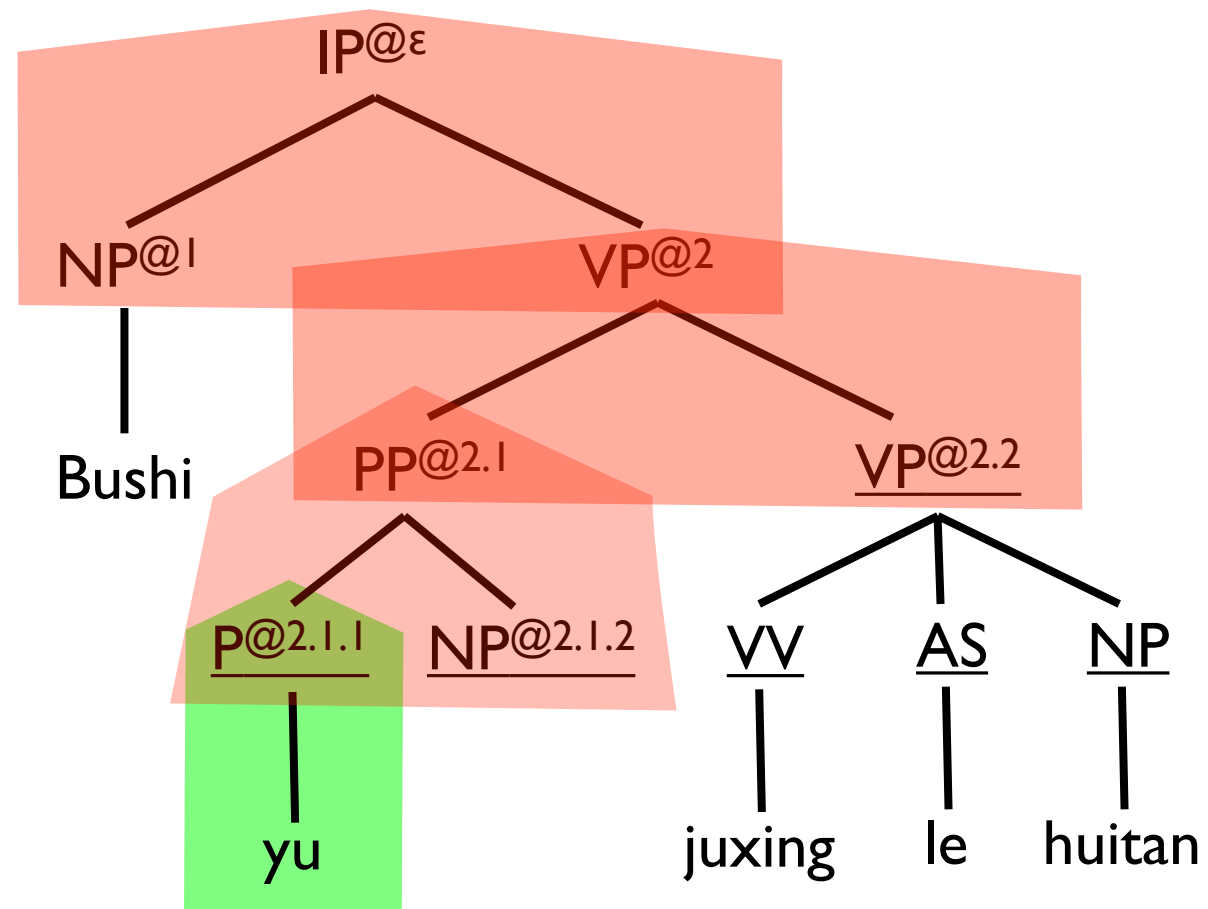
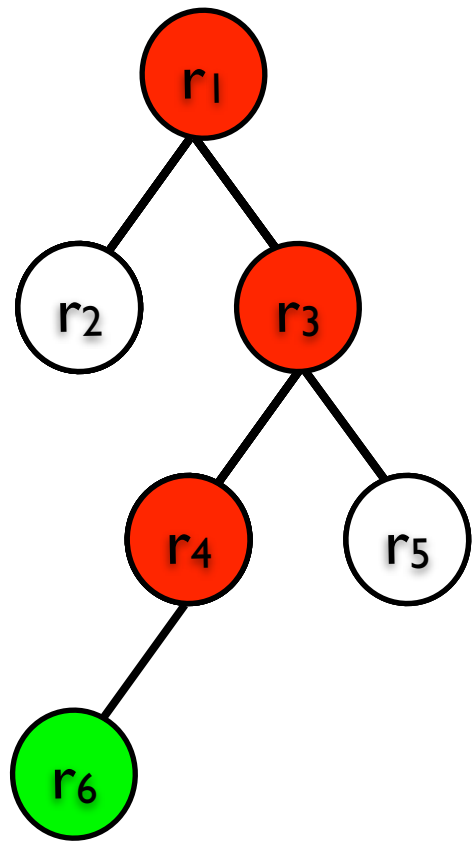
Generative Story



$P(r_4|r_1, r_3)$

Bush held talks P@2.1.1 NP@2.1.2

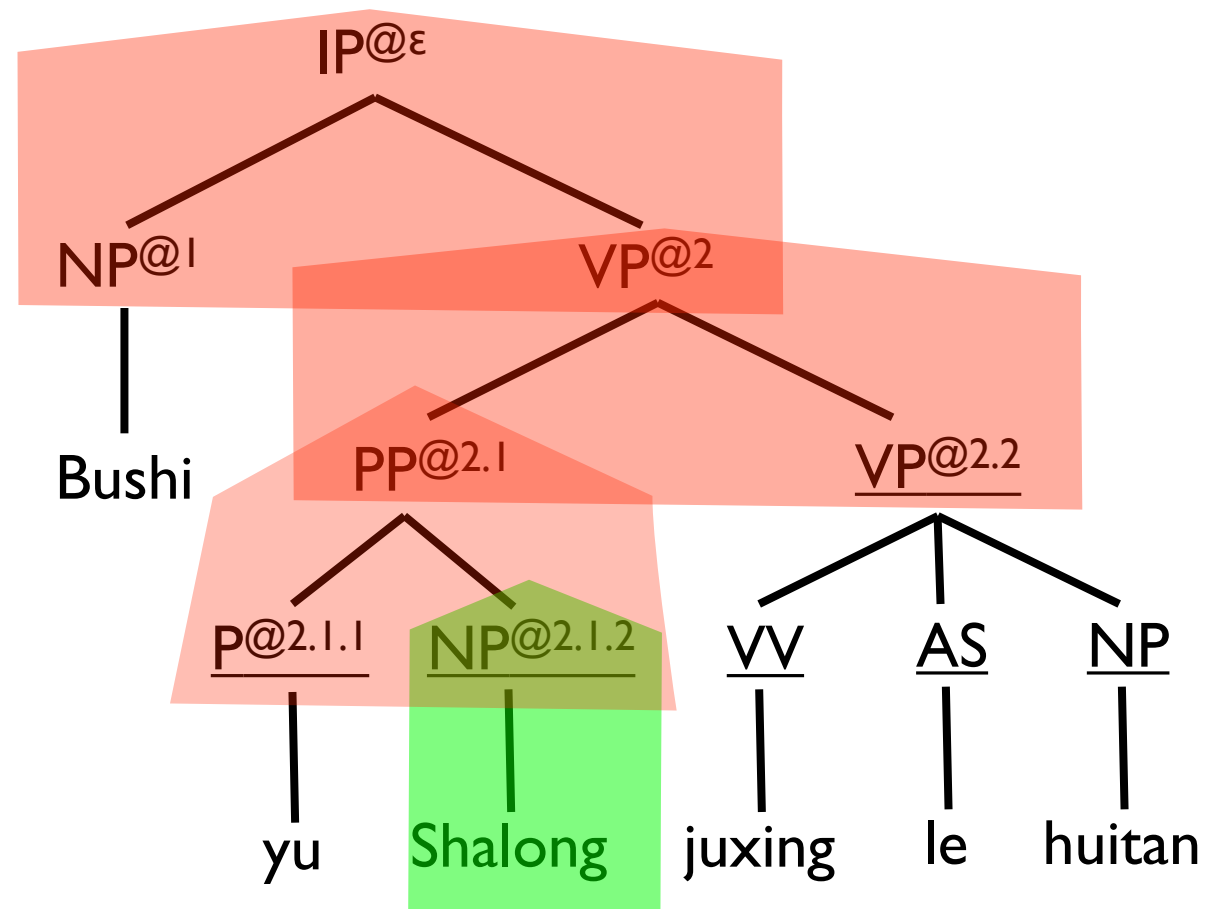
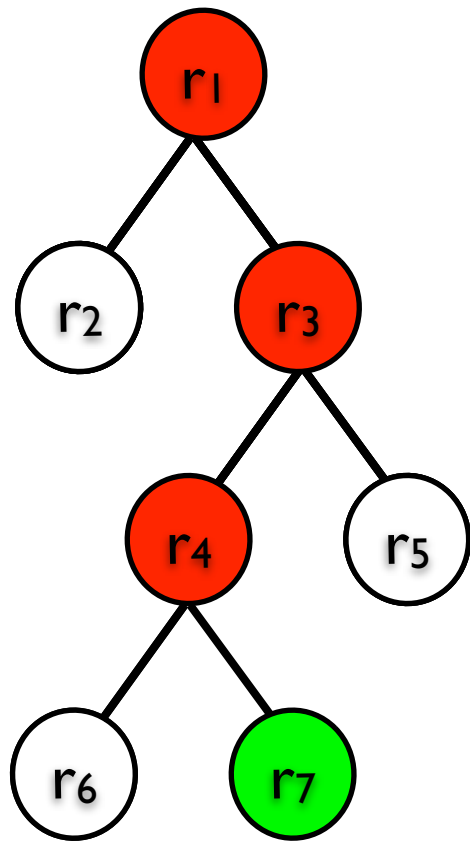
Generative Story



$P(r_6|r_1, r_3, r_4)$

Bush held talks **with** NP@2.1.2

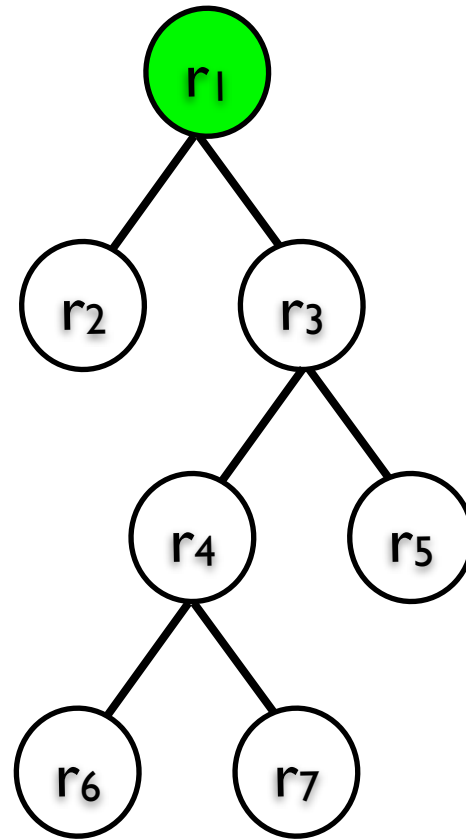
Generative Story



$P(r_7|r_1, r_3, r_4)$

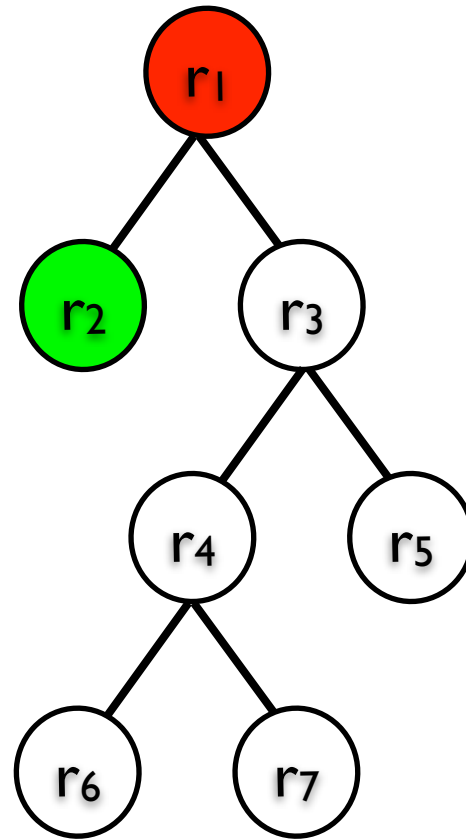
Bush held talks with Sharon

Probability of a Derivation Tree



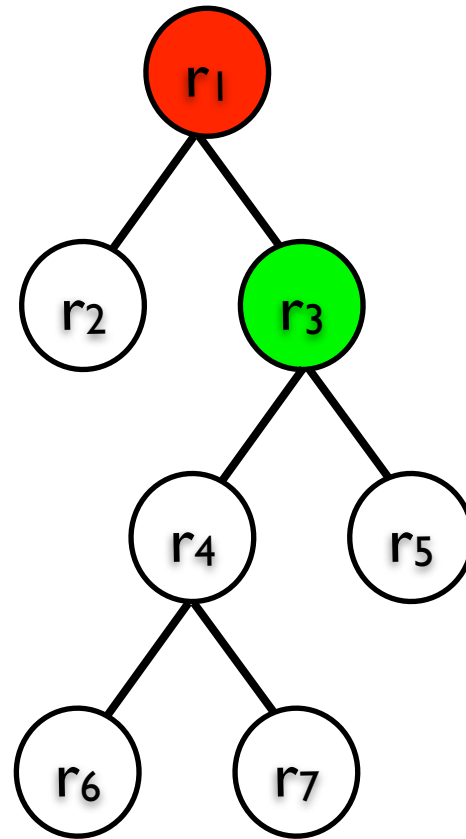
$$\begin{aligned} P(T) = & P(r_1 | \varepsilon) \cdot P(r_2 | r_1) \cdot P(r_3 | r_1) \\ & P(r_5 | r_1, r_3) \cdot P(r_4 | r_1, r_3) \\ & P(r_6 | r_1, r_3, r_4) \cdot P(r_7 | r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree



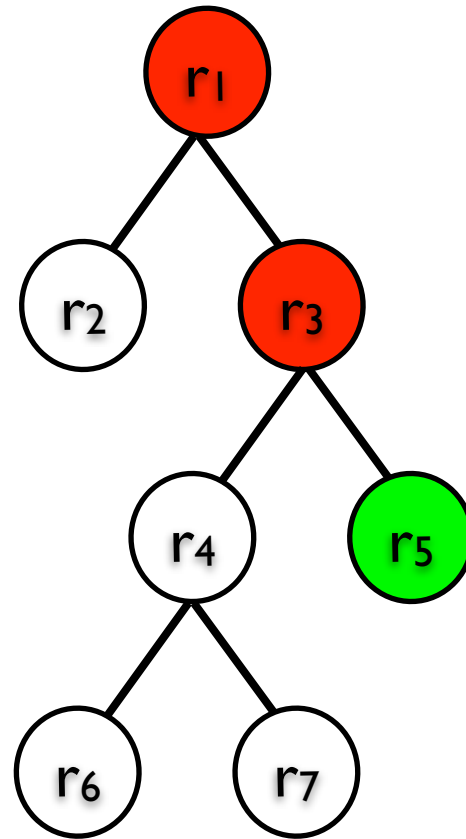
$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1, r_3) \cdot P(r_4|r_1, r_3) \\ & P(r_6|r_1, r_3, r_4) \cdot P(r_7|r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree



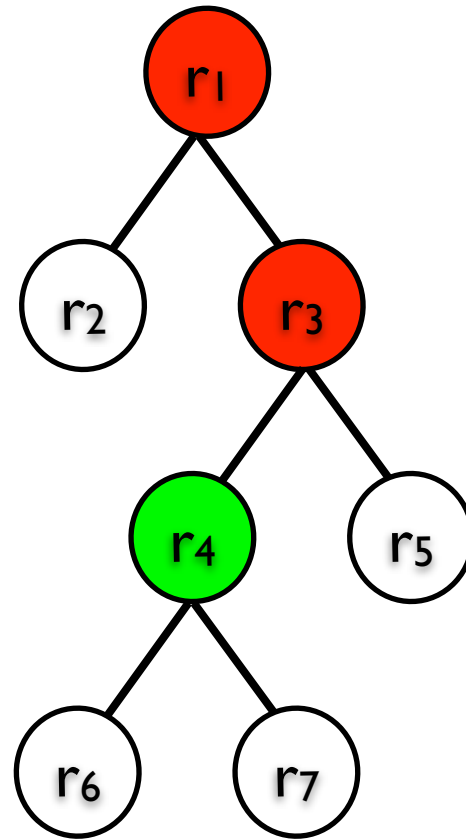
$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1, r_3) \cdot P(r_4|r_1, r_3) \\ & P(r_6|r_1, r_3, r_4) \cdot P(r_7|r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree



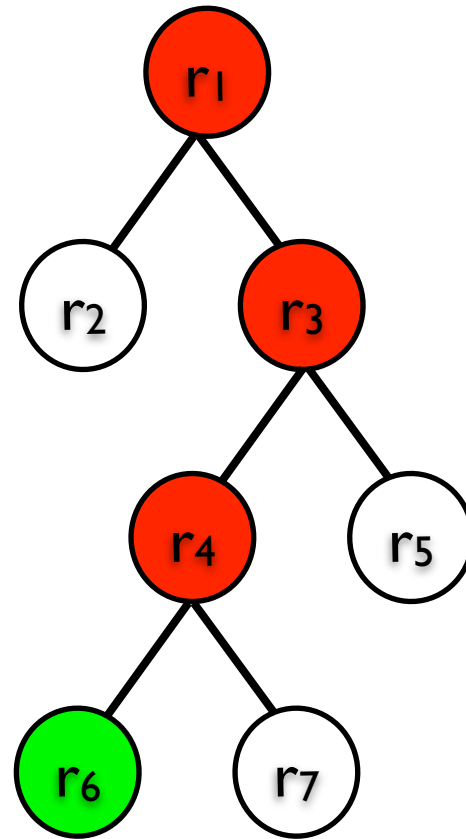
$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1, r_3) \cdot P(r_4|r_1, r_3) \\ & P(r_6|r_1, r_3, r_4) \cdot P(r_7|r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree



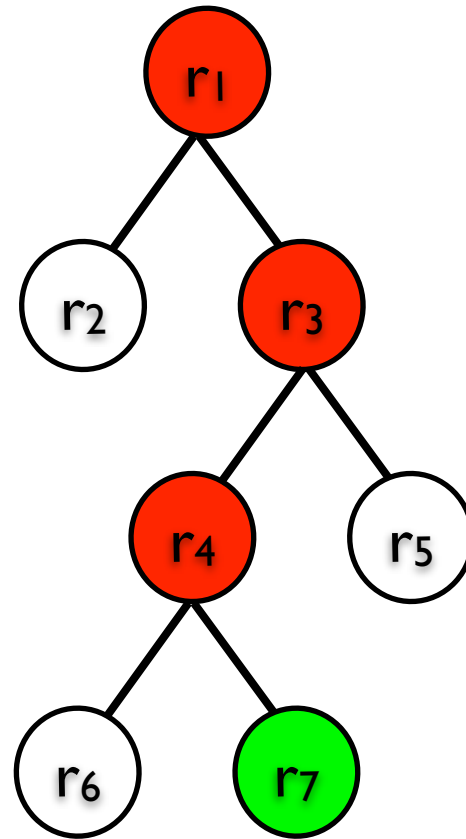
$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1, r_3) \cdot P(r_4|r_1, r_3) \\ & P(r_6|r_1, r_3, r_4) \cdot P(r_7|r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree



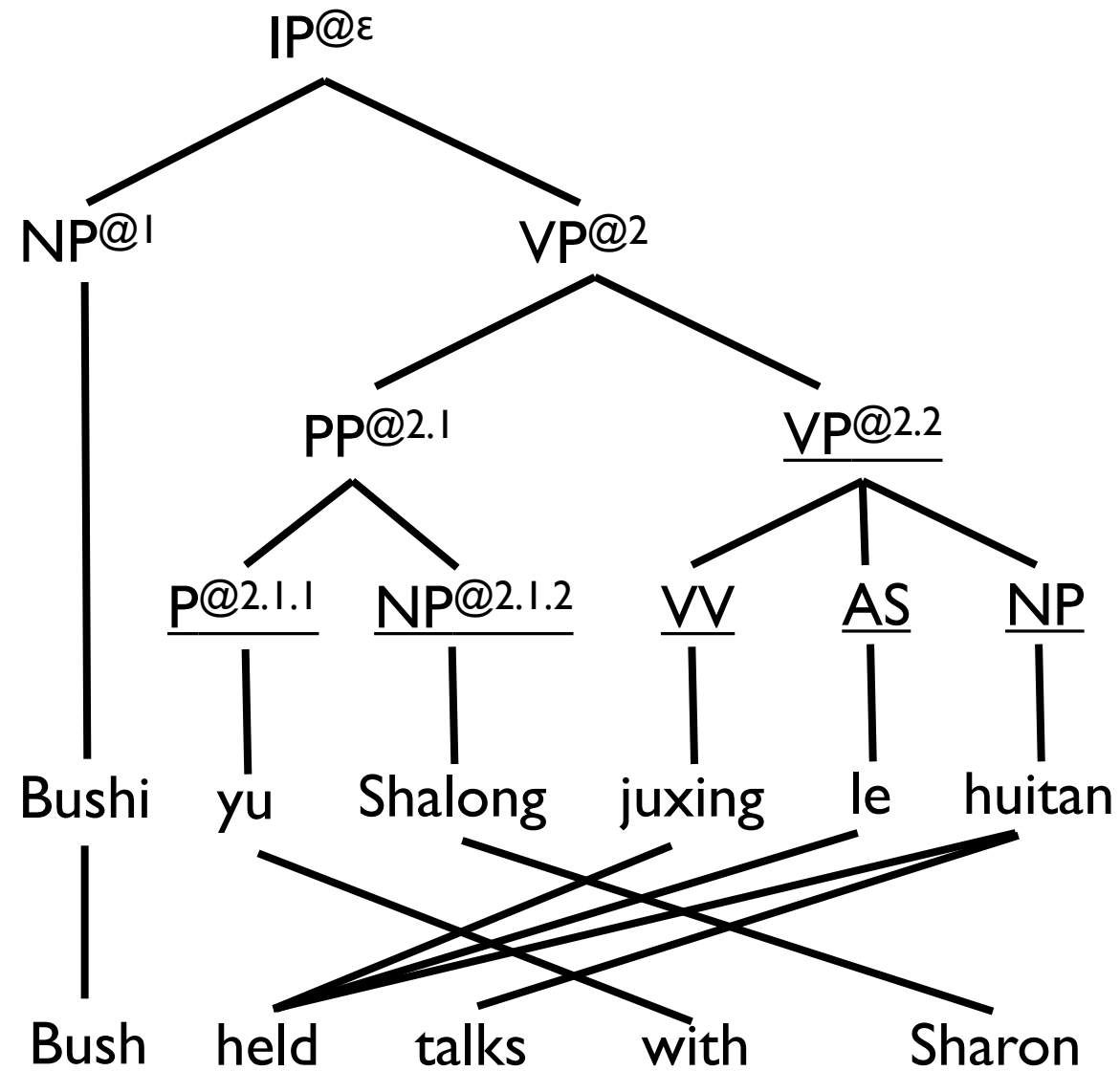
$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1, r_3) \cdot P(r_4|r_1, r_3) \\ & P(r_6|r_1, r_3, r_4) \cdot P(r_7|r_1, r_3, r_4) \end{aligned}$$

Probability of a Derivation Tree

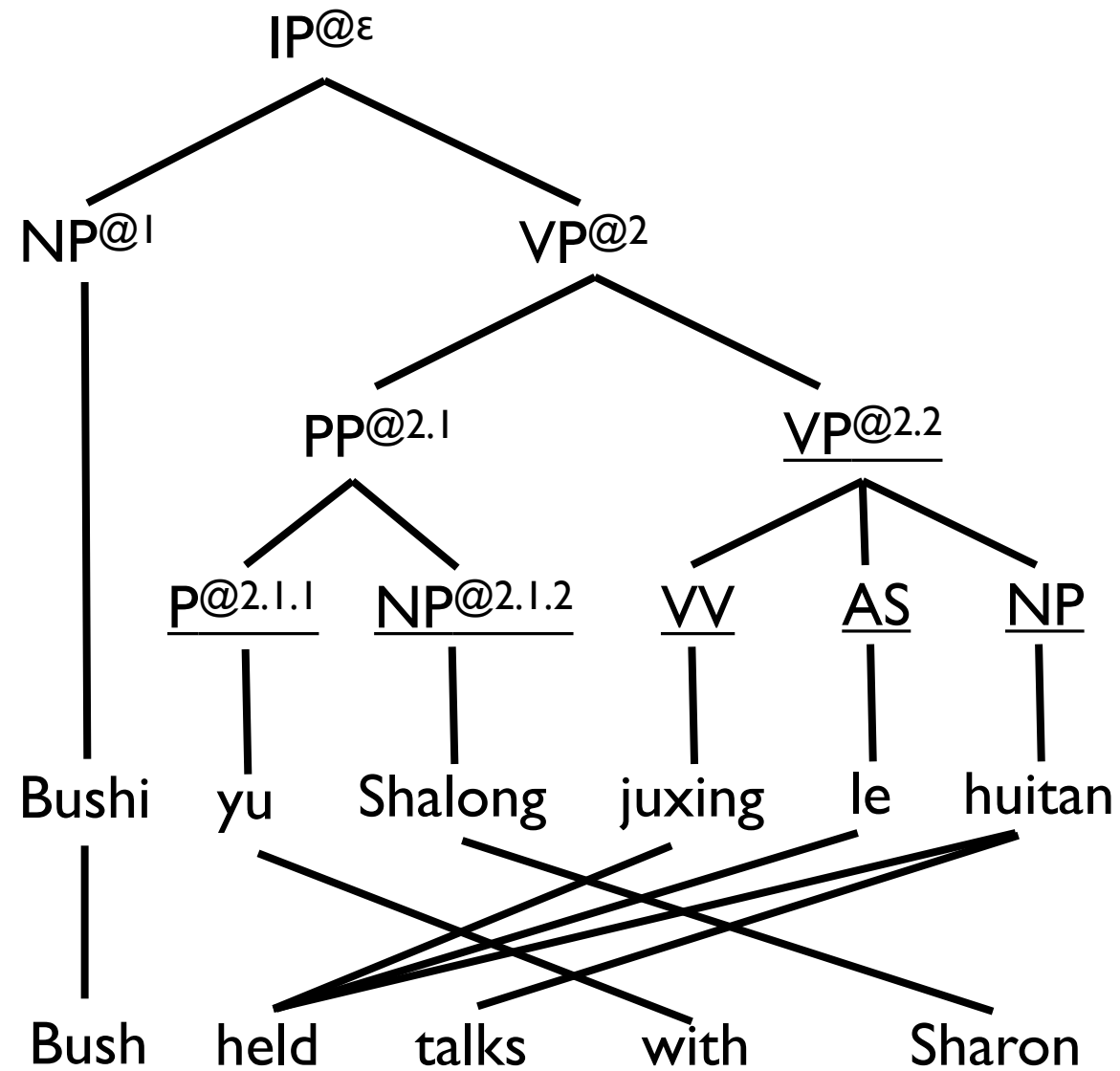


$$\begin{aligned} P(T) = & P(r_1|\varepsilon) \cdot P(r_2|r_1) \cdot P(r_3|r_1) \\ & P(r_5|r_1,r_3) \cdot P(r_4|r_1,r_3) \\ & P(r_6|r_1,r_3,r_4) \cdot P(r_7|r_1,r_3,r_4) \end{aligned}$$

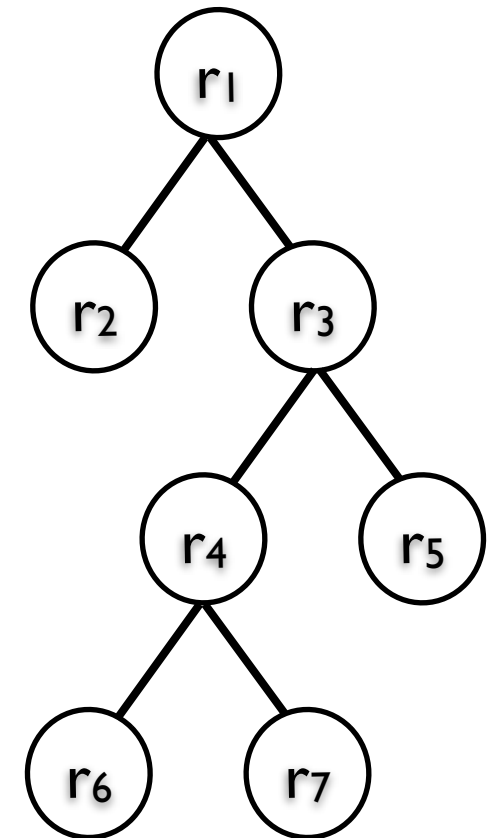
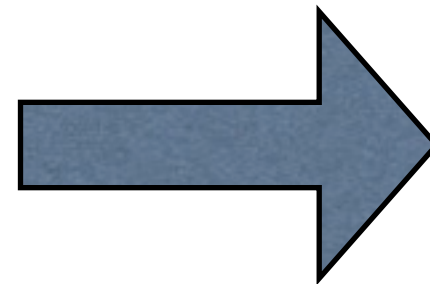
Learning Rule Markov Models



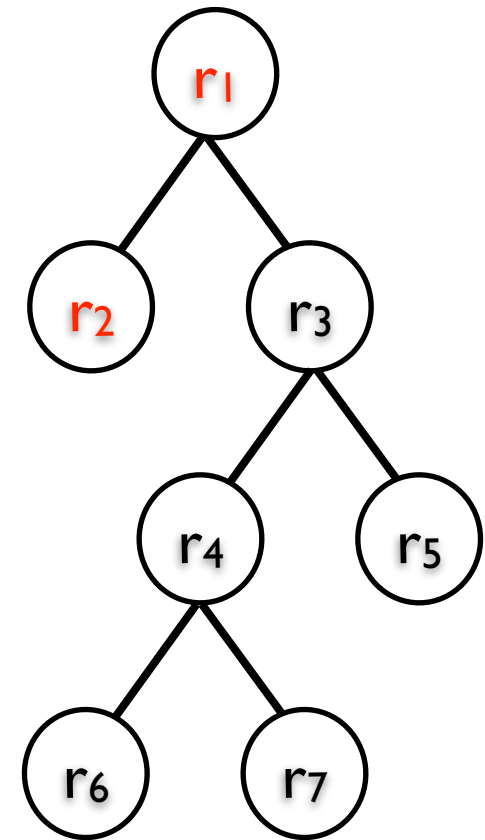
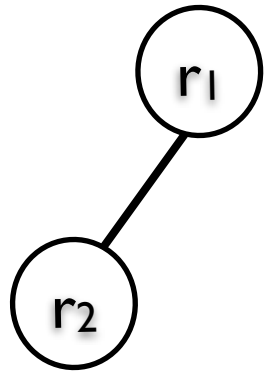
Learning Rule Markov Models



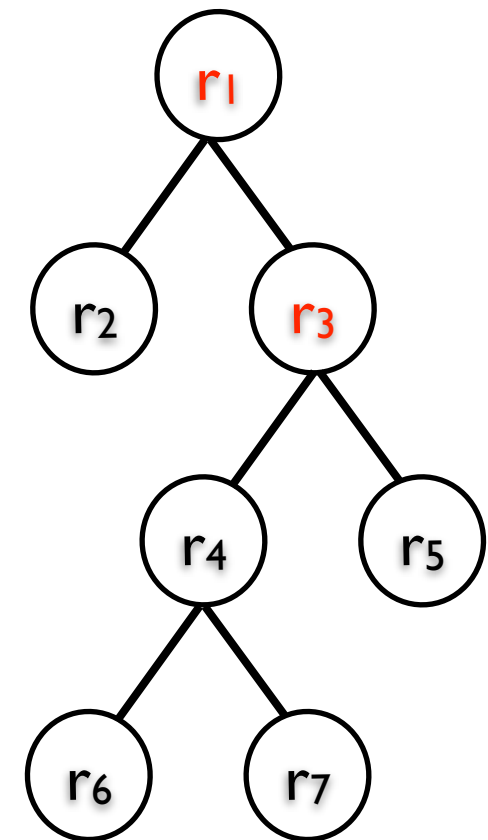
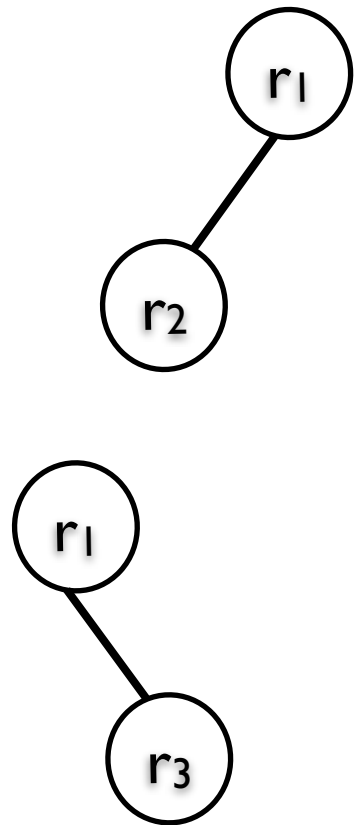
Galley et al.



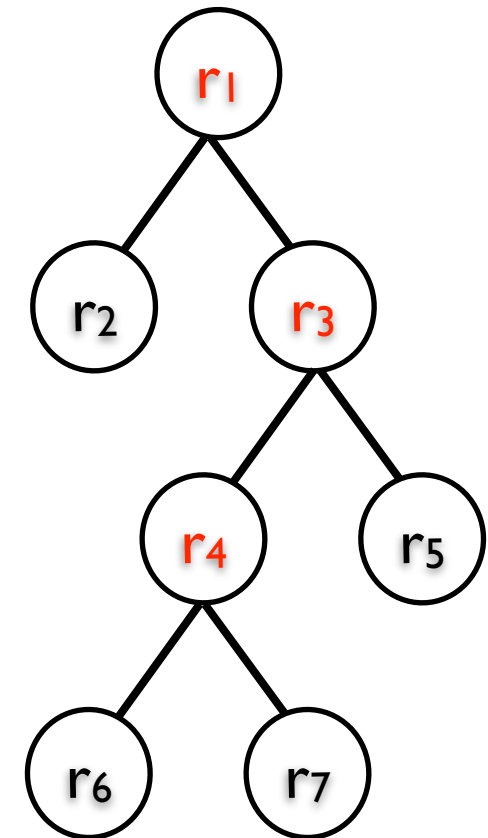
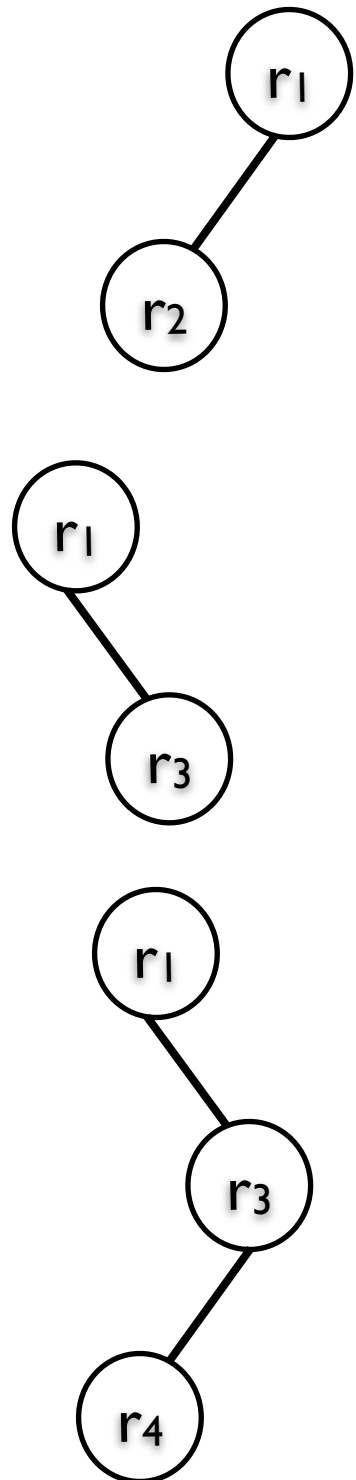
Learning Rule Markov Models



Learning Rule Markov Models



Learning Rule Markov Models



Smaller rule Markov models

Smaller rule Markov models

- RM-A: Keep only those contexts in which more than P items were observed. Tune P for BLEU on dev

Smaller rule Markov models

- RM-A: Keep only those contexts in which more than P items were observed. Tune P for BLEU on dev
- RM-B: Keep only those contexts which were observed more than P times. Tune P for BLEU on dev

Smaller rule Markov models

- RM-A: Keep only those contexts in which more than P items were observed. Tune P for BLEU on dev
- RM-B: Keep only those contexts which were observed more than P times. Tune P for BLEU on dev
- RM-C: Build Prediction Suffix Trees using the approach of Bejerano and Yona (1999)

Smaller rule Markov models

- RM-A: Keep only those contexts in which more than P items were observed. Tune P for BLEU on dev
- RM-B: Keep only those contexts which were observed more than P times. Tune P for BLEU on dev
- RM-C: Build Prediction Suffix Trees using the approach of Bejerano and Yona (1999)

Incremental decoding with RMMs

Incremental decoding with RMMs

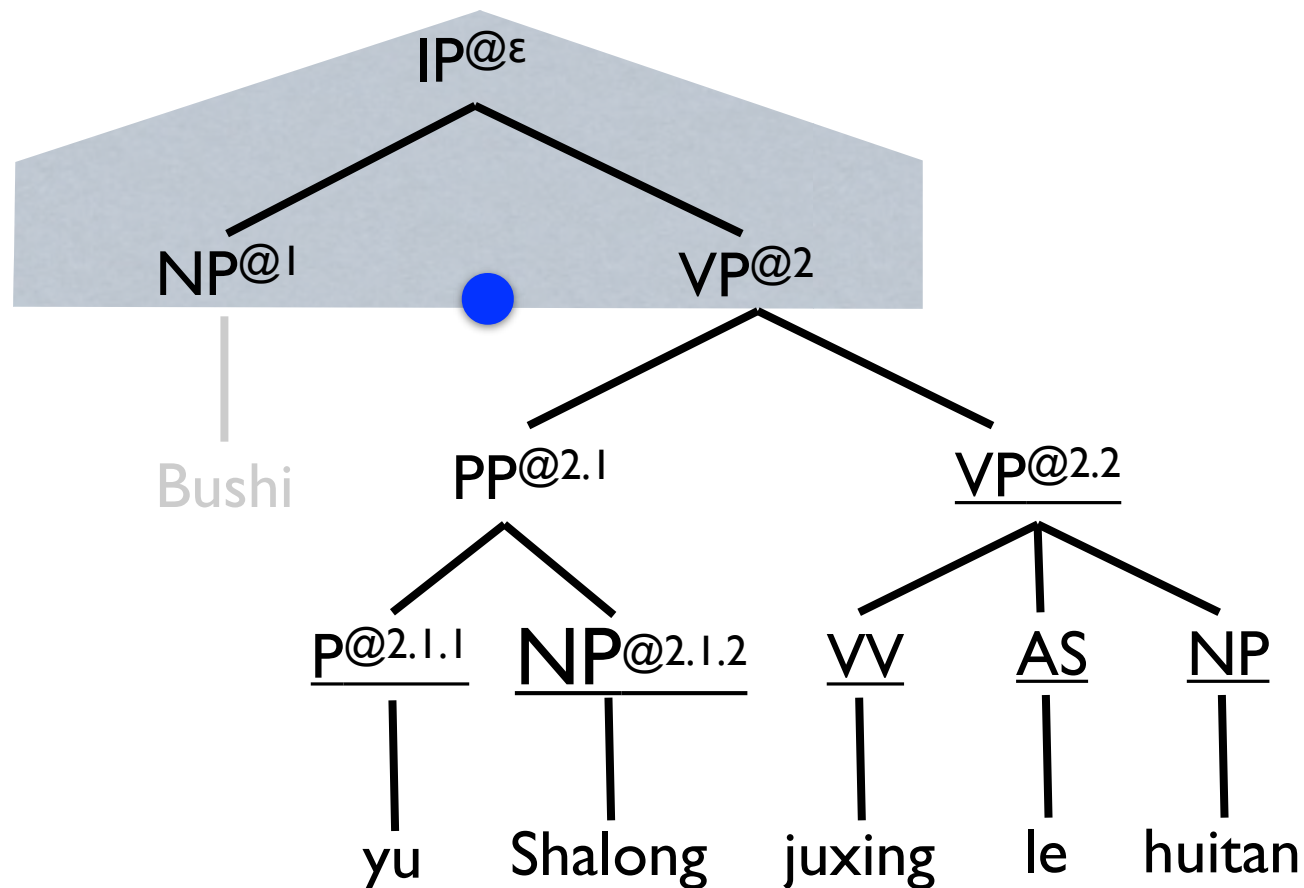
- rule Markov models using only vertical context are a natural fit for Incremental decoding (Huang and Mi, 2010)

Incremental decoding with RMMs

- rule Markov models using only vertical context are a natural fit for Incremental decoding (Huang and Mi, 2010)
- The top-down decoder maintains the vertical context as part of its state which is used in predicting the next rule

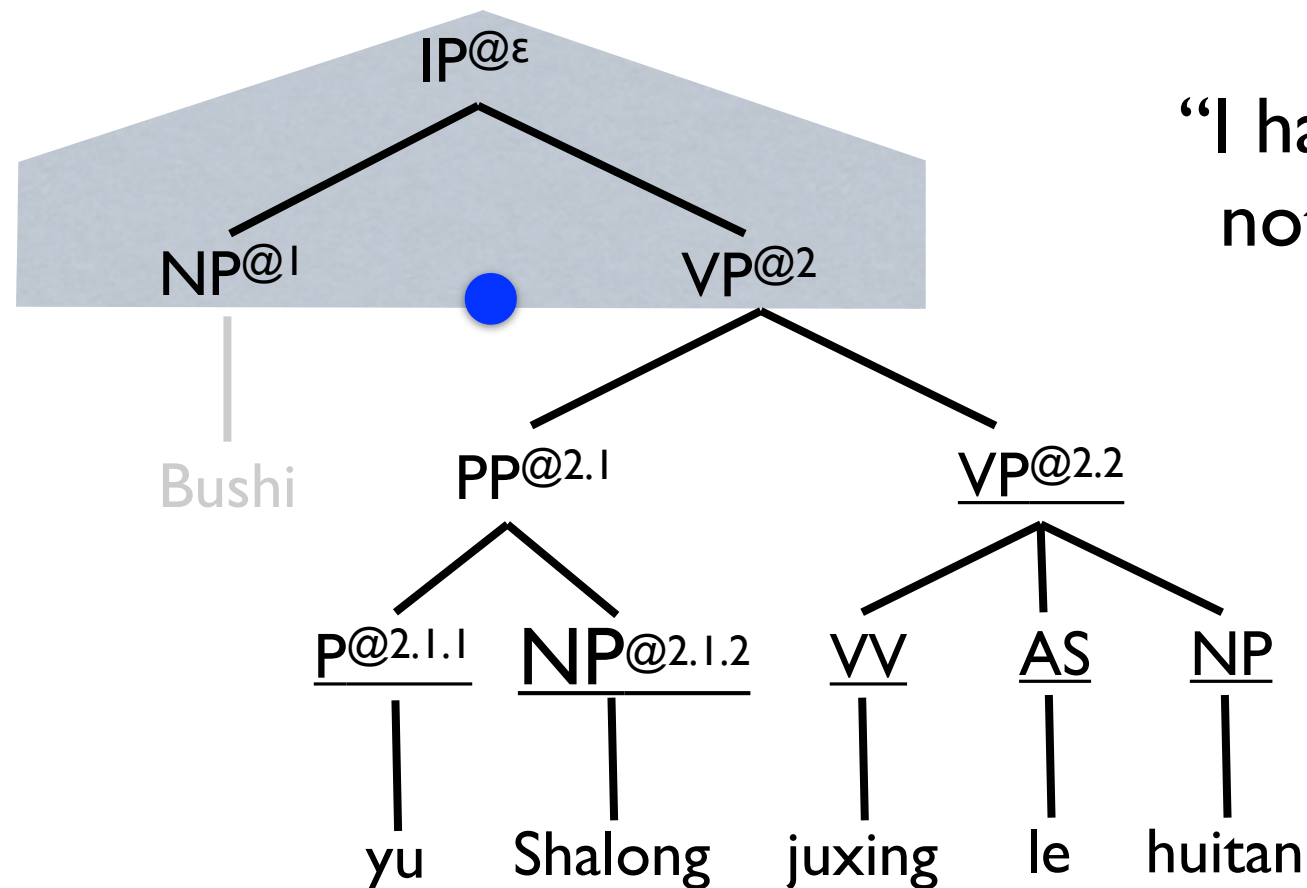
Incremental Decoding

- stack (*active* derivation history): $[\cdot IP@^\epsilon] [NP@^1 \cdot VP@^2]$
- three colors for nodes: white (uncovered), grey (partially covered), and black (covered)



Incremental Decoding

- stack (*active* derivation history): $[\cdot IP@^\epsilon] [NP@^1 \cdot VP@^2]$
- three colors for nodes: white (uncovered), grey (partially covered), and black (covered)



“I have finished NP subtree but not started with VP subtree”

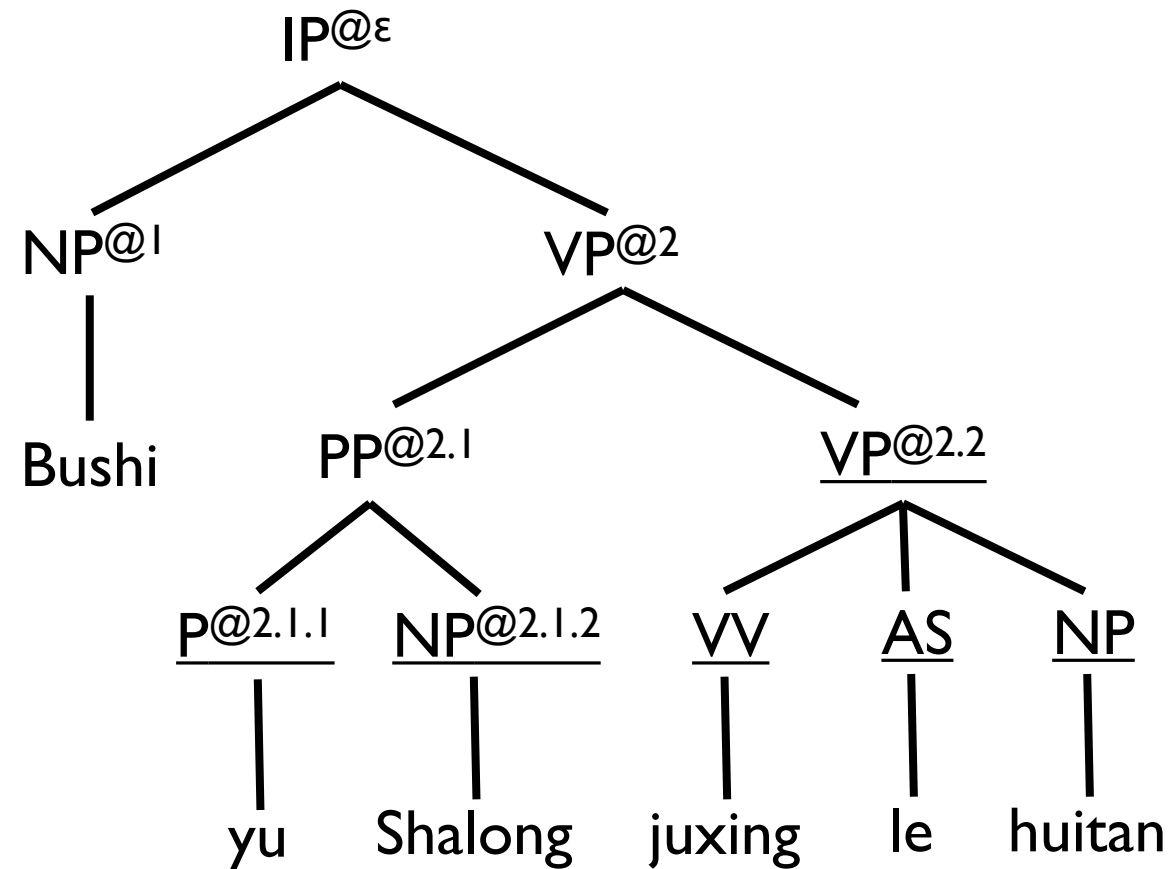
Example Incremental Decoding

[<s> . IP@ ϵ </s>]

stack

<s>

hypothesis



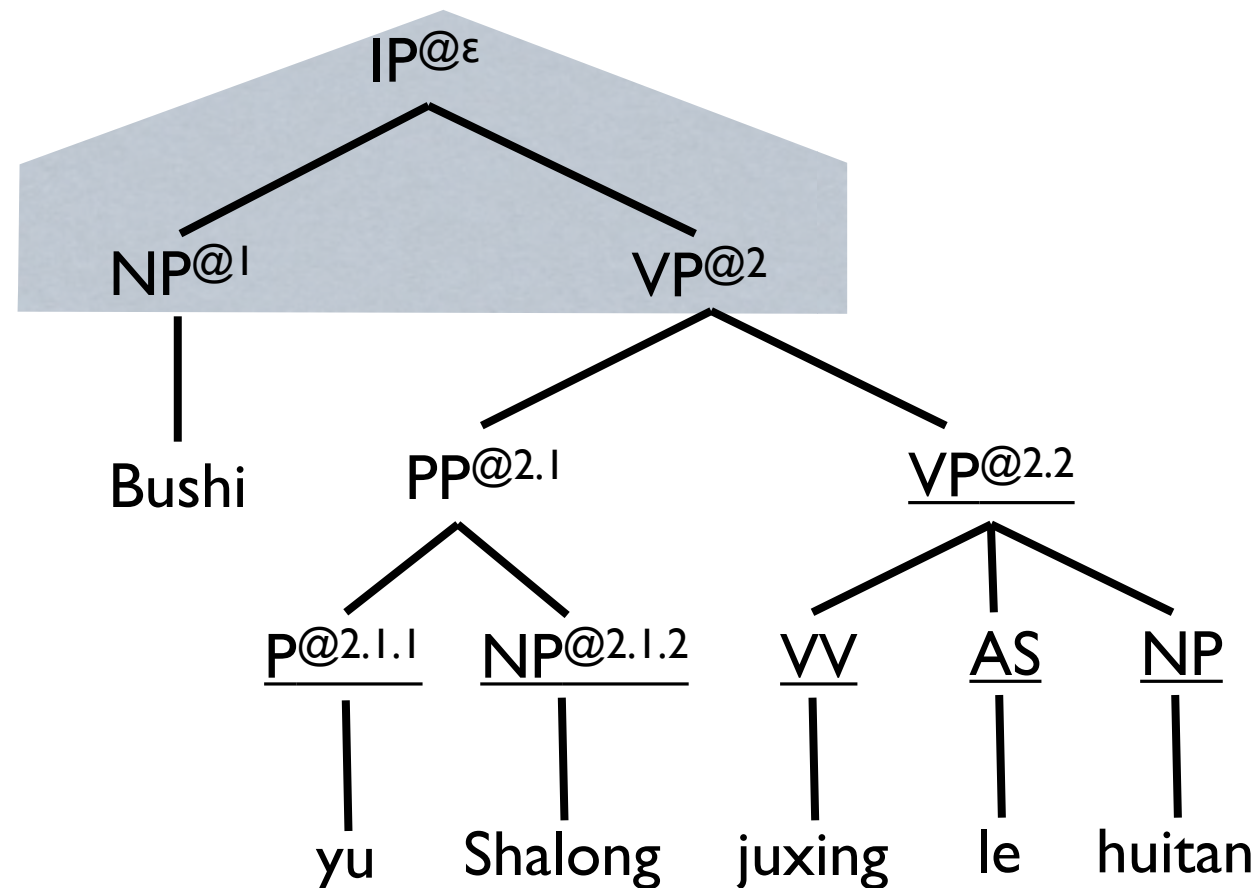
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2]
r₁

stack

<s>

hypothesis



action: predict (push)

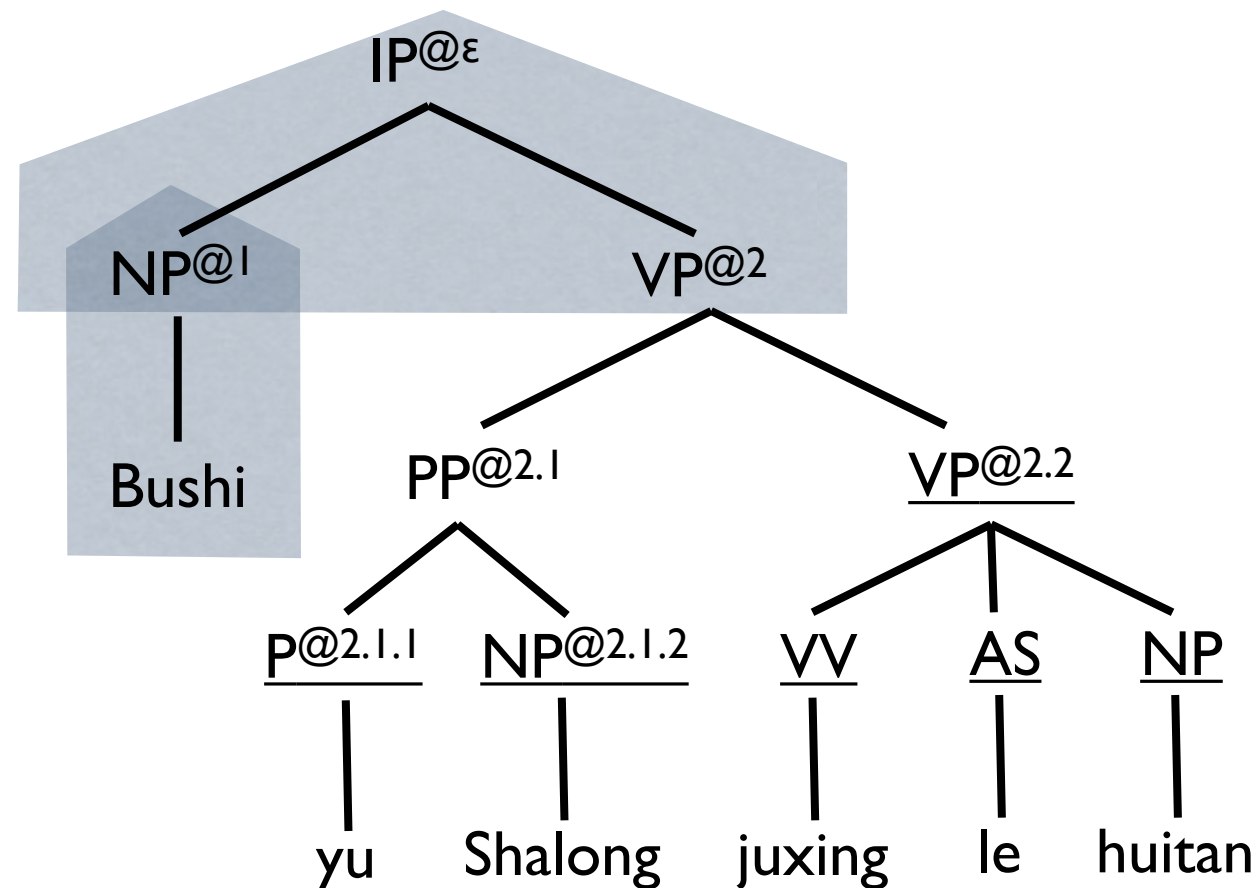
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2] [. Bush]

stack

<s>

hypothesis



action: predict

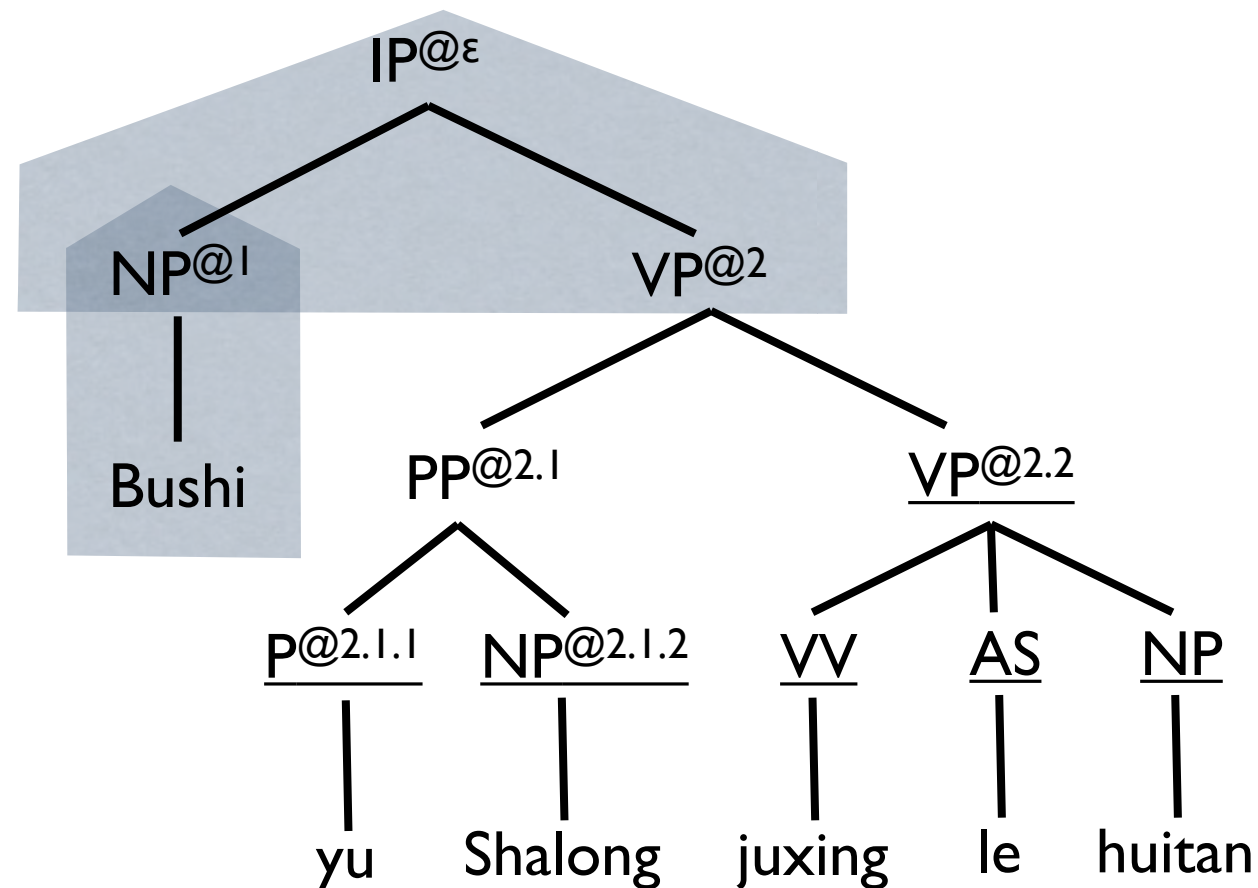
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2] [Bush .]
r1 r2

stack

<s> Bush

hypothesis



action: scan

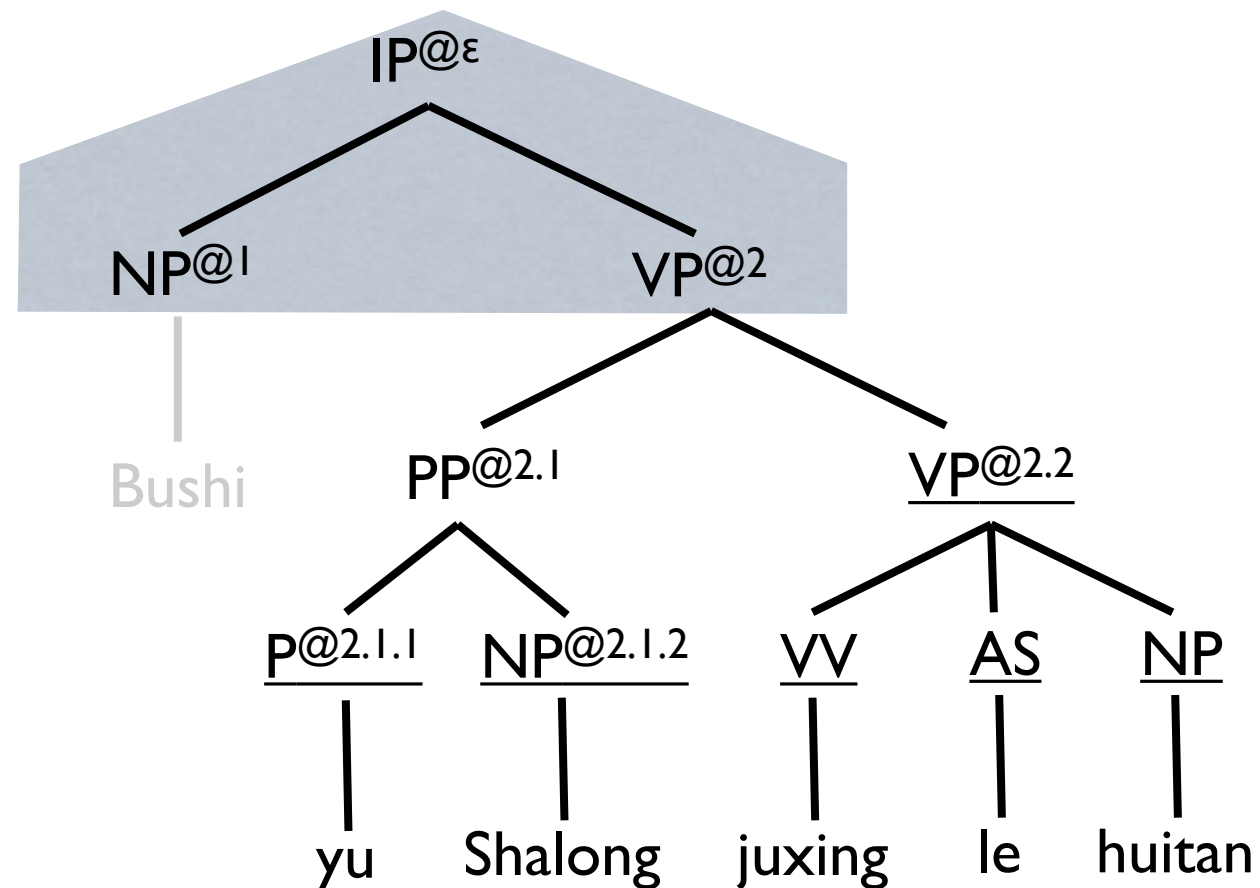
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2]
r1

stack

<s> Bush

hypothesis



action: pop

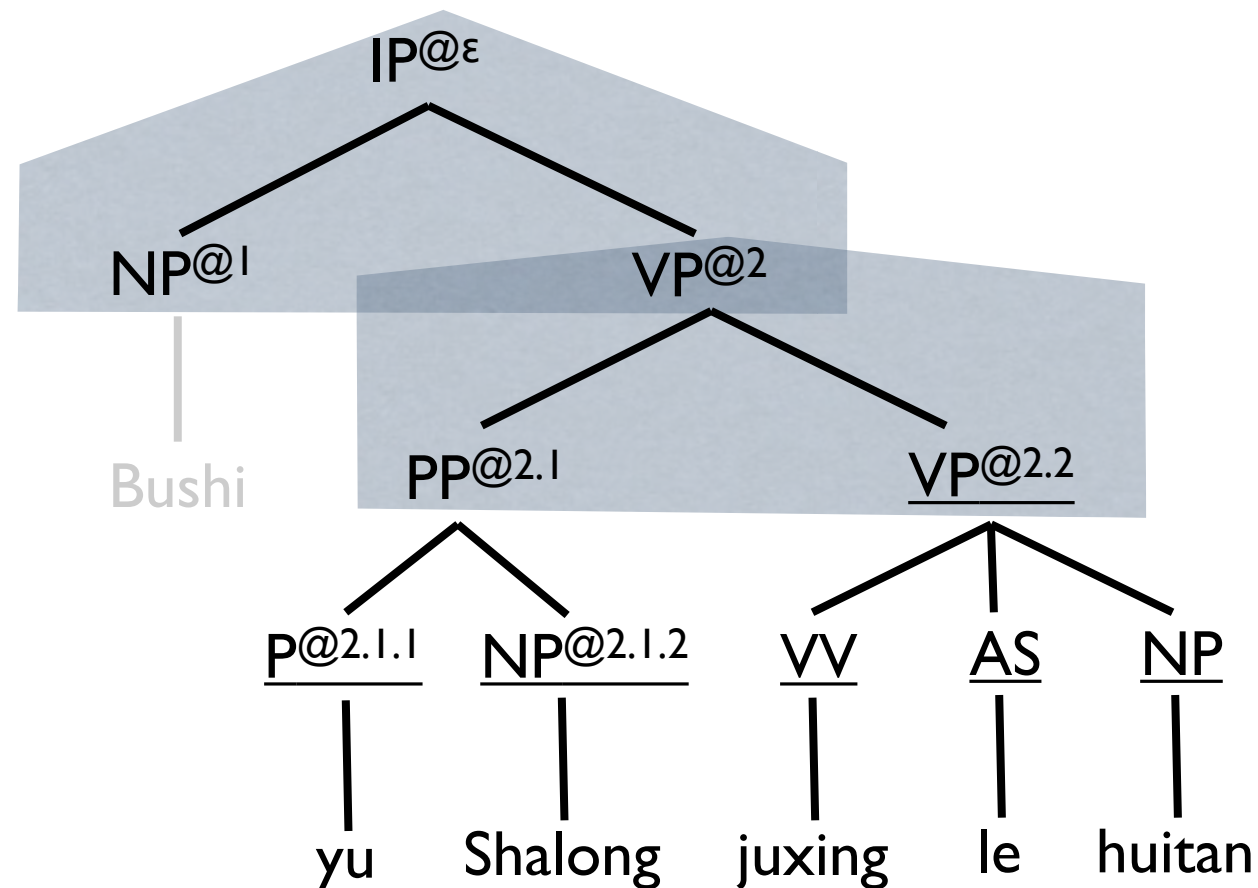
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1]

stack

<s> Bush

hypothesis



action: predict

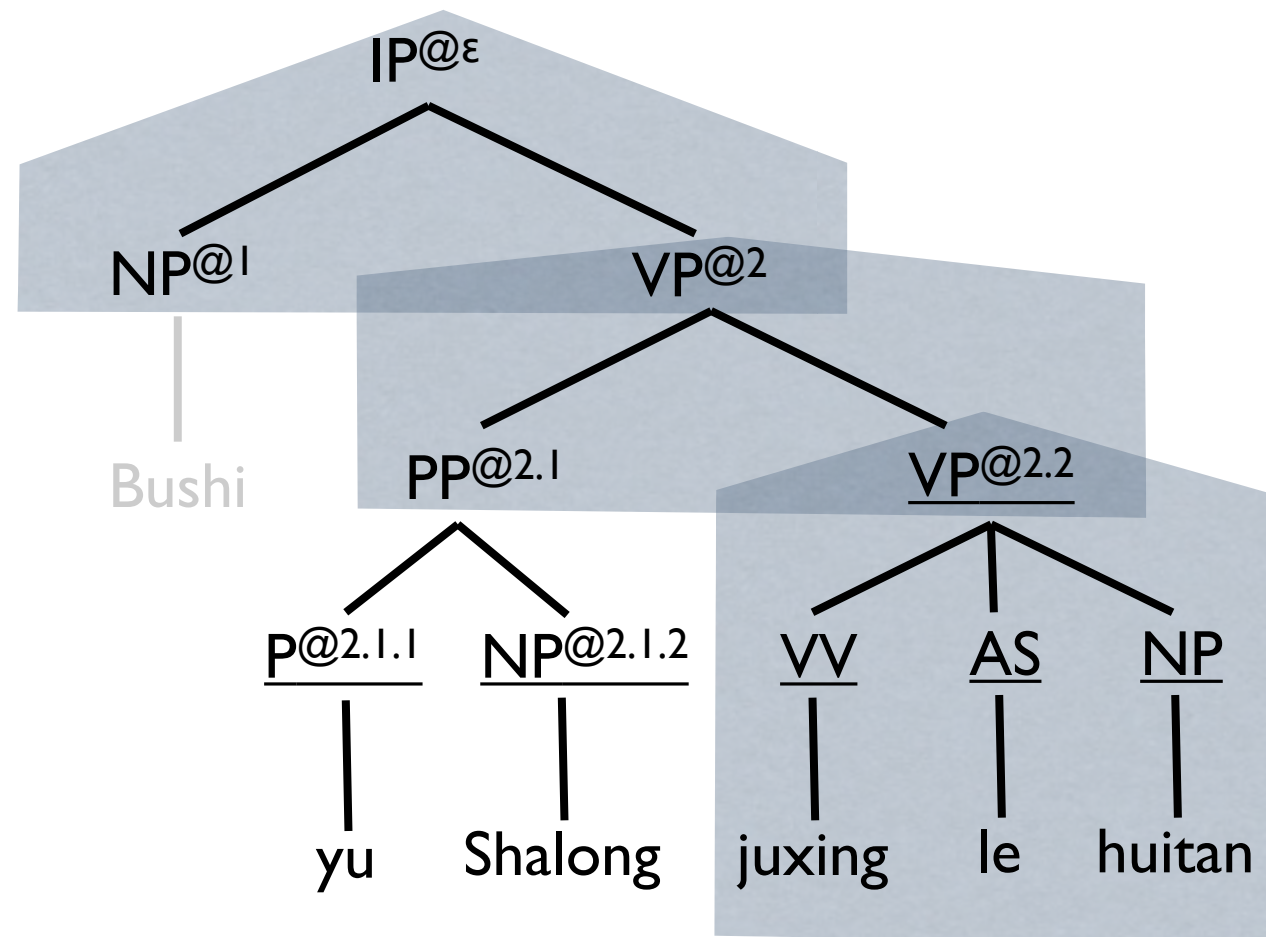
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [. VP@2.2 PP@2.1] [. held talks]

stack

<s> Bush

hypothesis



action: predict

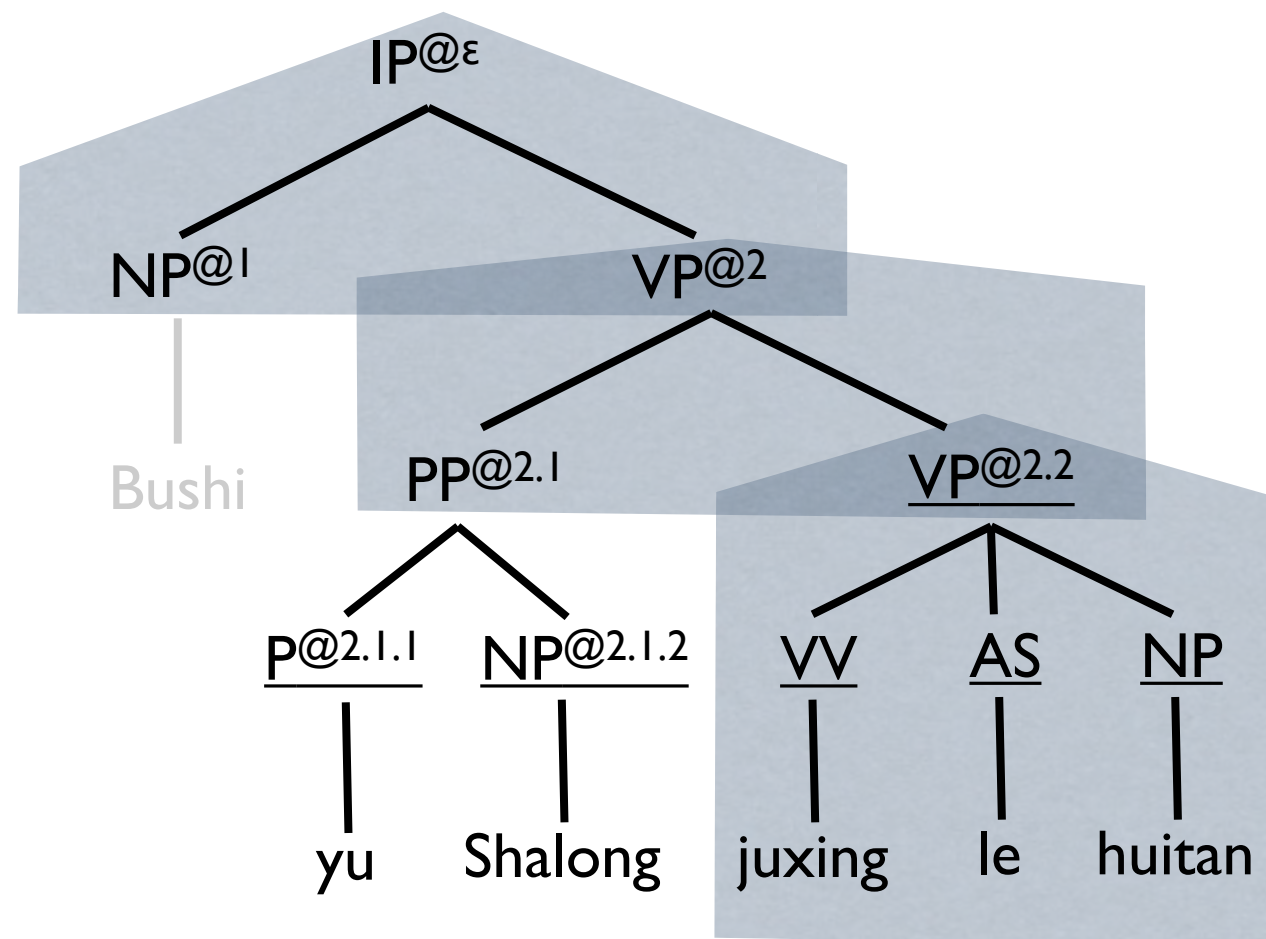
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1] [held talks .]
r1 r3 r5

stack

<s> Bush held talks

hypothesis



action: scan

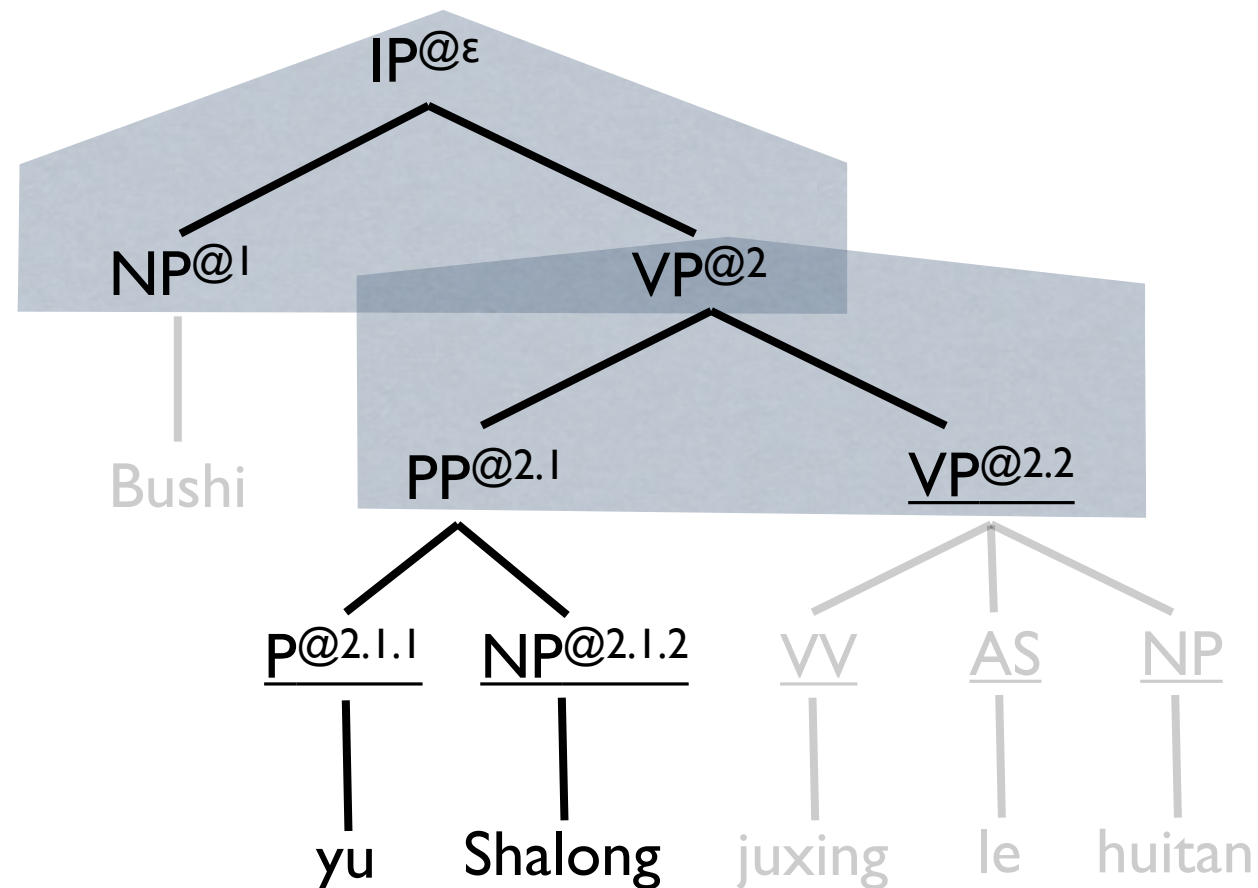
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1]

stack

<s> Bush held talks

hypothesis



action: pop

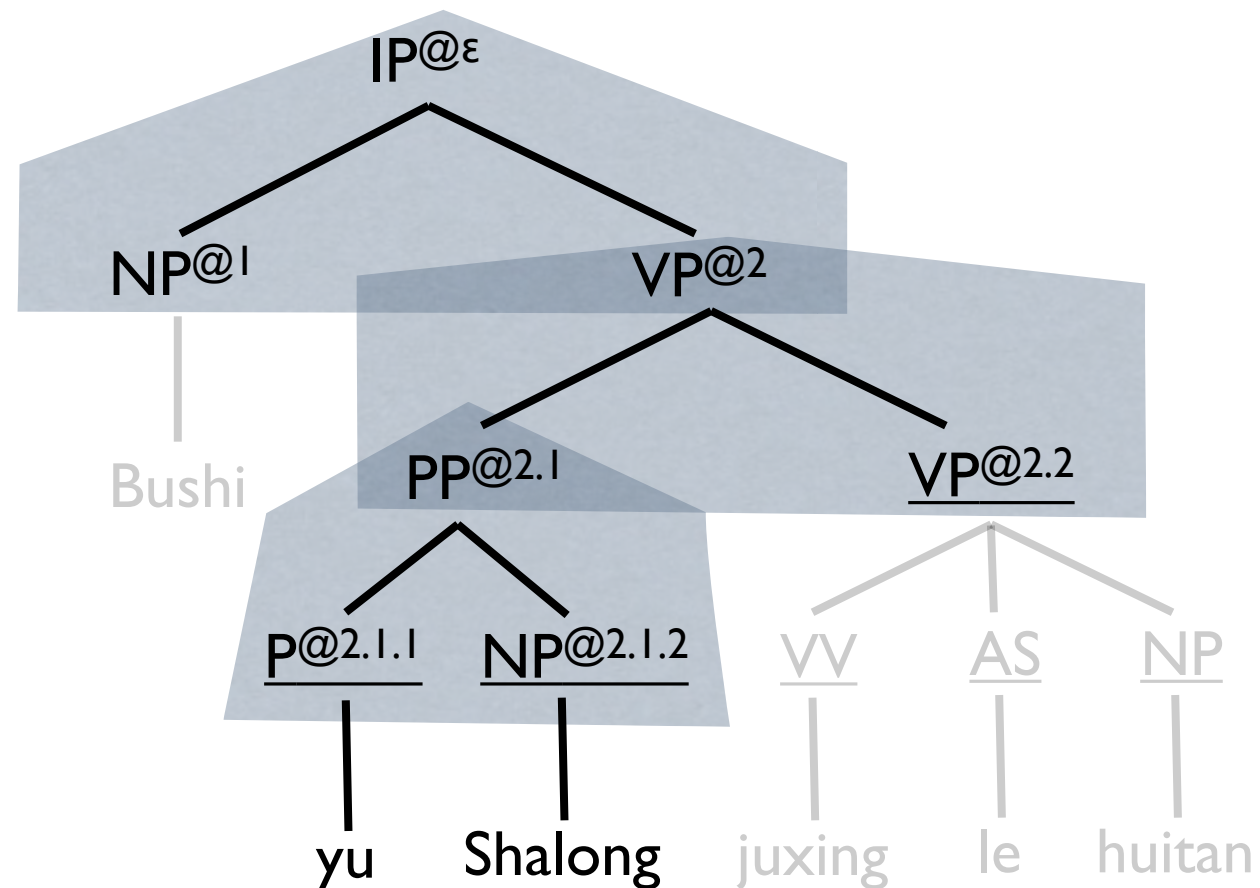
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [. P@2.1.1 NP@2.1.2]

stack

<s> Bush held talks

hypothesis



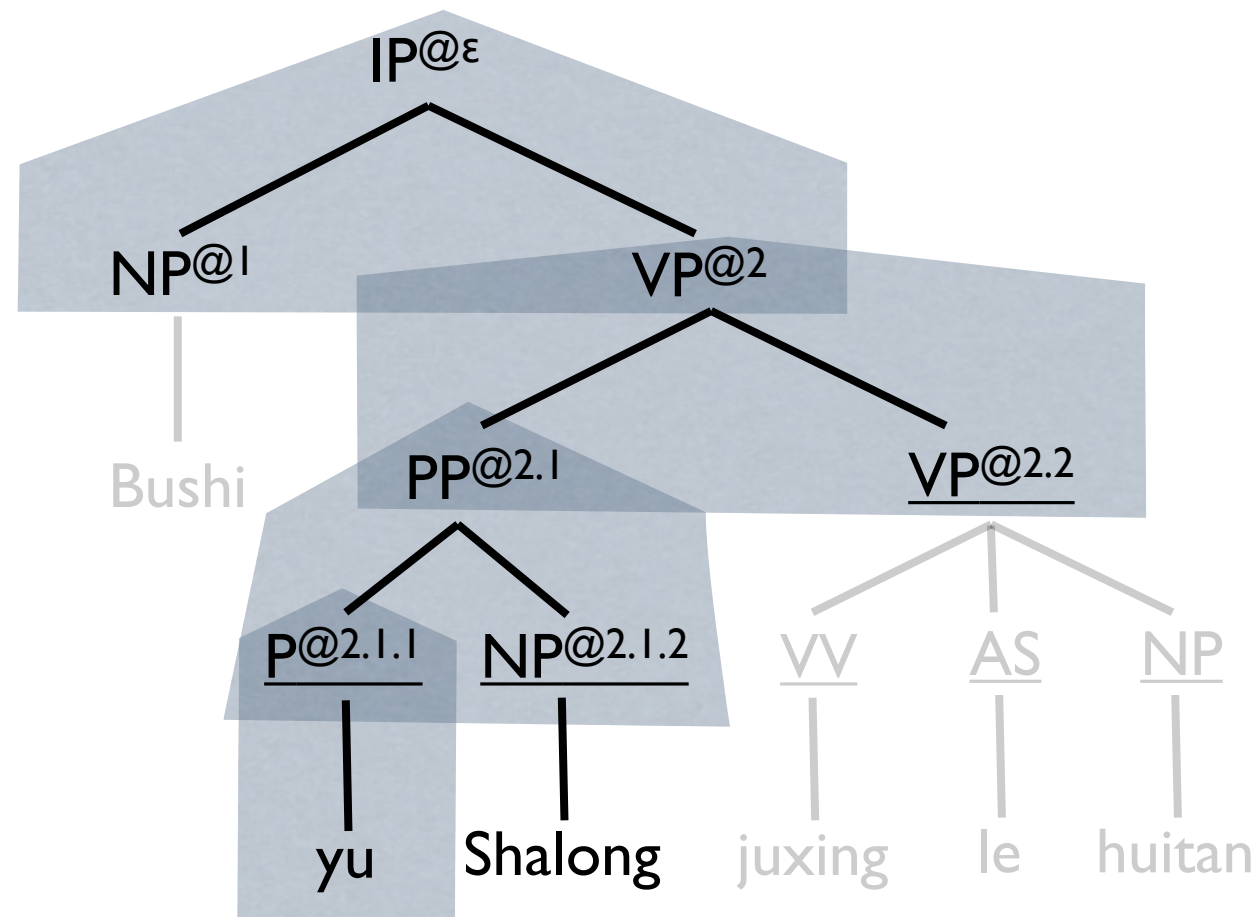
action: predict

Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [. P@2.1.1 NP@2.1.2] [. with/and] *stack*
r₁ r₃ r₄ r₆/r'₆

<s> Bush held talks

hypothesis



action: predict

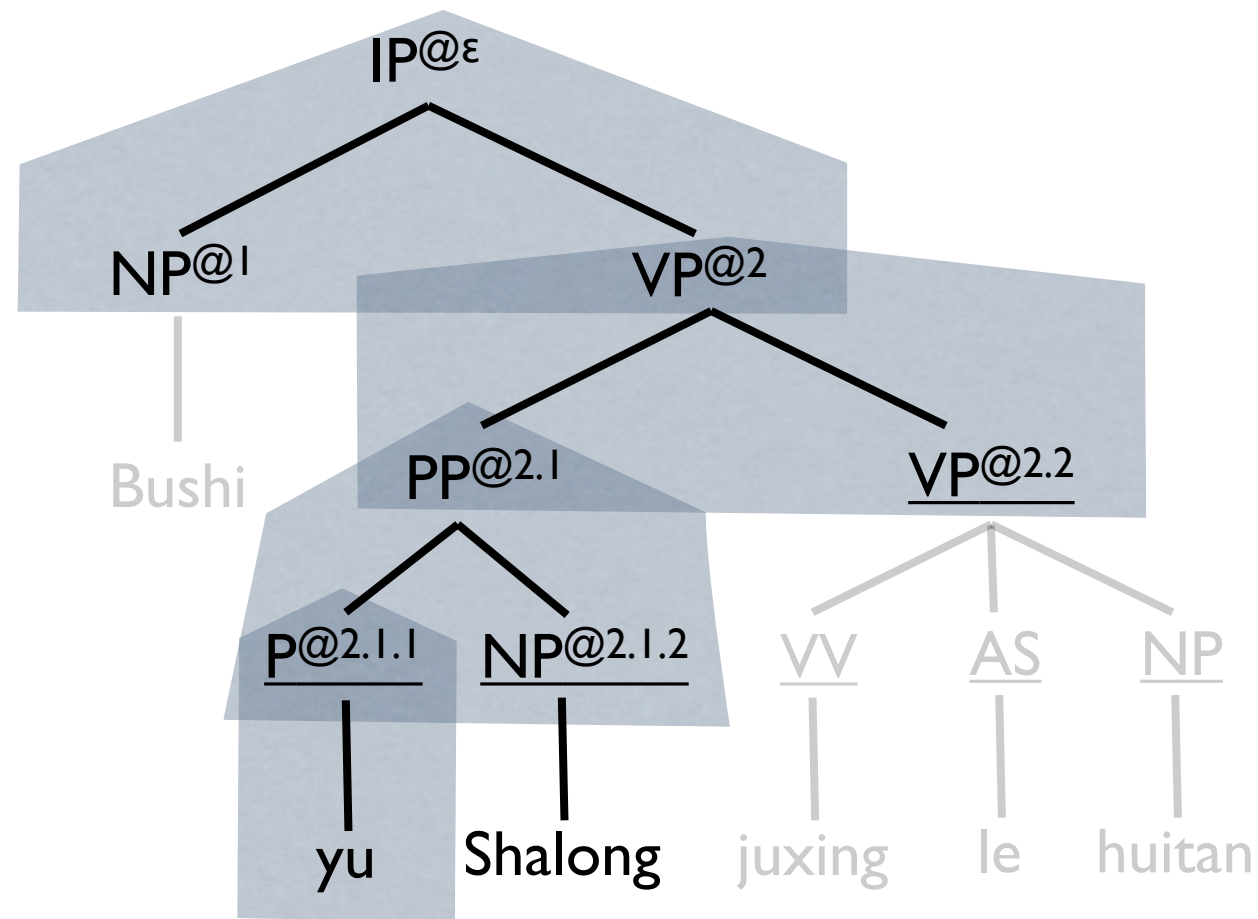
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [. P@2.1.1 NP@2.1.2] [and .]

stack

<s> Bush held talks and

hypothesis



action: scan

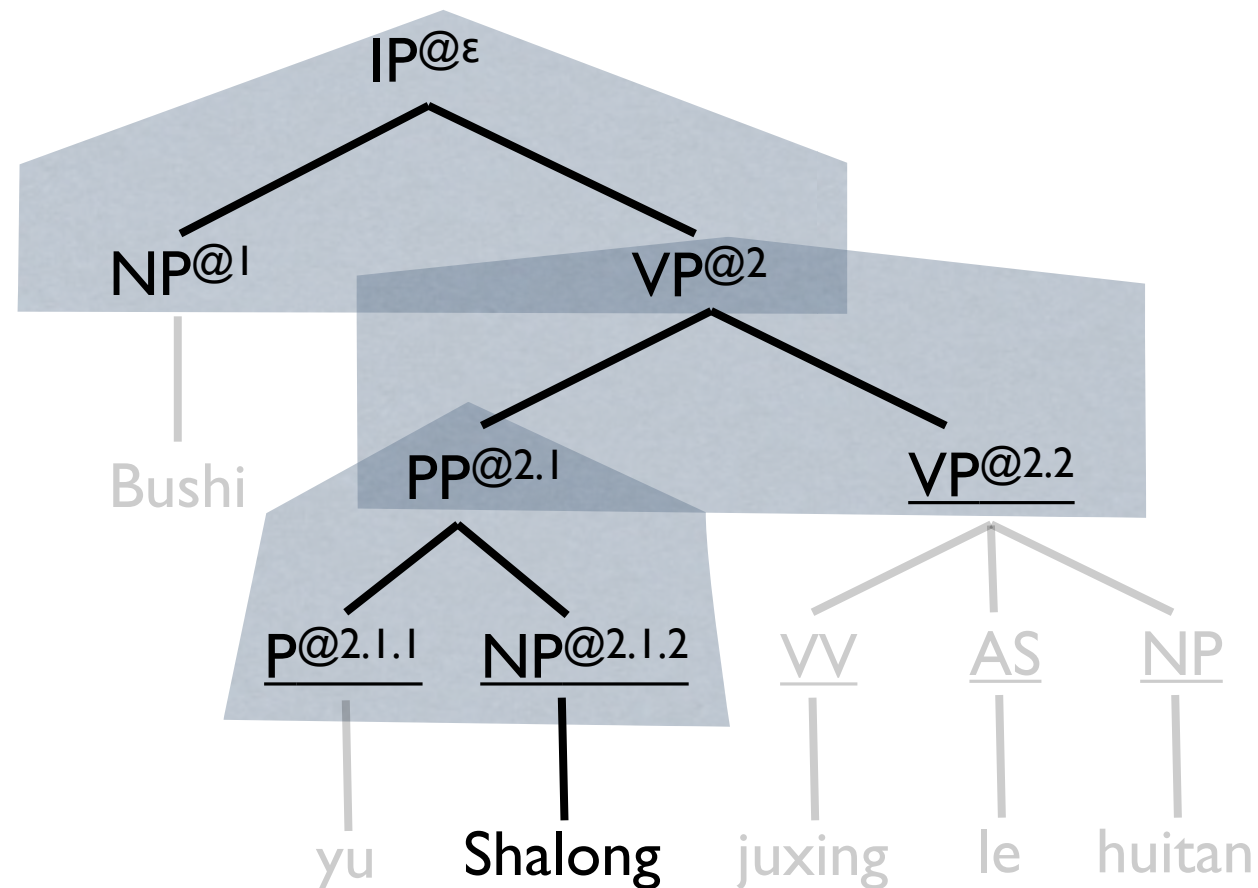
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2]

stack

<s> Bush held talks and

hypothesis



action: pop

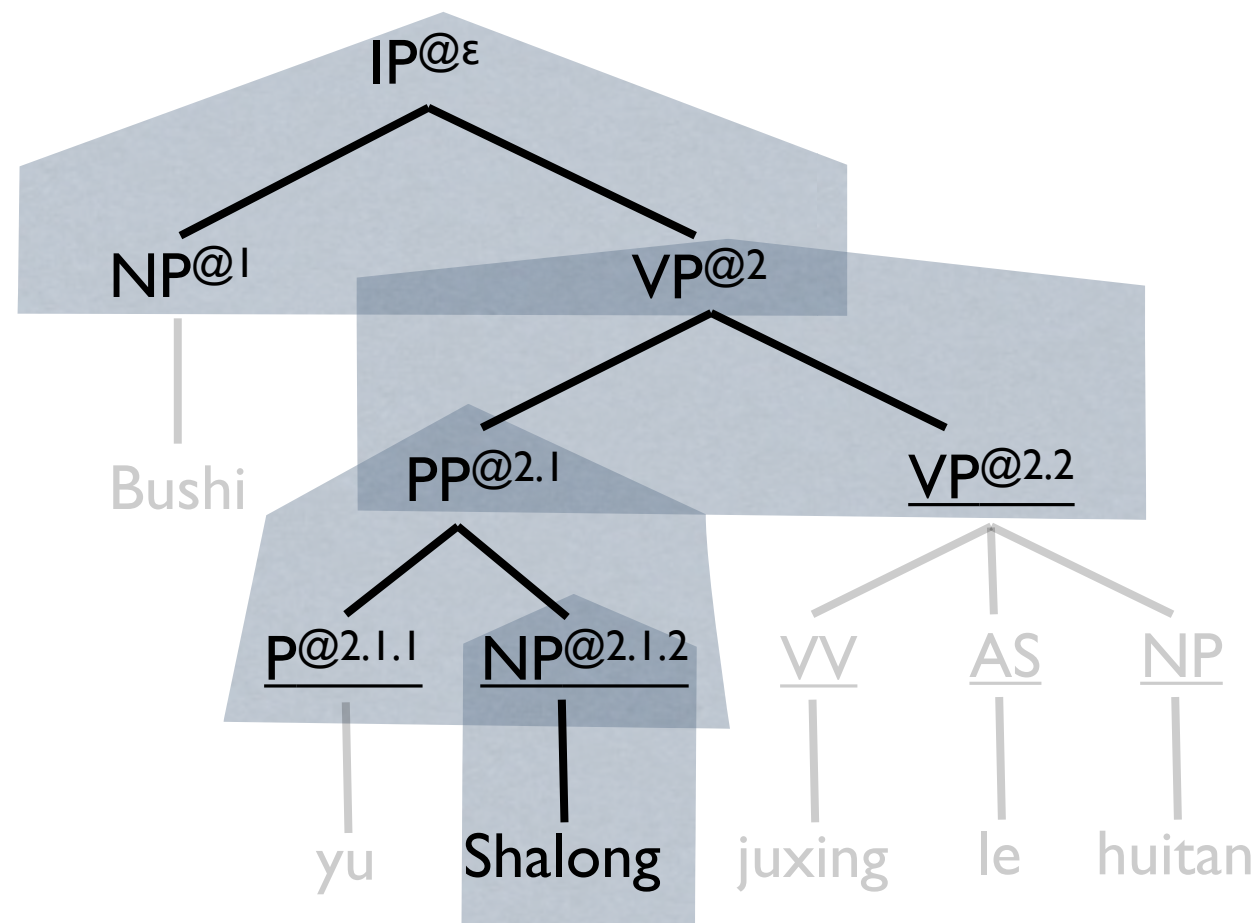
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2] [. Sharon] stack

r_1 r_3 r_4 r_7

<s> Bush held talks and

hypothesis



action: predict

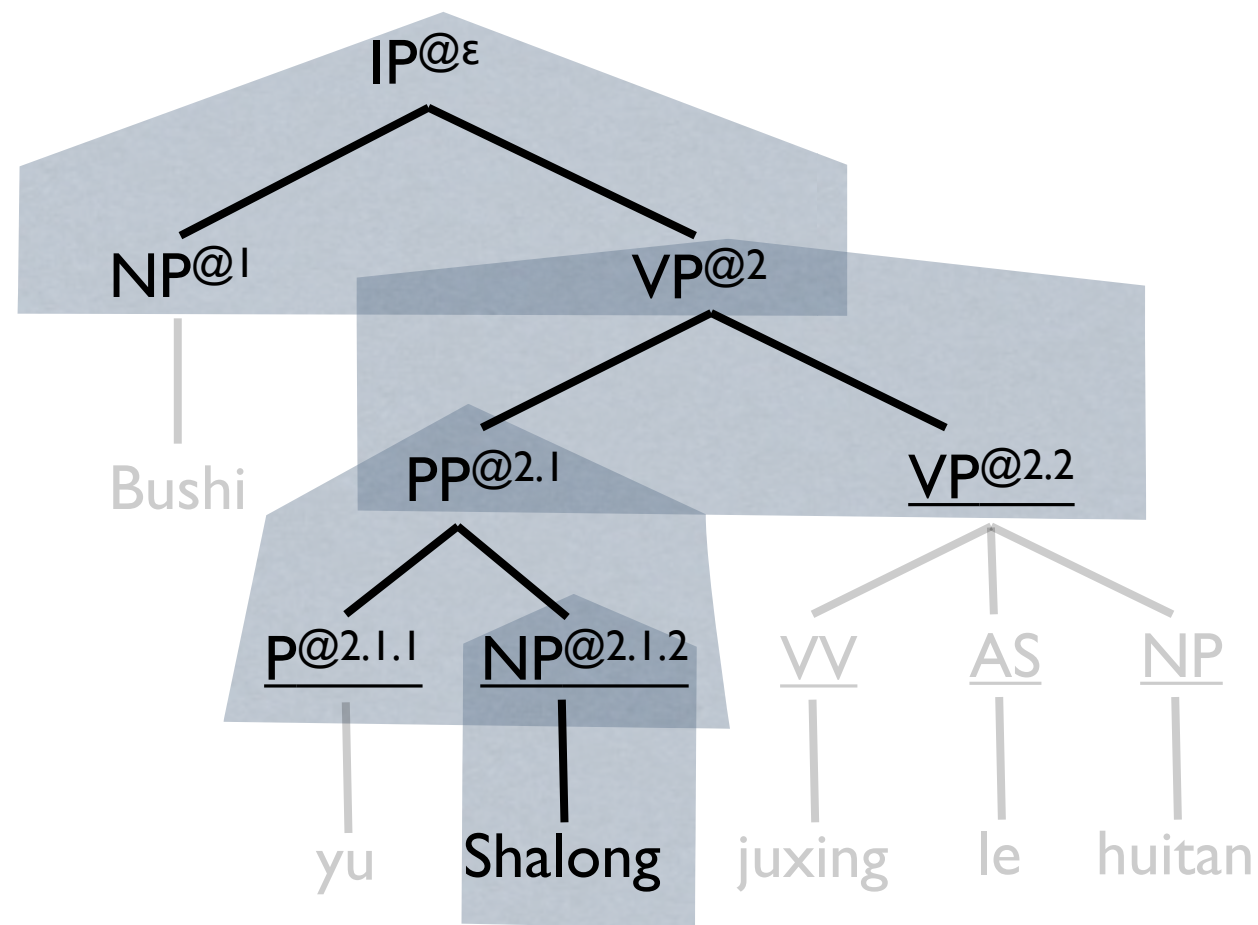
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2] [Sharon .] *stack*

r_1 r_3 r_4 r_7

<s> Bush held talks and Sharon

hypothesis



action: scan

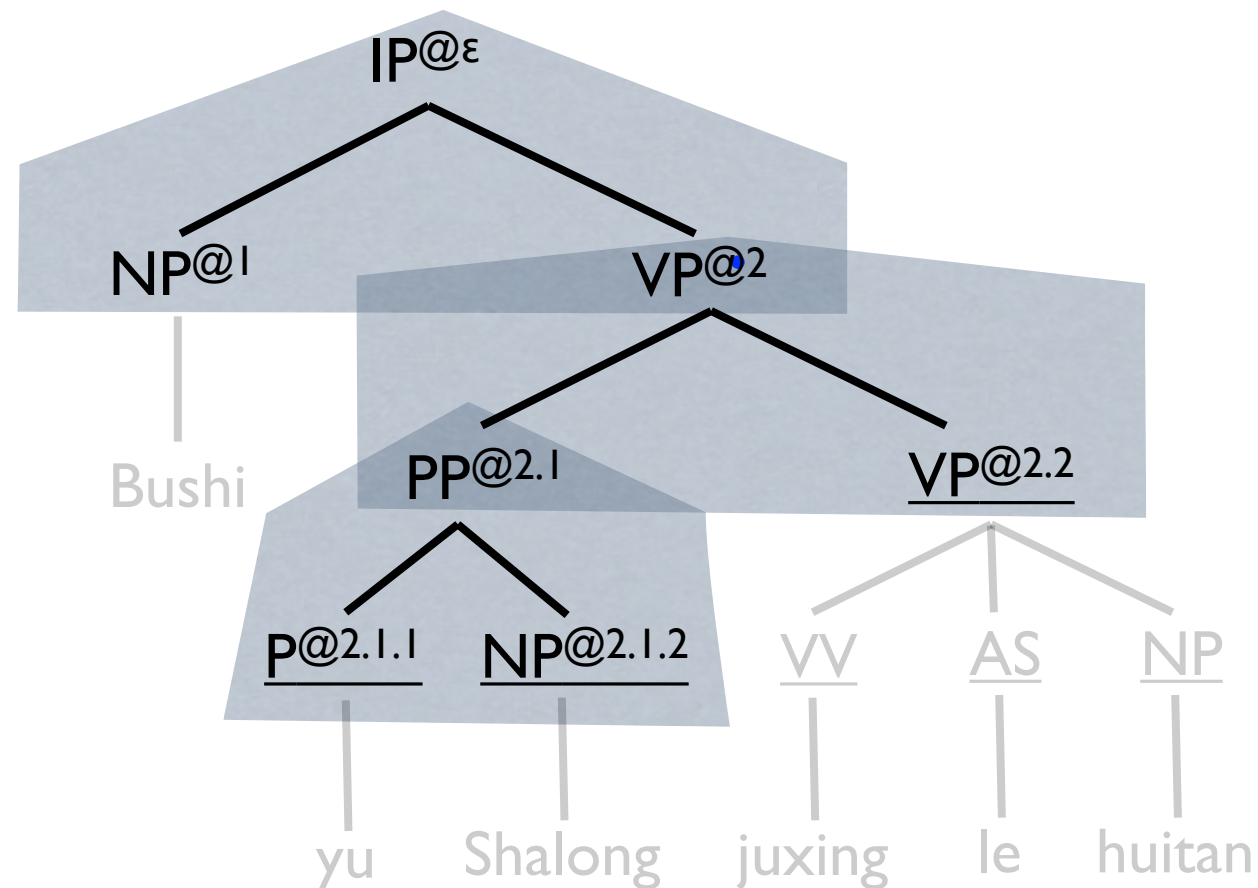
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 NP@2.1.2 .]

stack

<s> Bush held talks and Sharon

hypothesis



action: pop

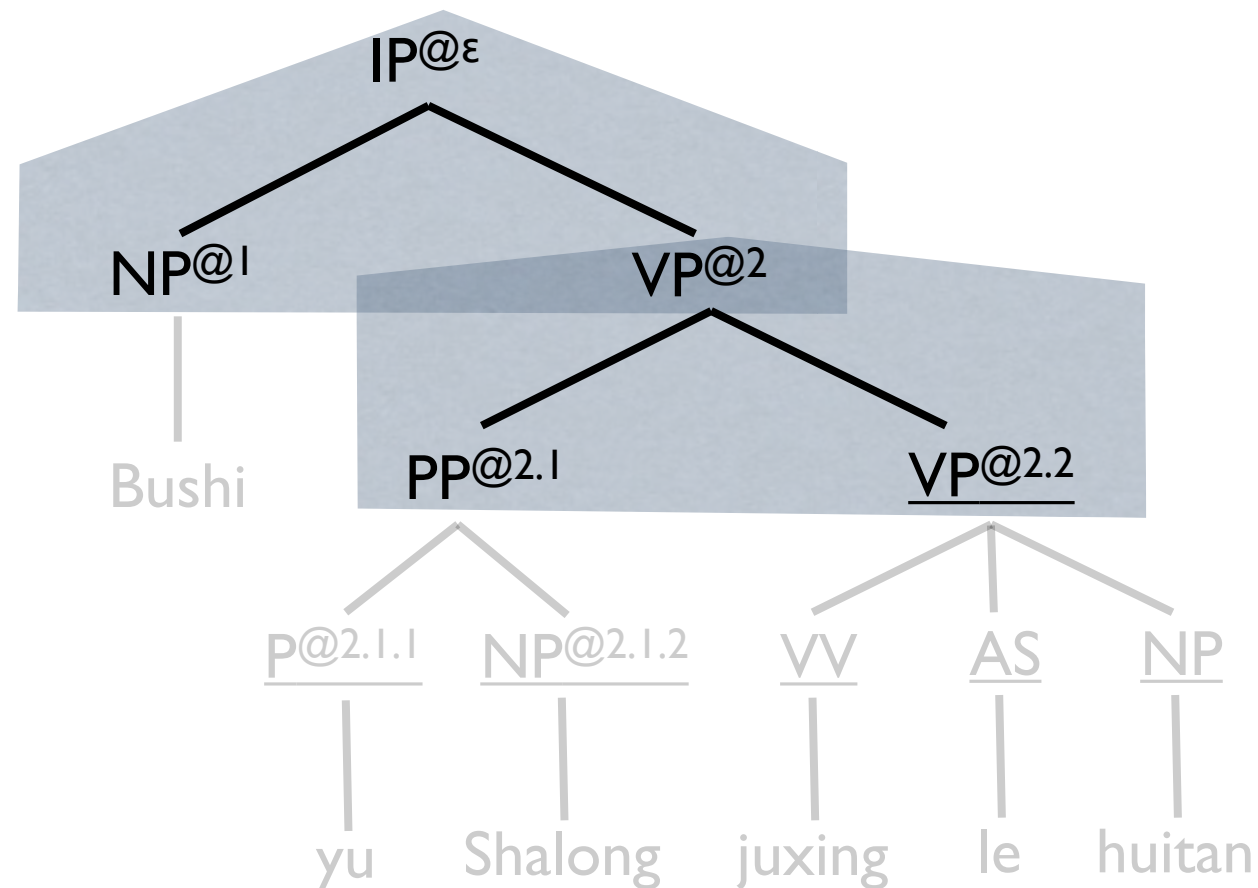
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1 .]
r1 r3

stack

<s> Bush held talks and Sharon

hypothesis



action: pop

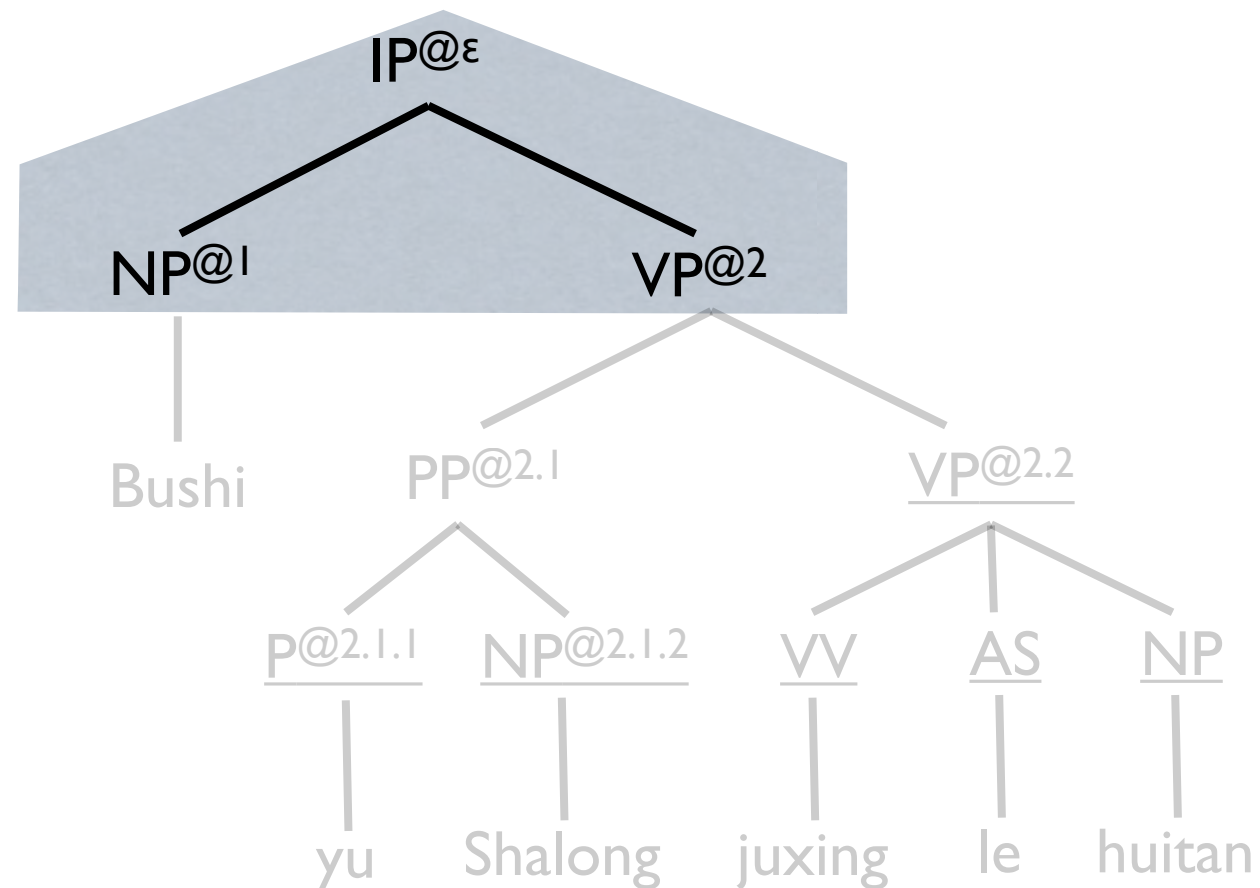
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 VP@2 .]
r₁

stack

<s> Bush held talks and Sharon

hypothesis



action: pop

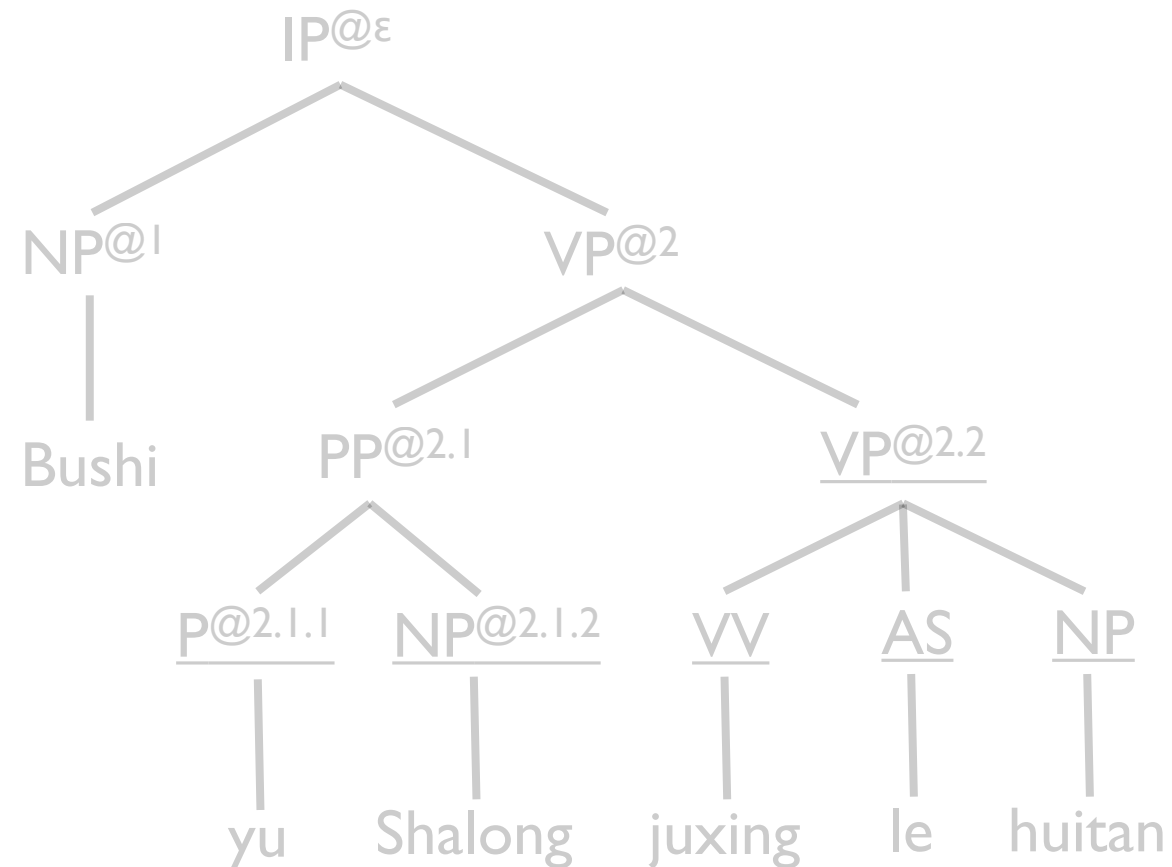
Example Incremental Decoding

[<s> IP@ ϵ .</s>]

stack

<s> Bush held talks and Sharon

hypothesis



action: pop

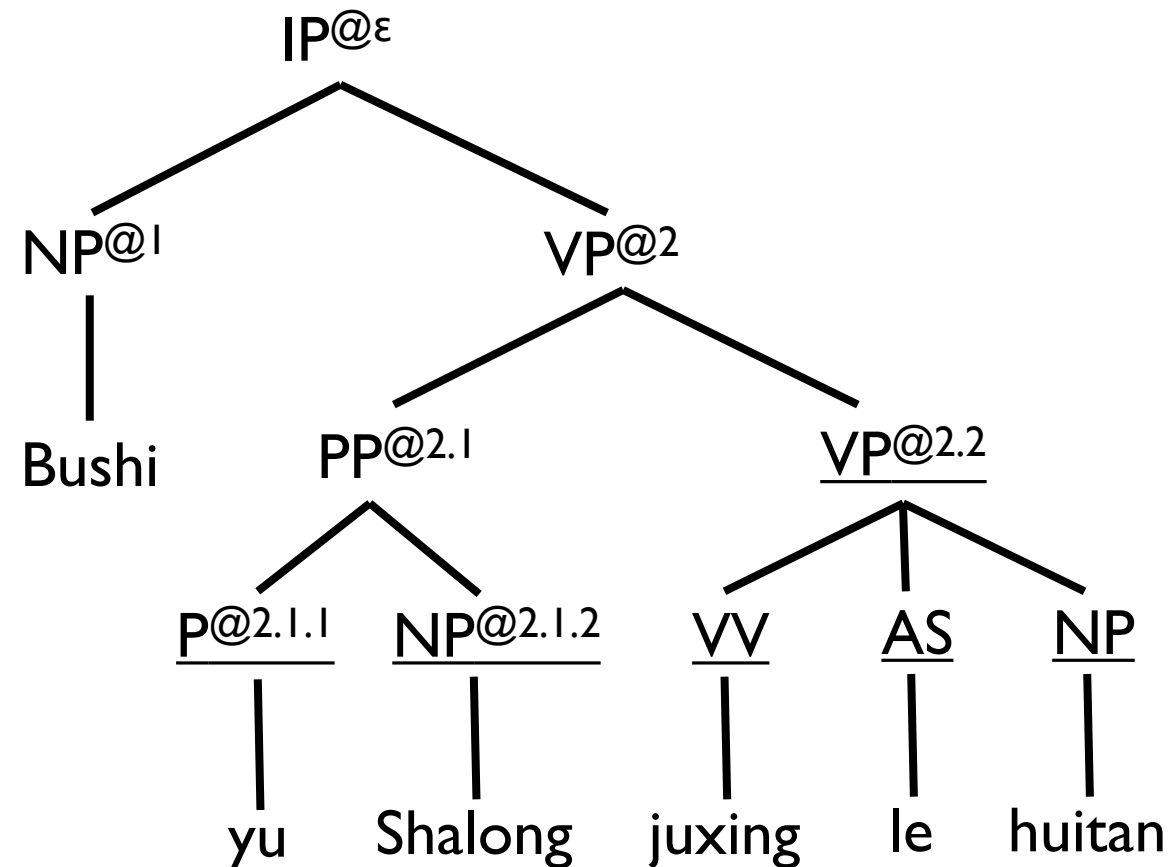
Example Incremental Decoding

[<s> . IP@ ϵ </s>]

stack

<s>

hypothesis



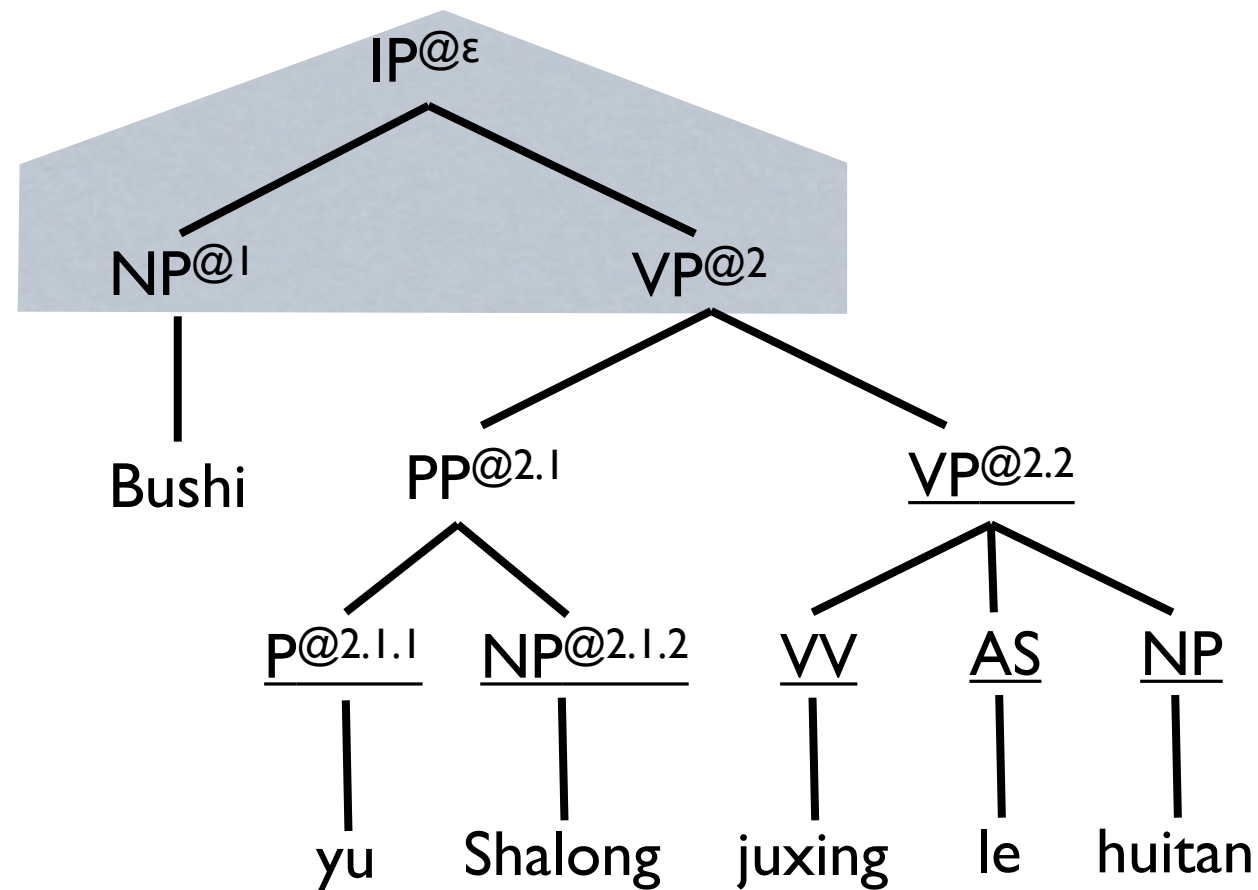
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2]
r₁

stack

<s>

hypothesis



rule probability

$P(r_1|\epsilon)$

action: predict (push)

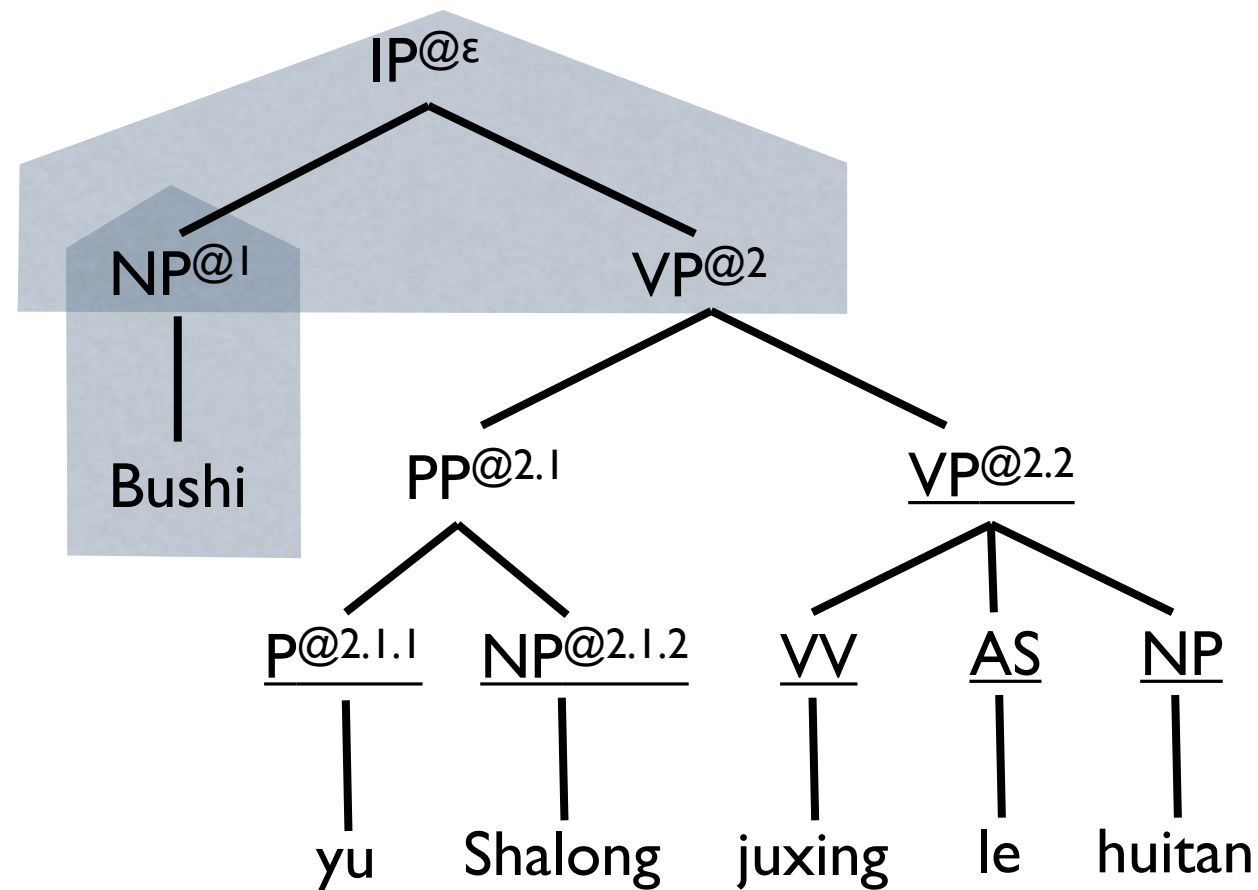
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2] [. Bush]

stack

<s>

hypothesis



rule probability
 $P(r_2|r_1)$

action: predict

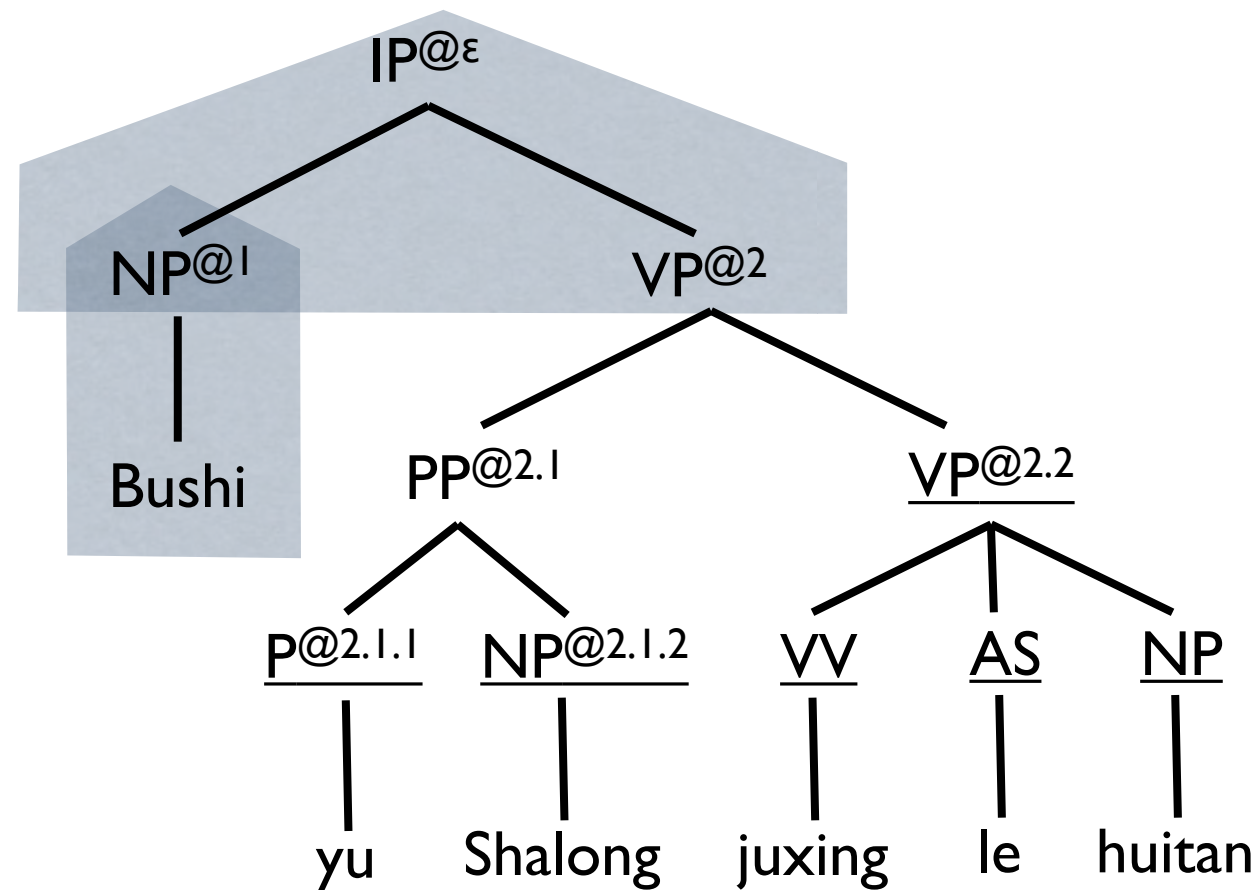
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [. NP@1 VP@2] [Bush .]
r₁ r₂

stack

<s> Bush

hypothesis



action: scan

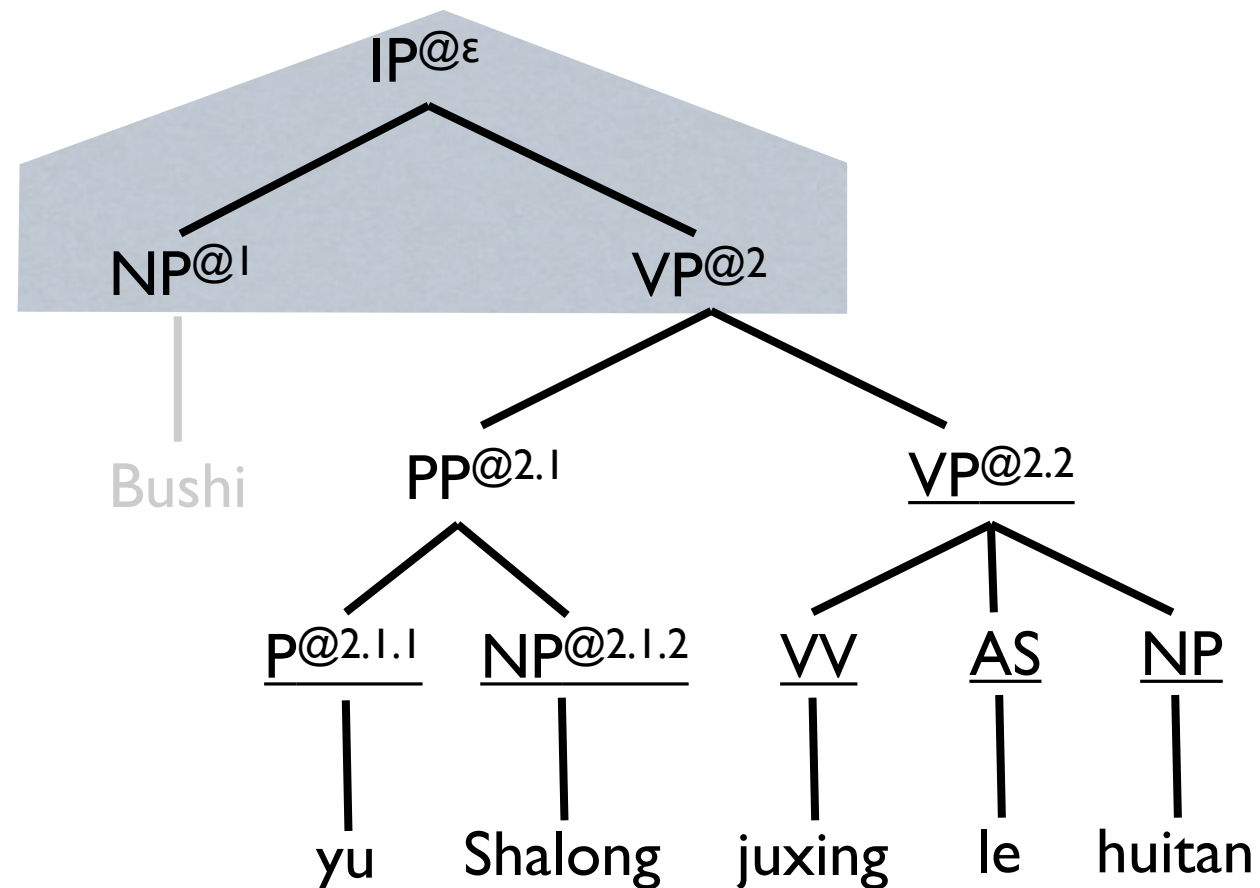
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2]
r1

stack

<s> Bush

hypothesis



action: pop

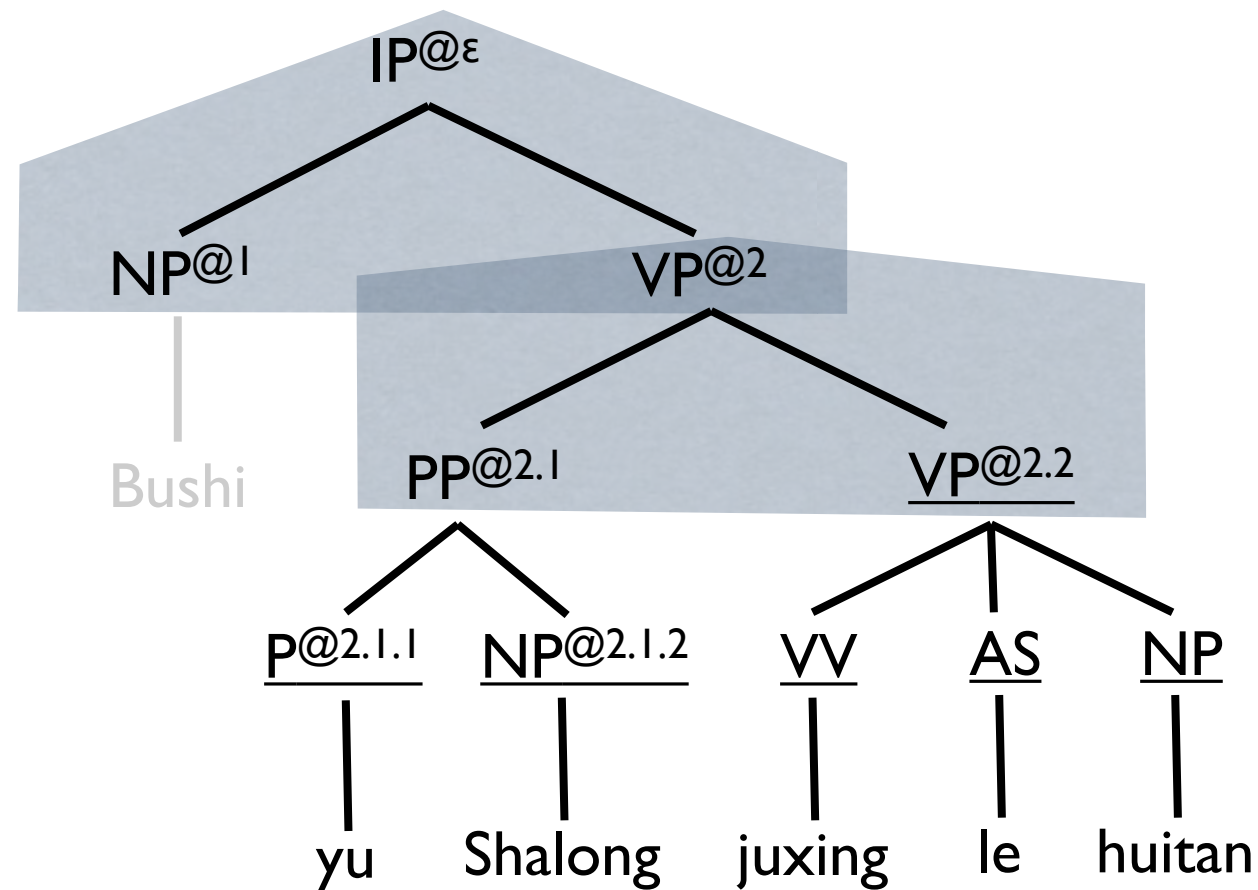
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1]
r₁ r₃

stack

<s> Bush

hypothesis



action: predict

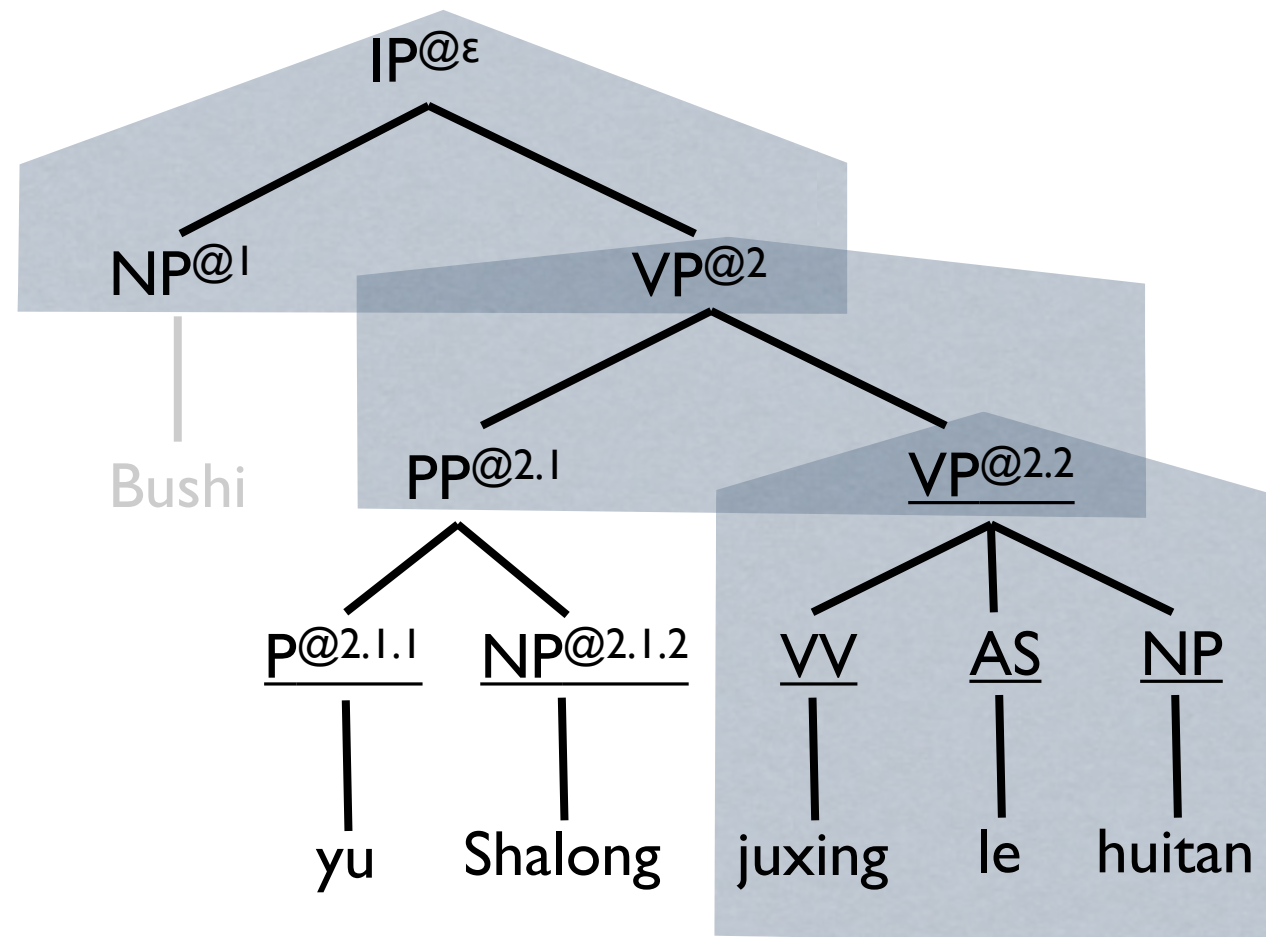
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [. VP@2.2 PP@2.1] [. held talks]
r1 r3 r5

stack

<s> Bush

hypothesis



rule probability

$$P(r_5|r_1, r_3)$$

action: predict

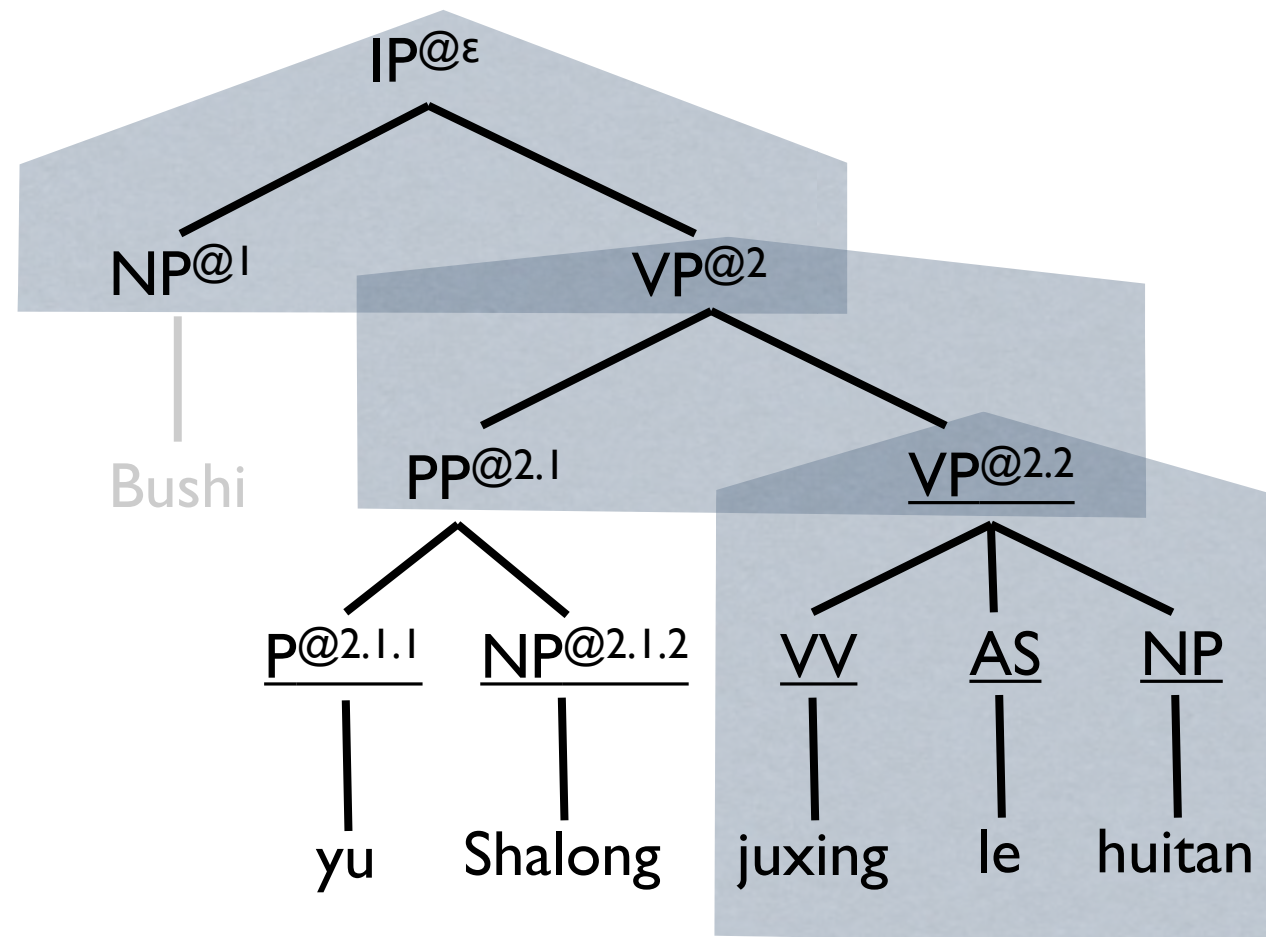
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1] [held talks .]
r1 r3 r5

stack

<s> Bush held talks

hypothesis



action: scan

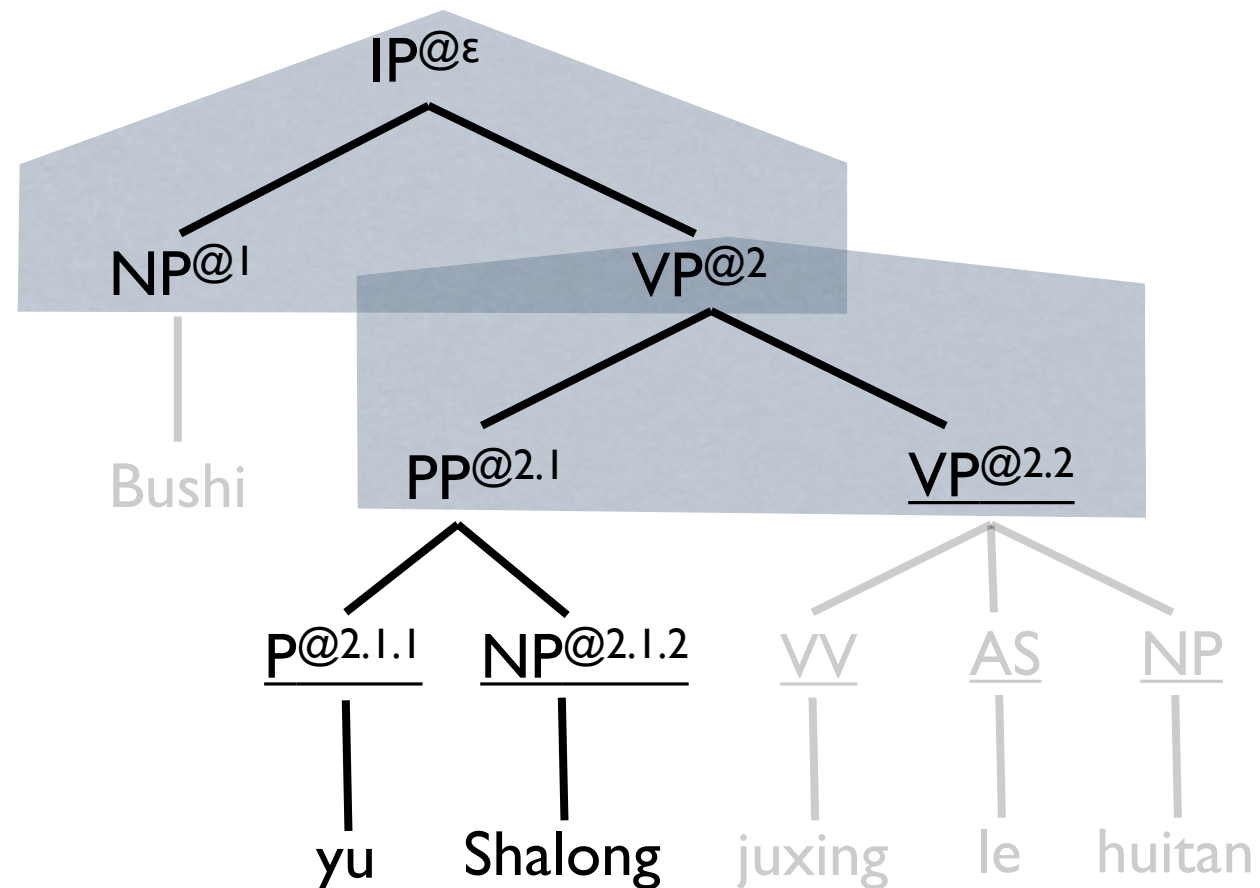
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1]
r1 r3

stack

<s> Bush held talks

hypothesis



action: pop

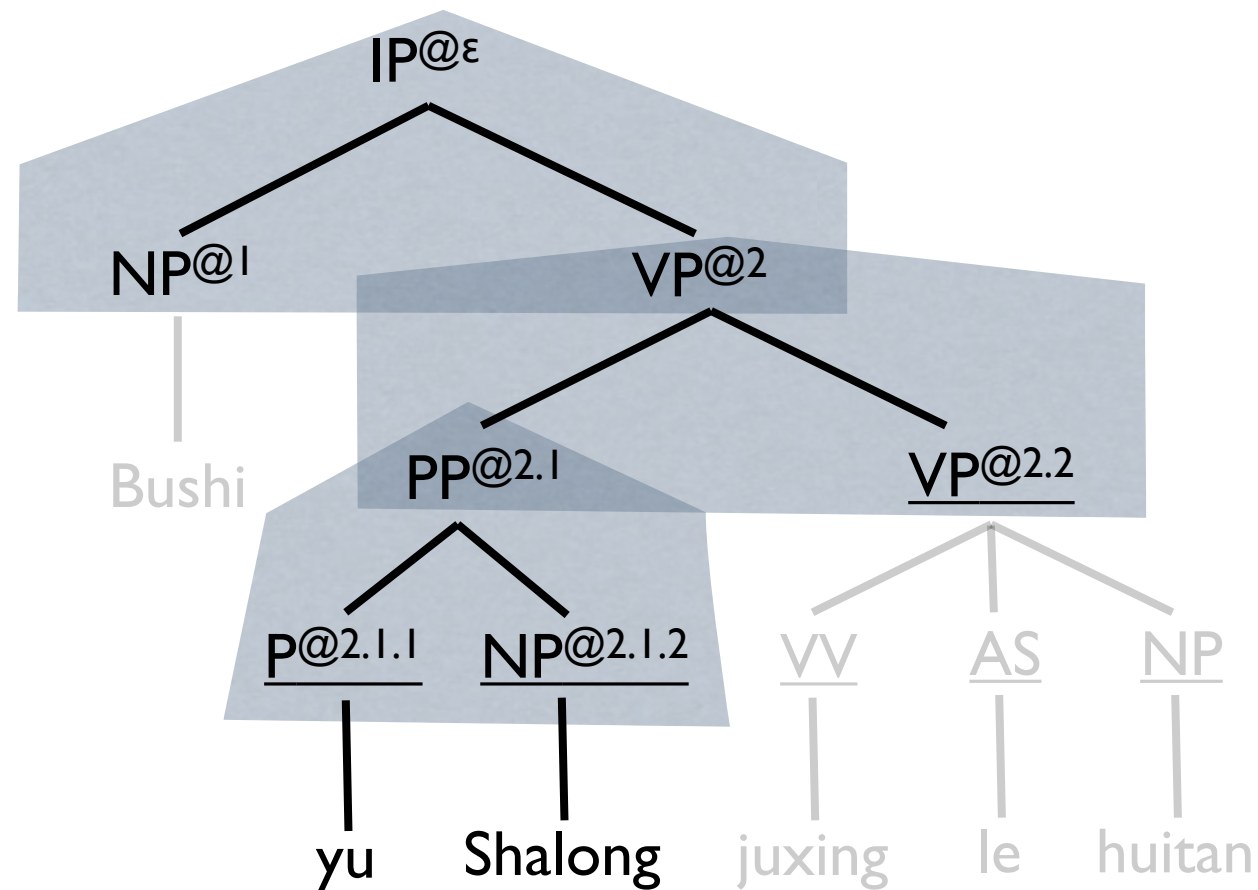
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 NP@2.1.2]

stack

<s> Bush held talks

hypothesis



rule probability

$$P(r_4|r_1, r_3)$$

action: predict

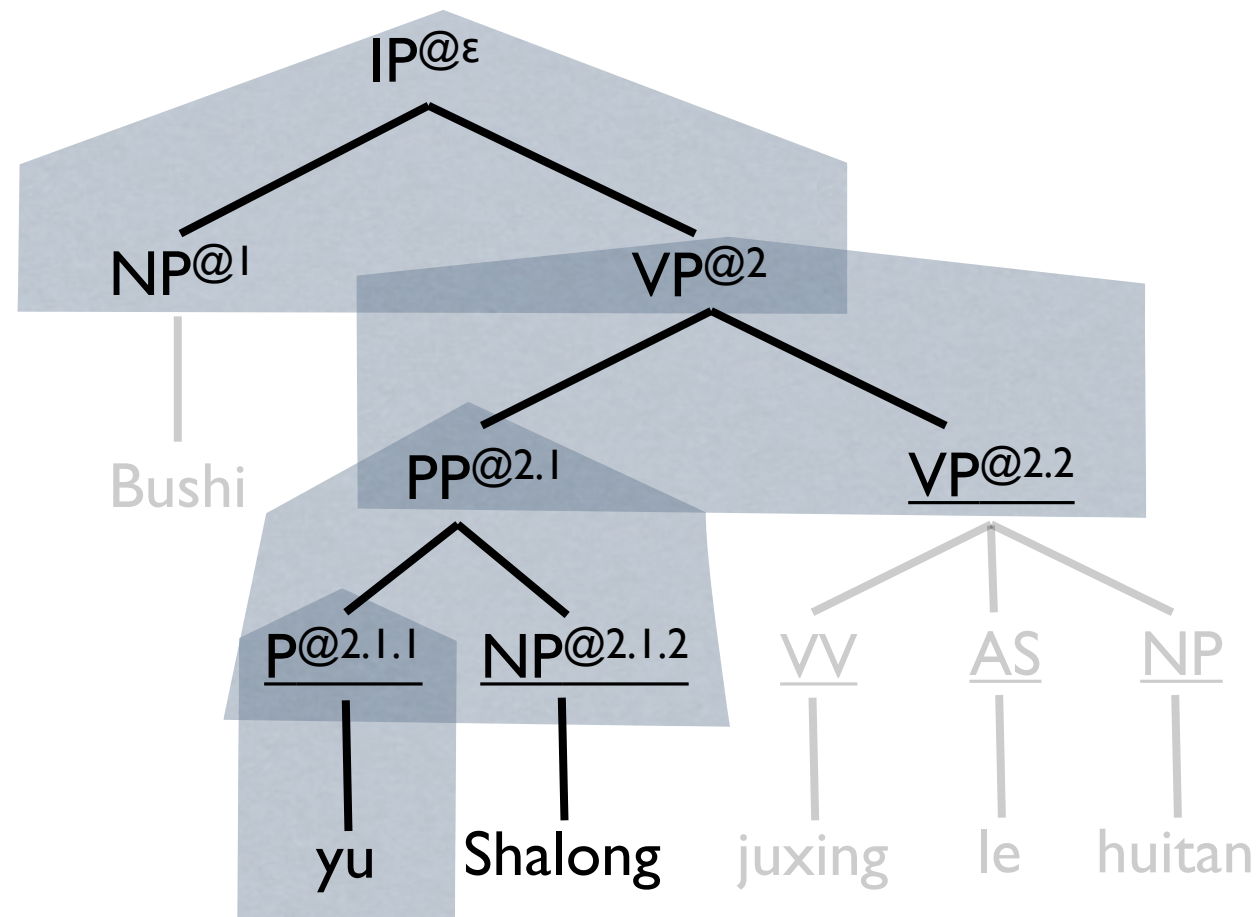
Example Incremental Decoding

$[\langle s \rangle \cdot \text{IP}@^\epsilon \langle /s \rangle]$ $[\text{NP}@^1 \cdot \text{VP}@^2]$ $[\text{VP}@^{2.2} \cdot \text{PP}@^{2.1}]$ $[\cdot \text{P}@^{2.1.1} \text{NP}@^{2.1.2}]$ $[\cdot \text{with/and}]$ *stack*

r_1 r_3 r_4 r_6/r'_6

$\langle s \rangle$ Bush held talks

hypothesis



rule probability
 $P(r_6/r'_6 | r_3, r_4)$

action: predict

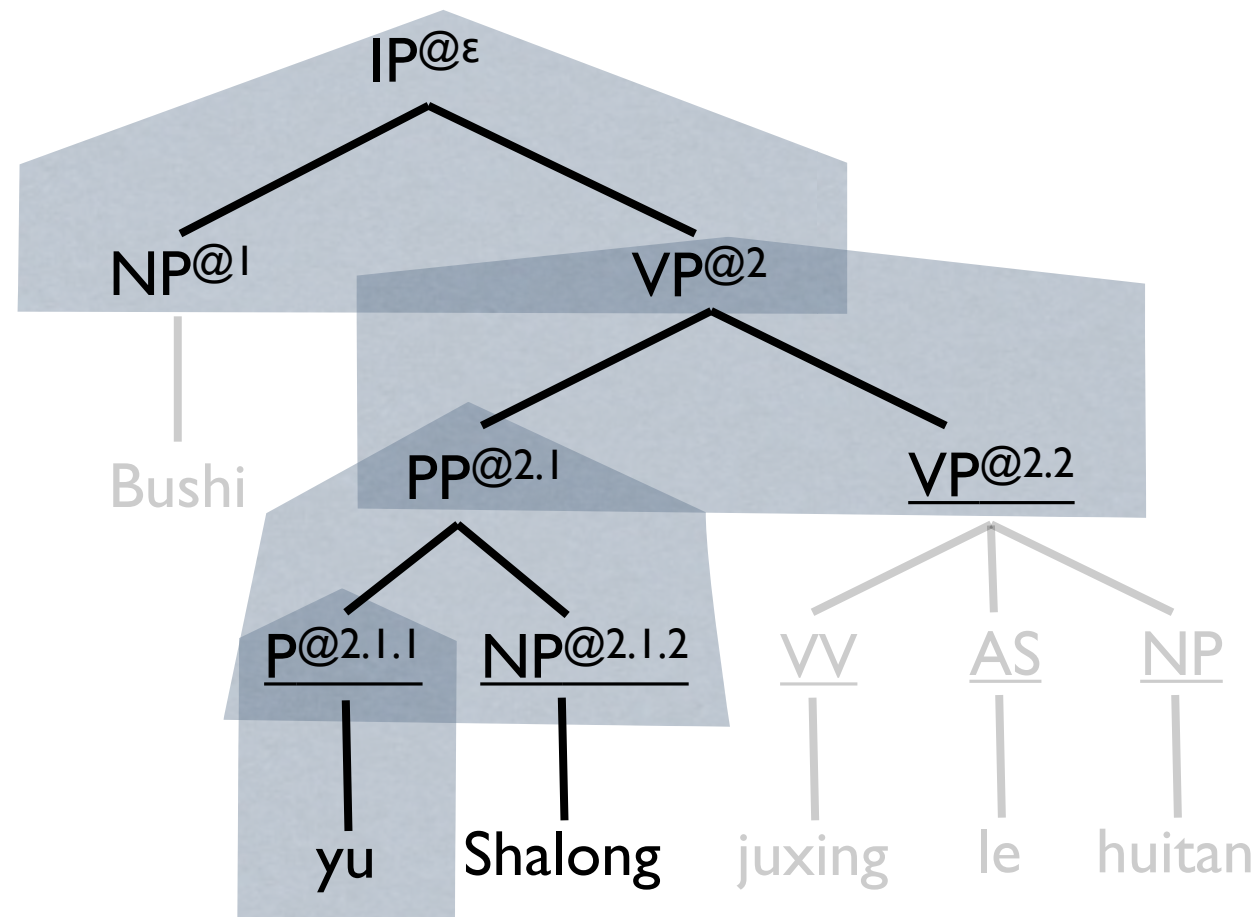
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [. P@2.1.1 NP@2.1.2] [with .]

stack

<s> Bush held talks with

hypothesis



action: scan

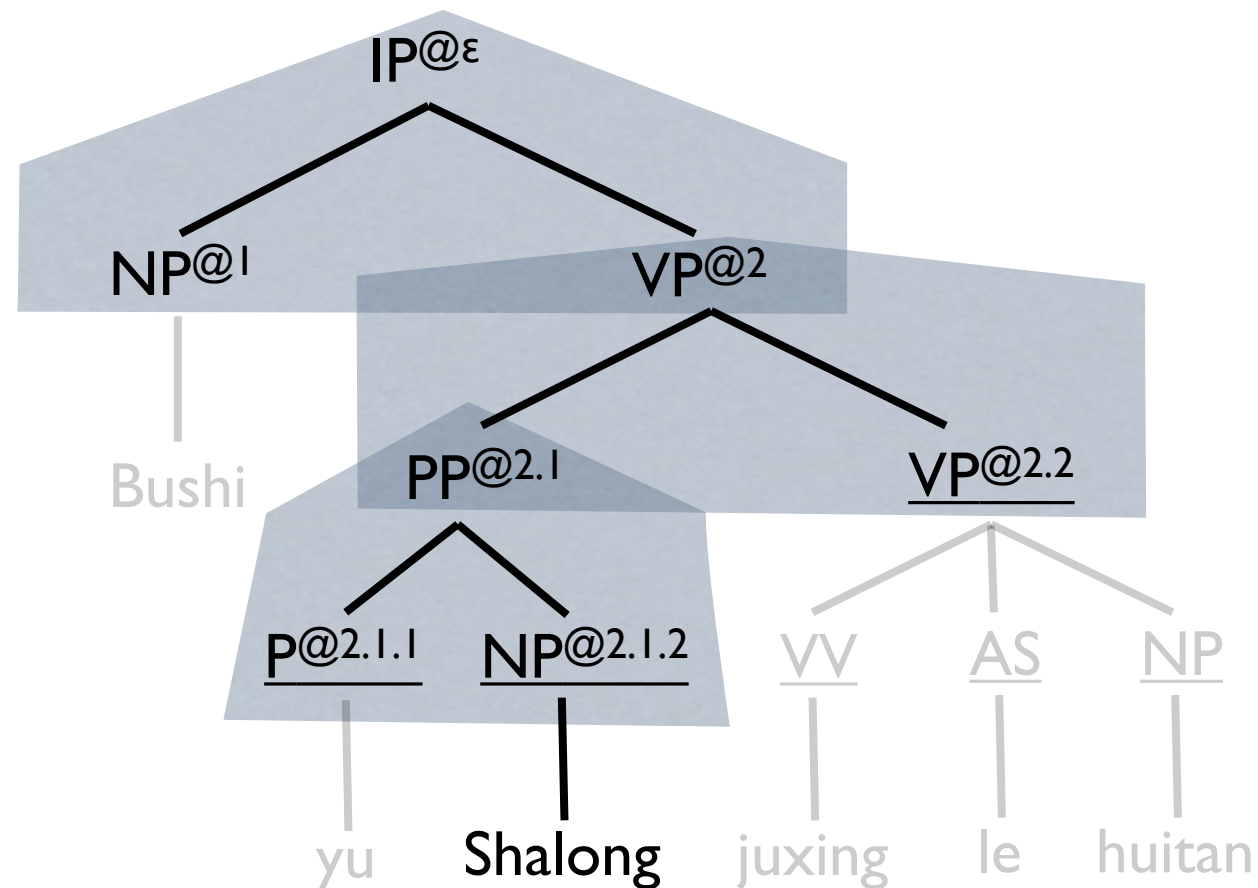
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2]

stack

<s> Bush held talks with

hypothesis



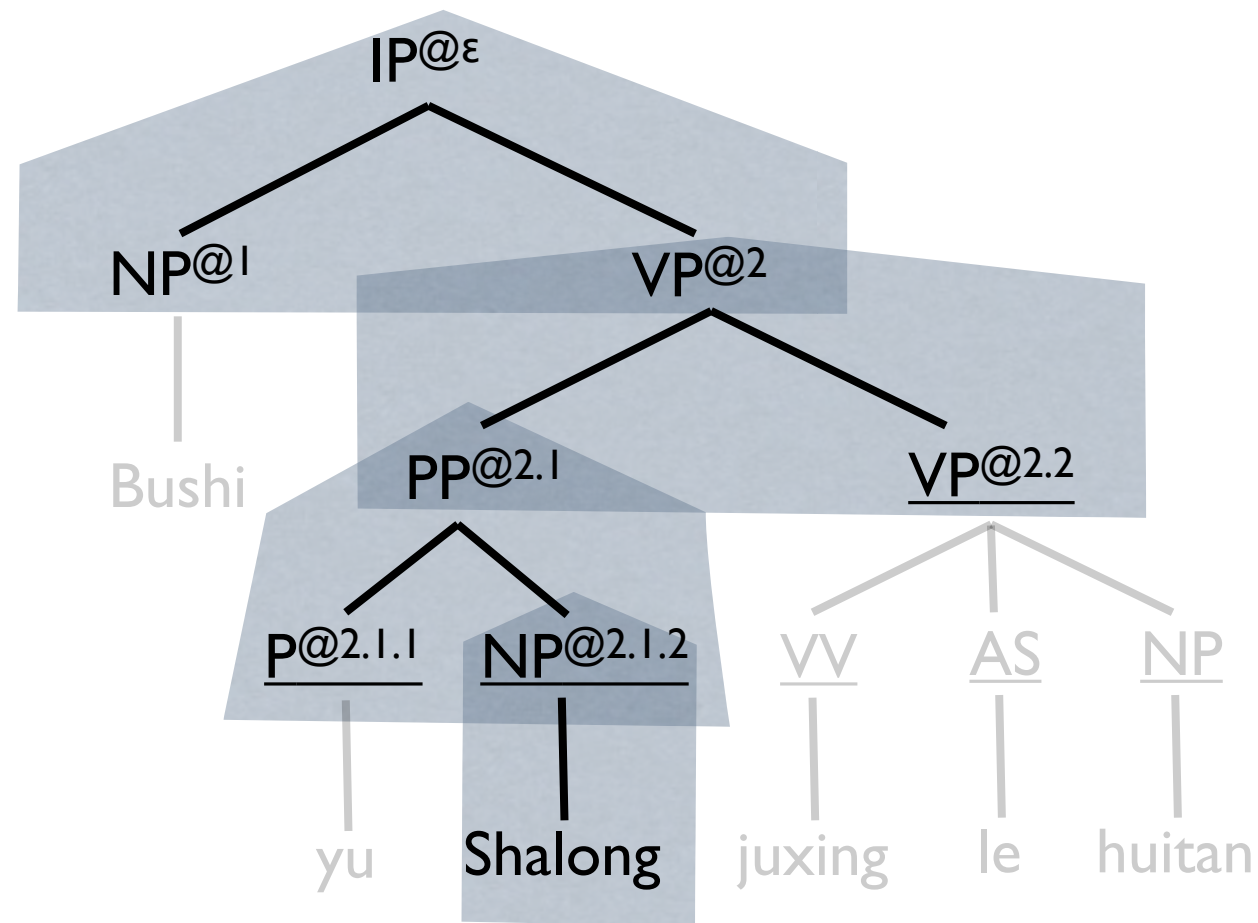
action: pop

Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2] [. Sharon] *stack*

r_1 r_3 r_4 r_7

<s> Bush held talks with *hypothesis*



rule probability
 $P(r_7 | r_3, r_4)$

action: predict

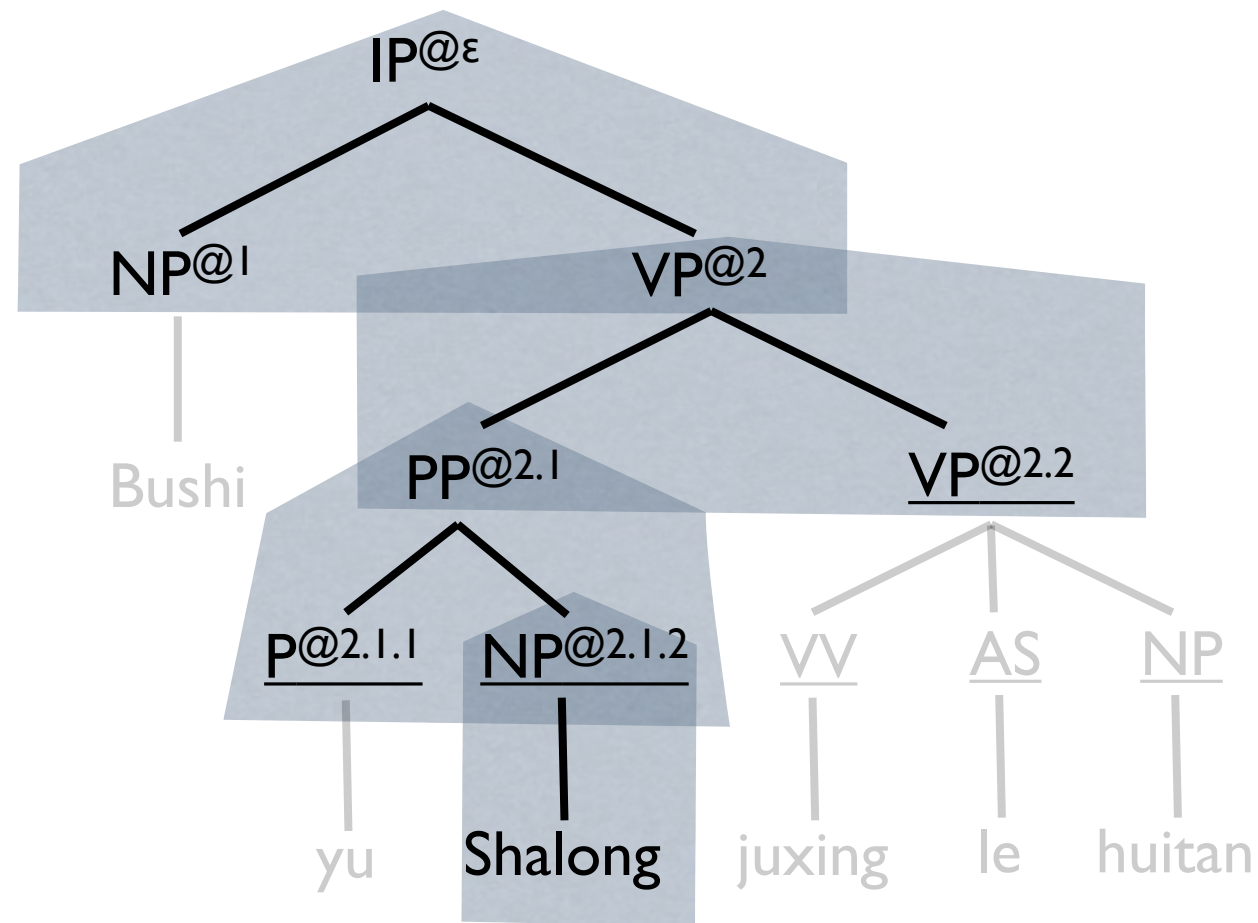
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 . NP@2.1.2] [Sharon .] *stack*

r_1 r_3 r_4 r_7

<s> Bush held talks with Sharon

hypothesis



action: scan

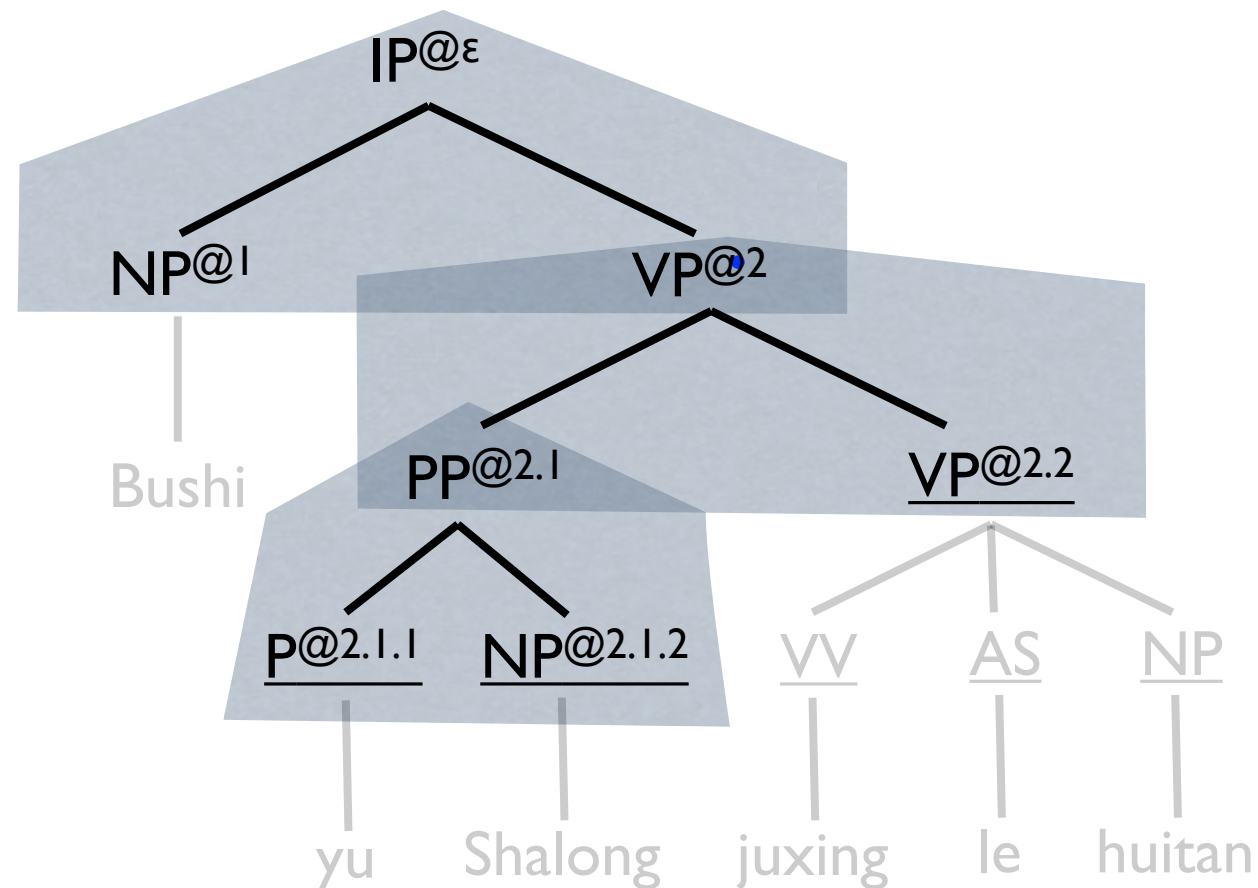
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 . PP@2.1] [P@2.1.1 NP@2.1.2 .]

stack

<s> Bush held talks with Sharon

hypothesis



action: pop

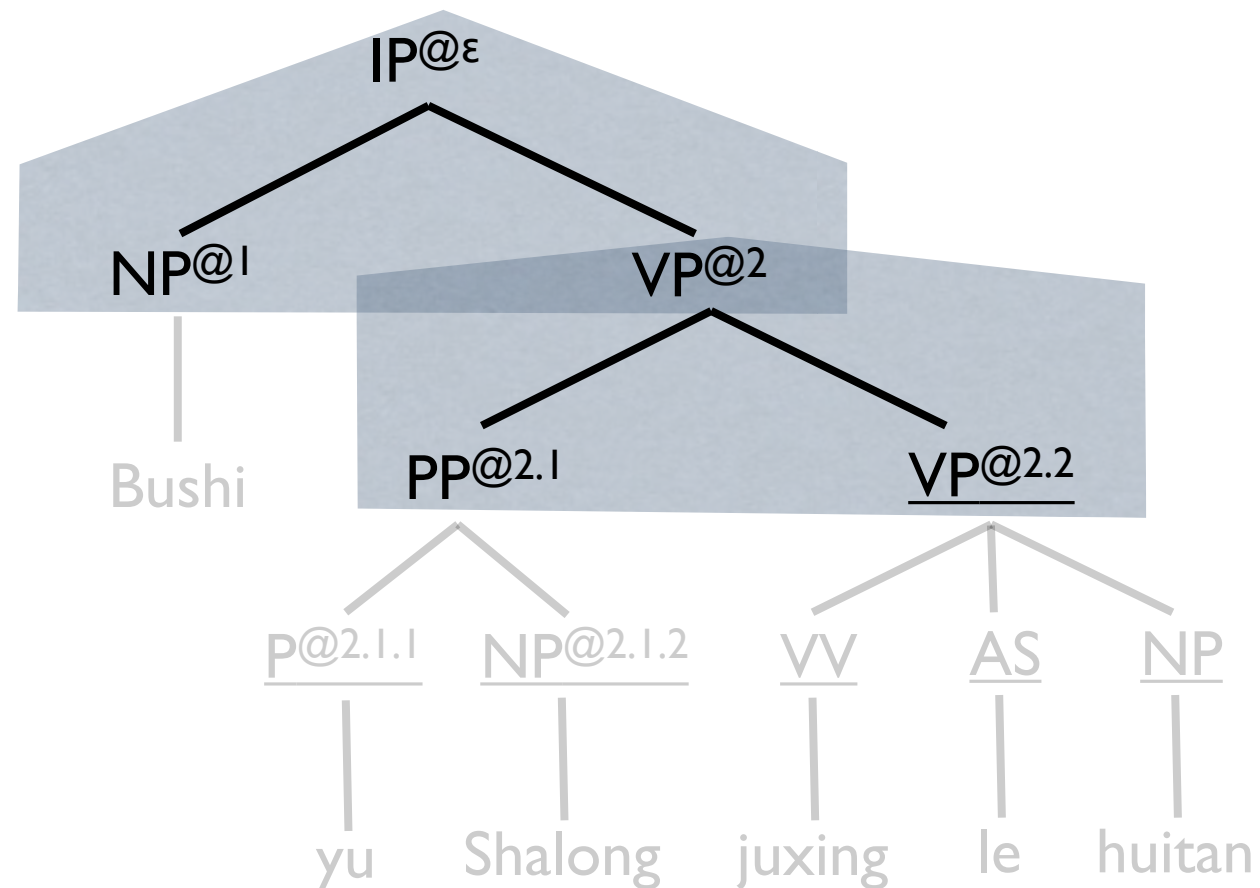
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 . VP@2] [VP@2.2 PP@2.1 .]
r1 r3

stack

<s> Bush held talks with Sharon

hypothesis



action: pop

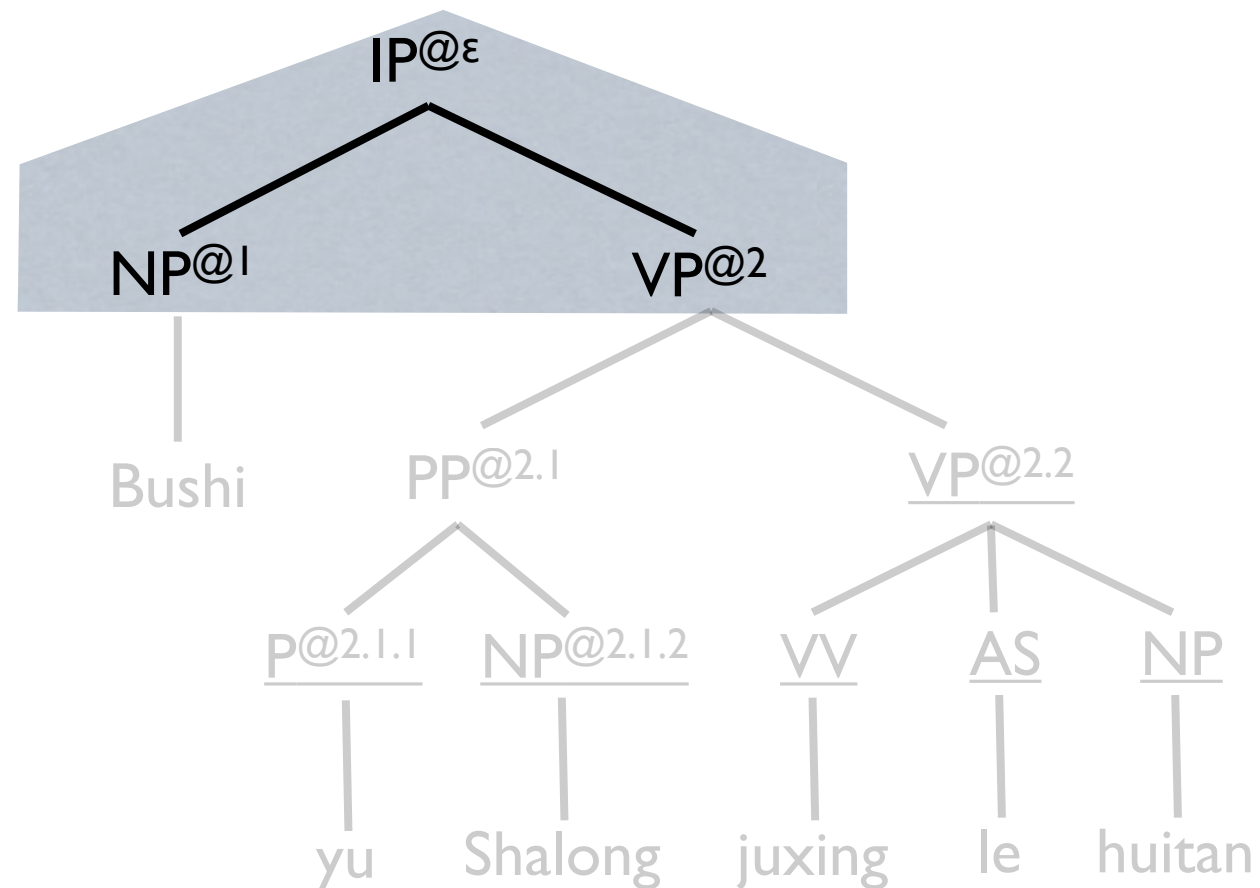
Example Incremental Decoding

[<s> . IP@ ϵ </s>] [NP@1 VP@2 .]
r1

stack

<s> Bush held talks with Sharon

hypothesis



action: pop

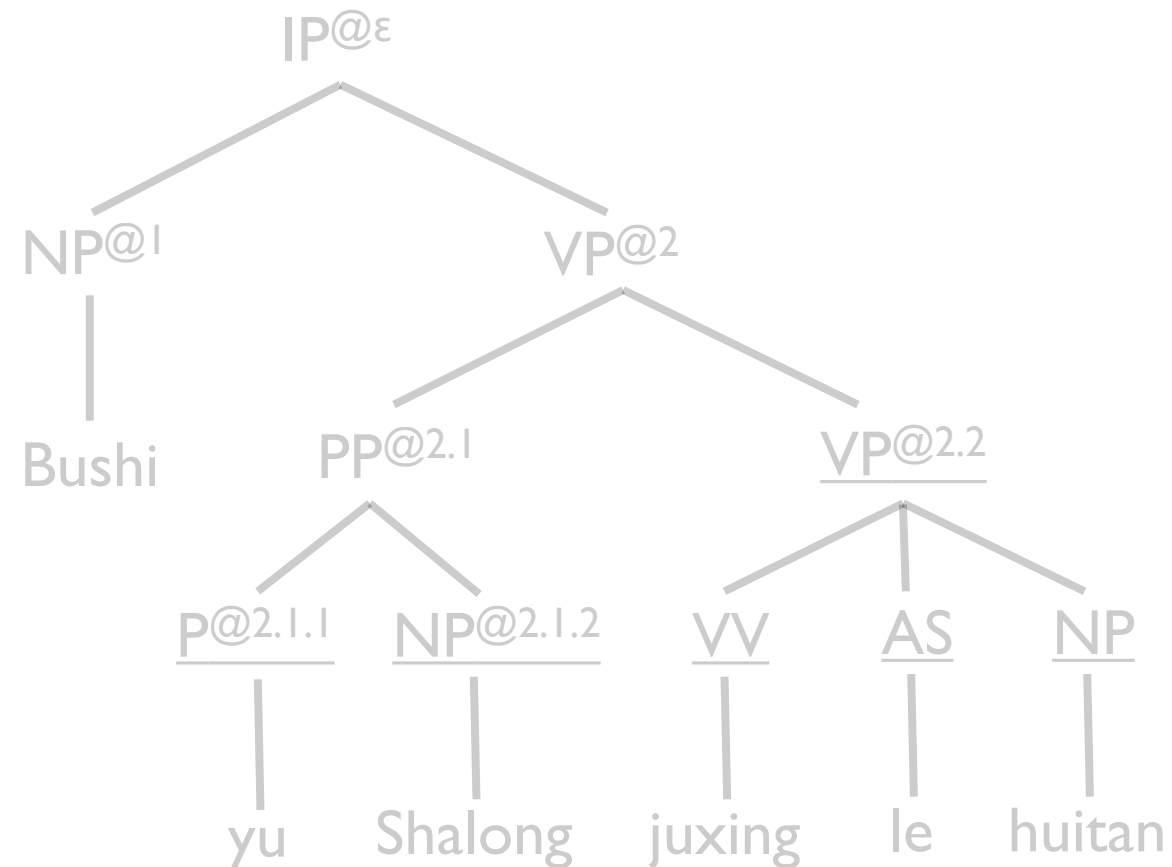
Example Incremental Decoding

[<s> IP@ ϵ .</s>]

stack

<s> Bush held talks with Sharon

hypothesis



action: pop

Experiments

- 1.5 million sentence pairs with 38/32 million words of Chinese/English
- Dev set : 616 sentences of the Newswire portion of 2006 NIST MT evaluation test set
- Test set: 619 sentences of the Newswire portion of 2006 NIST MT evaluation test set

Results

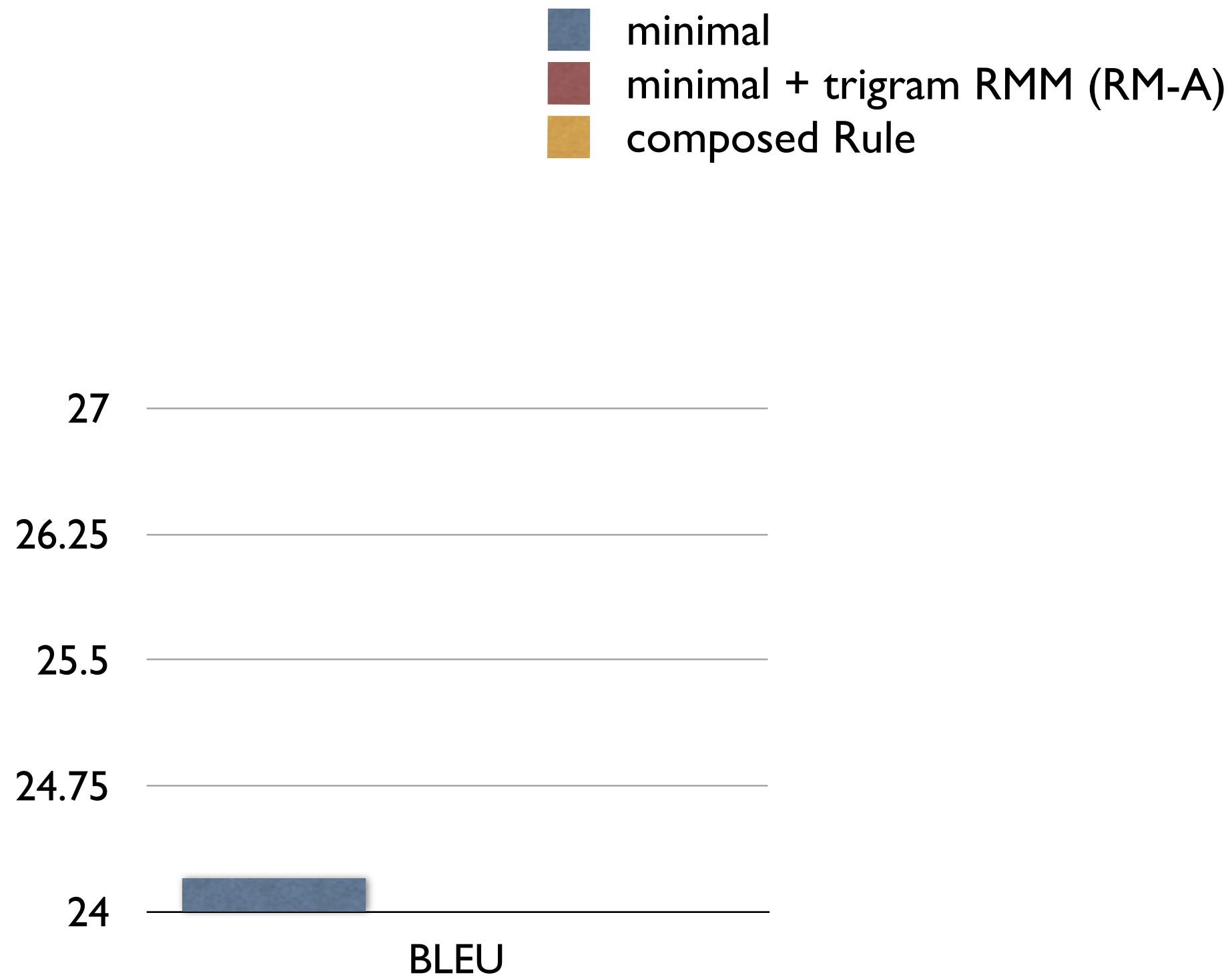
- minimal
- minimal + trigram RMM (RM-A)
- composed Rule

Results

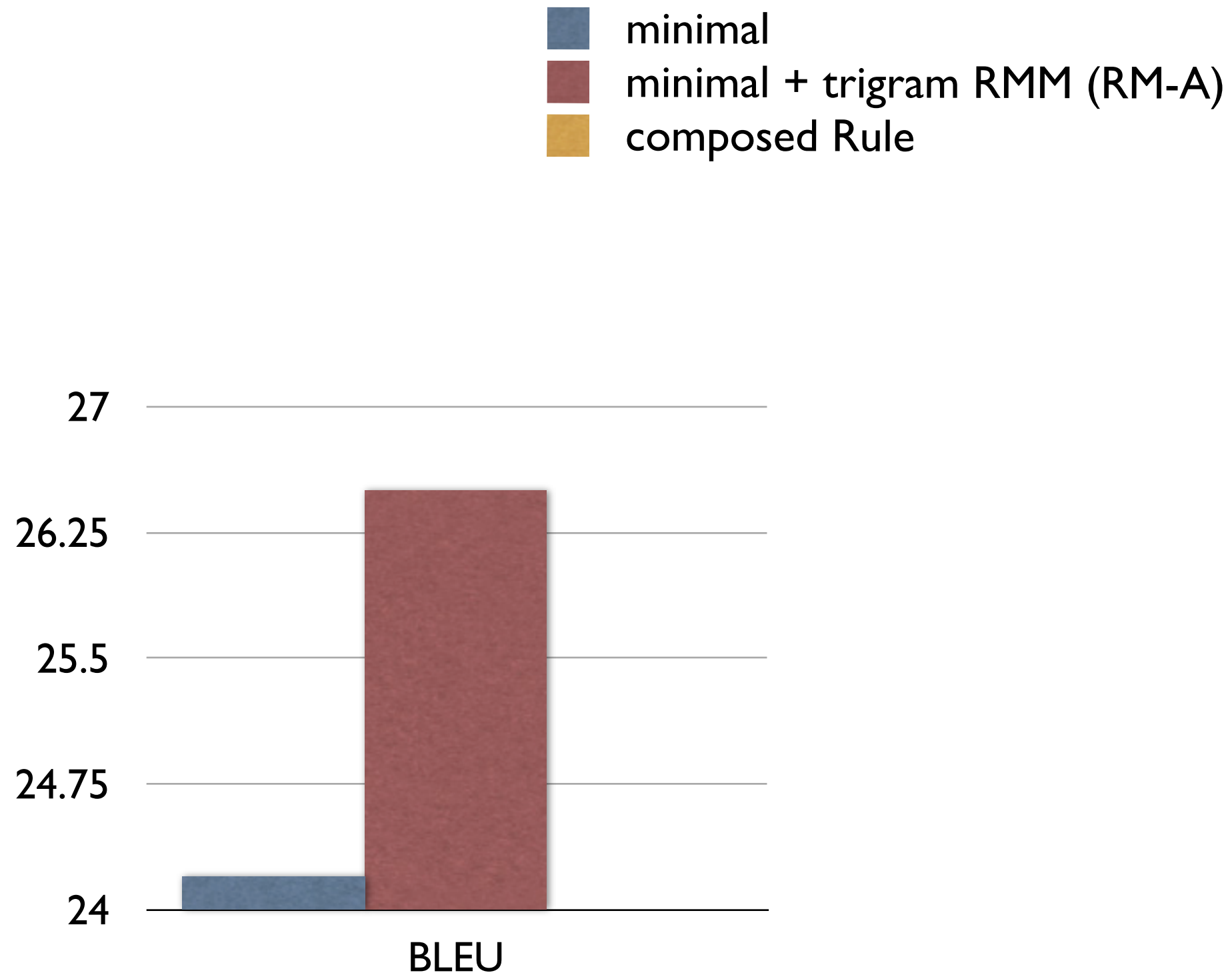
- minimal
- minimal + trigram RMM (RM-A)
- composed Rule



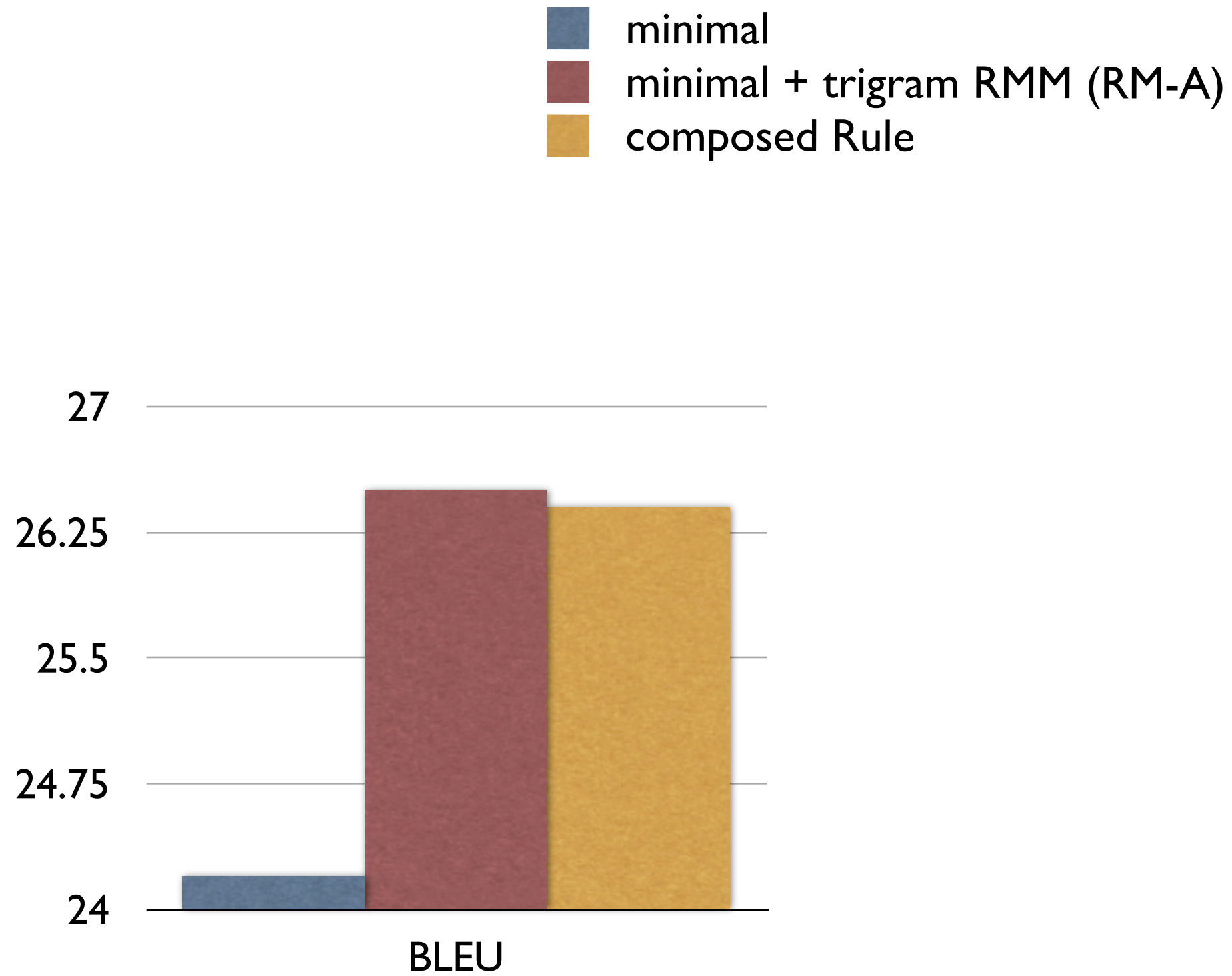
Results



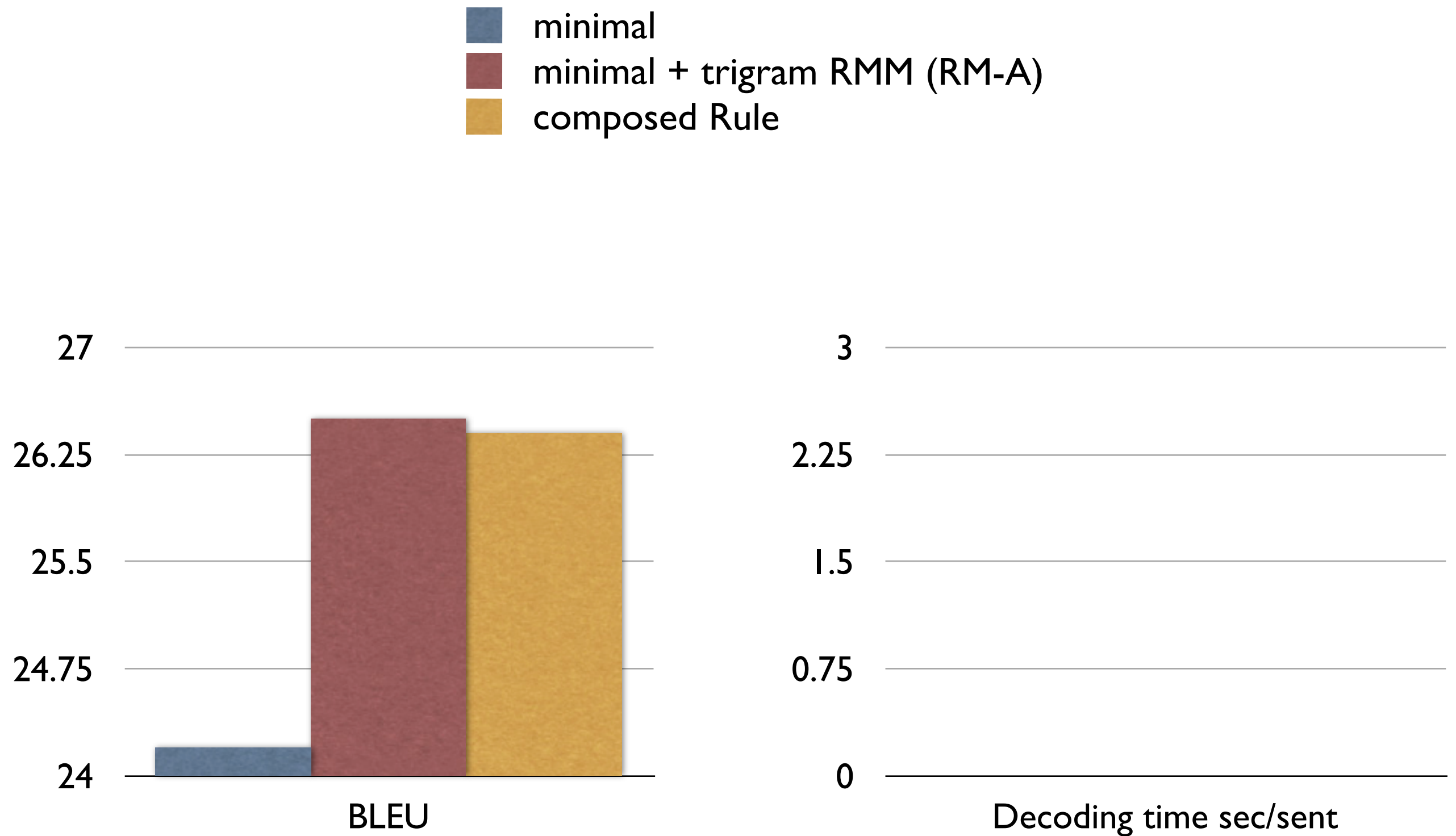
Results



Results

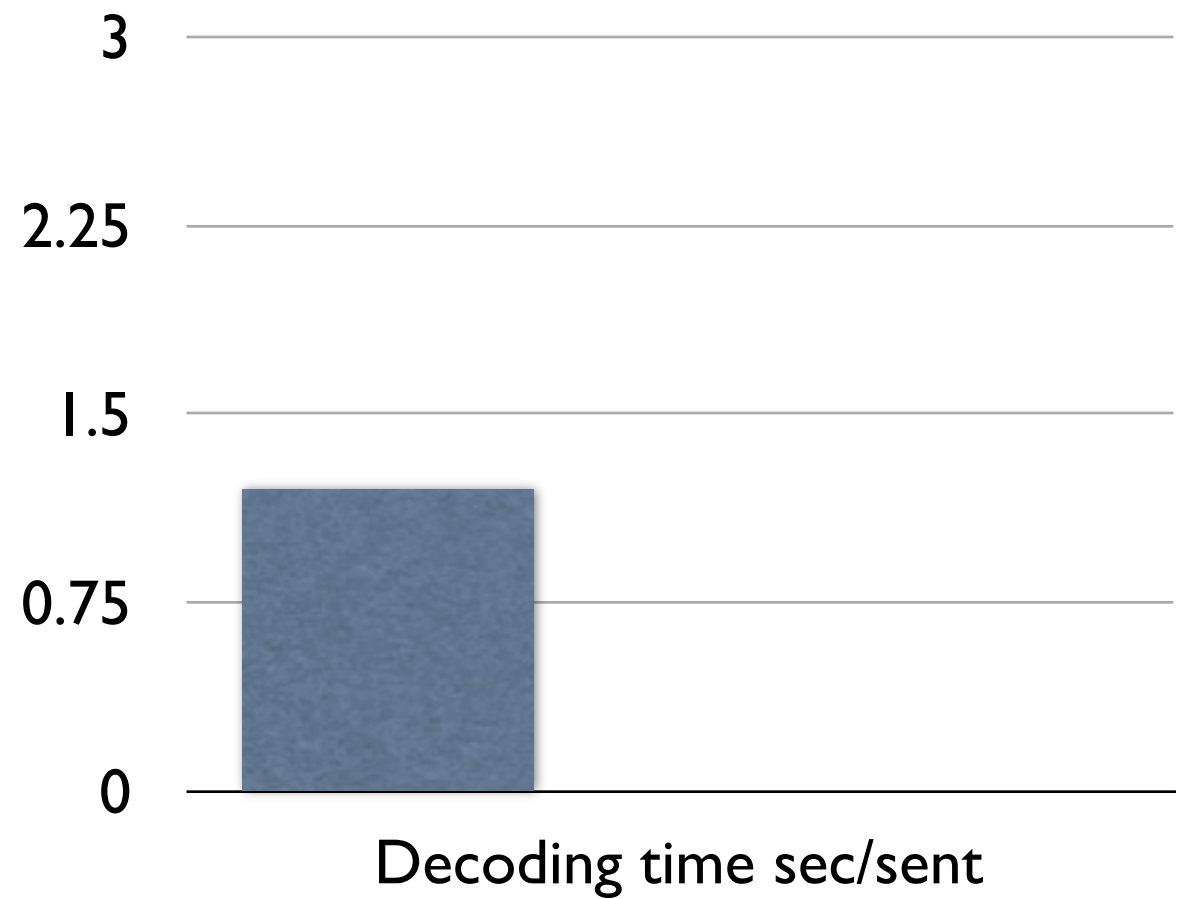
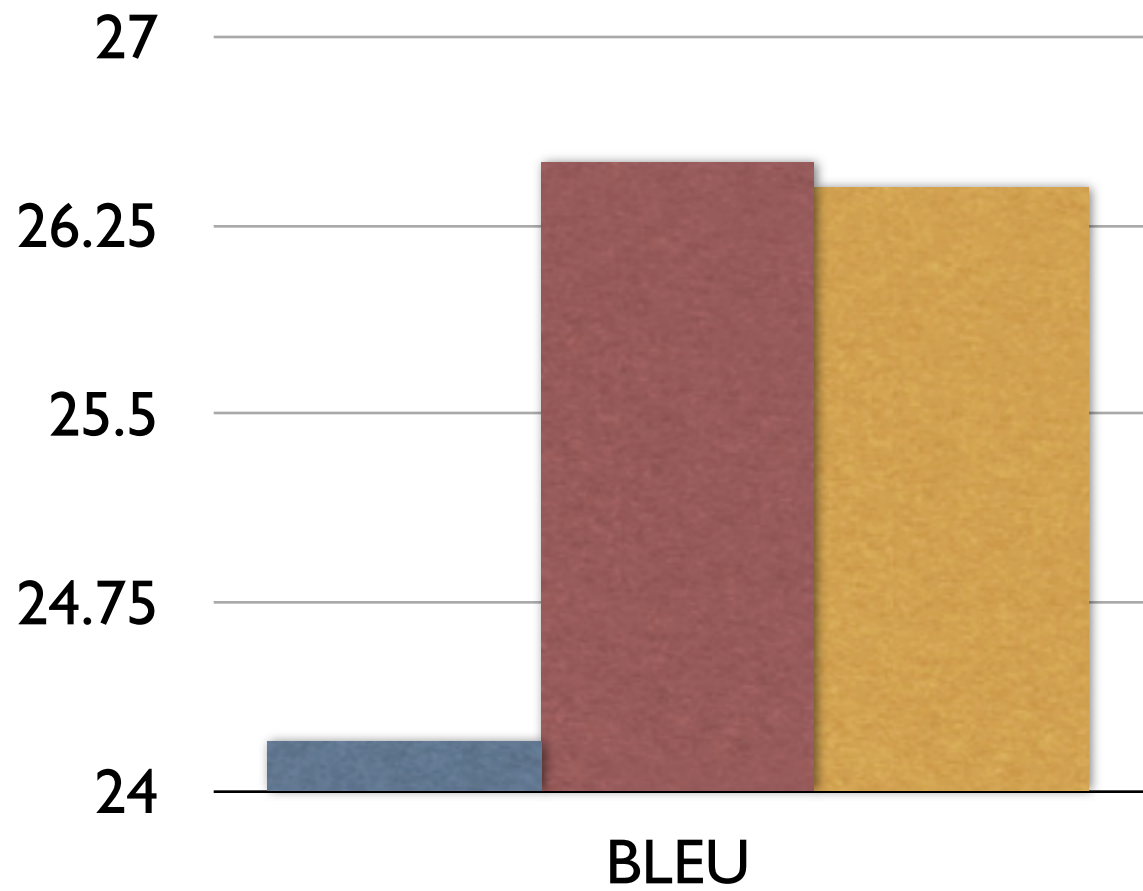


Results



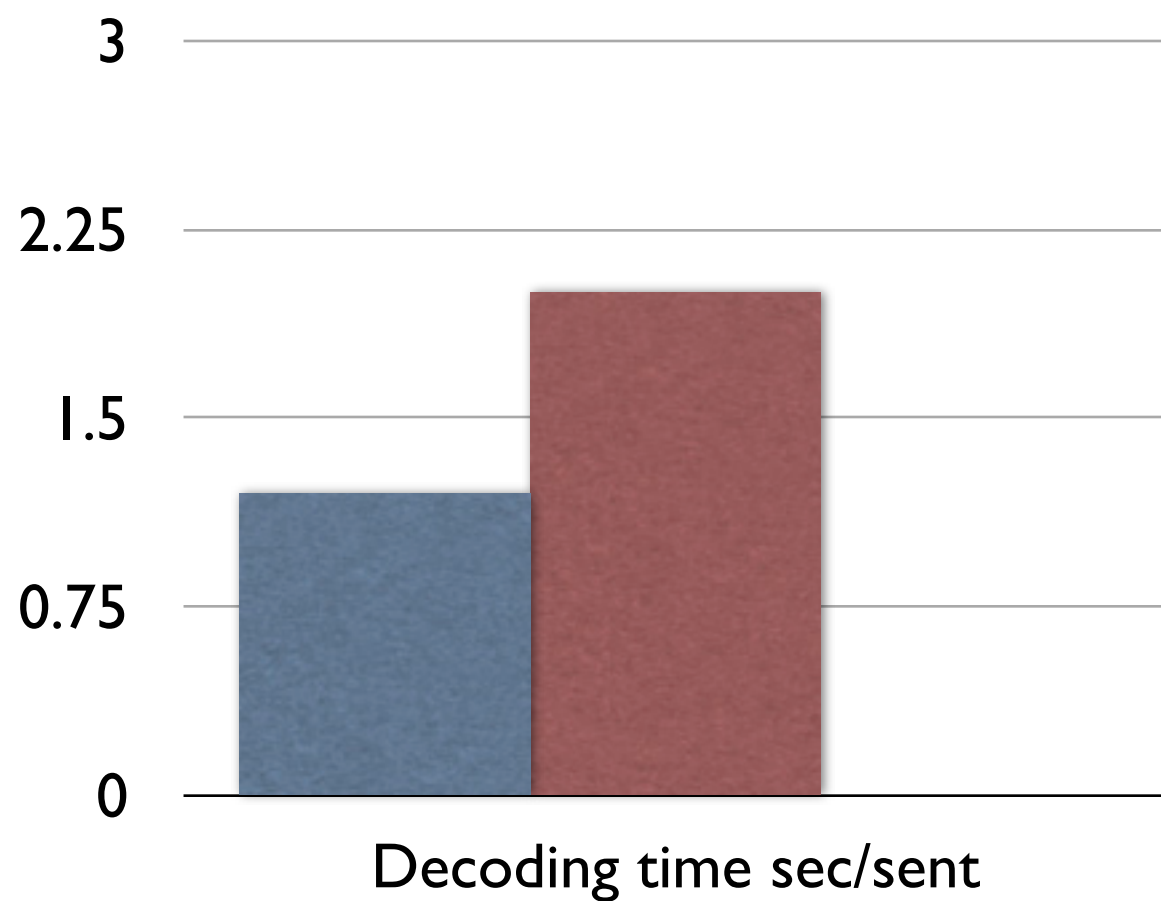
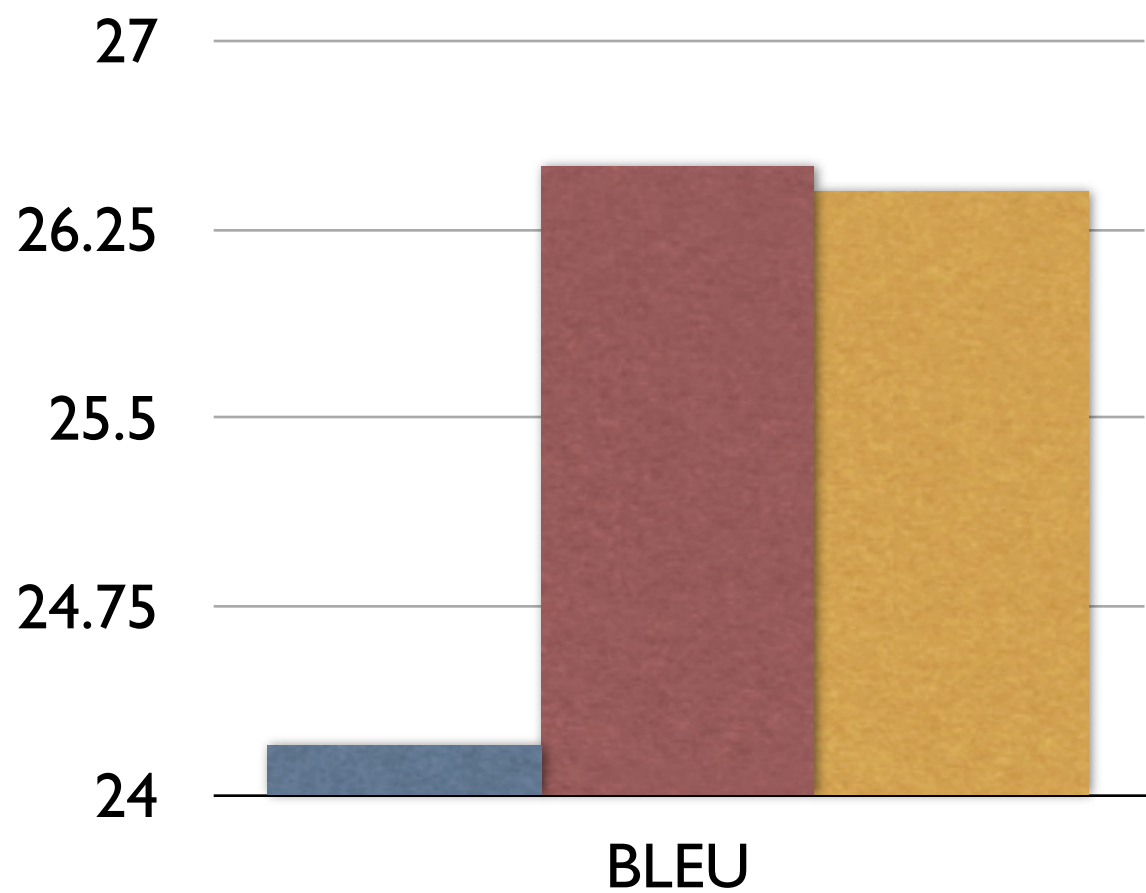
Results

- minimal
- minimal + trigram RMM (RM-A)
- composed Rule



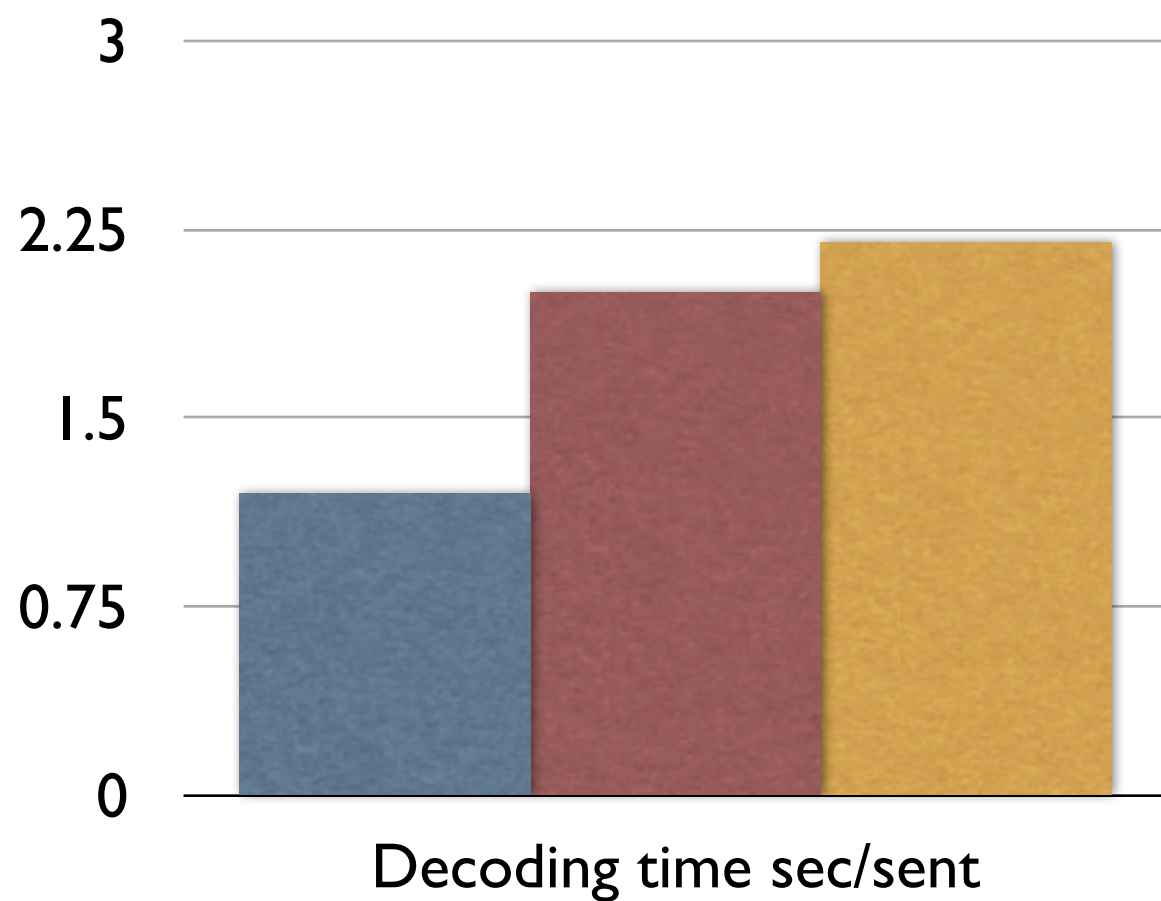
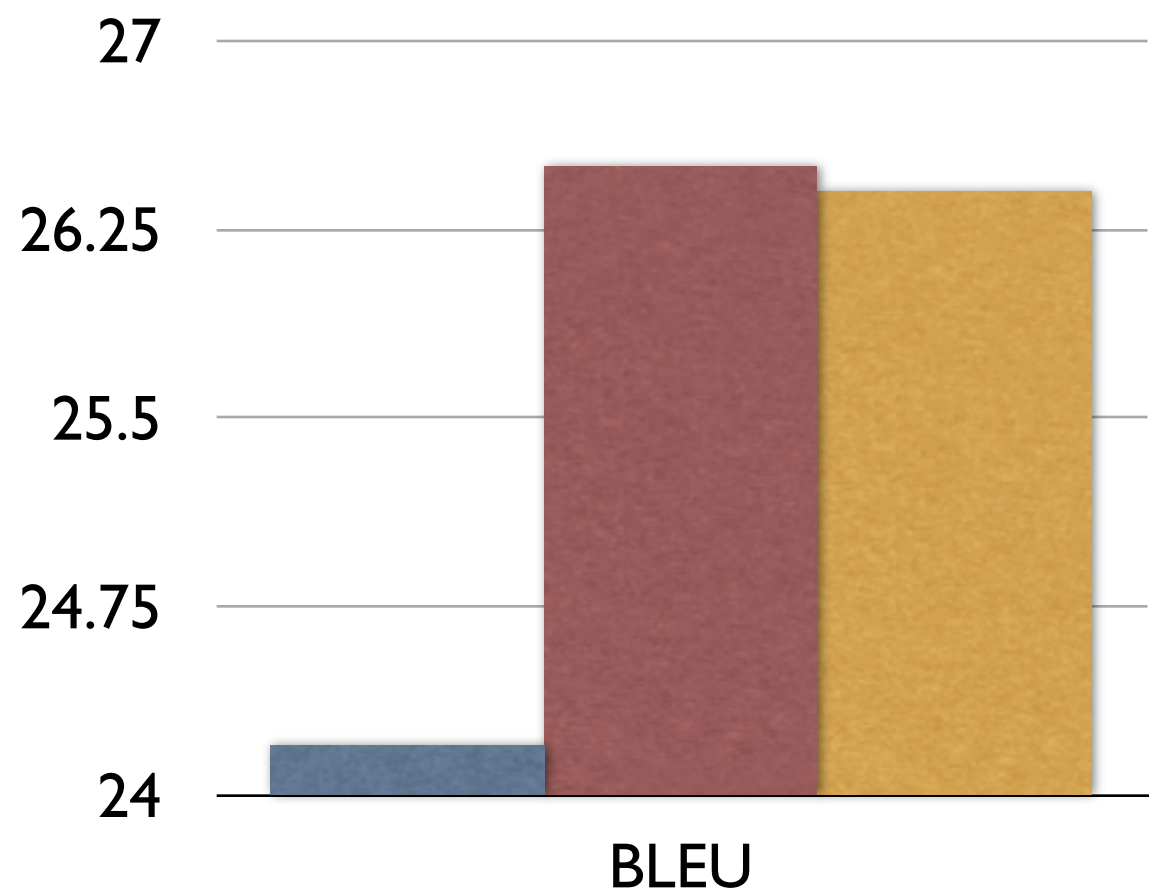
Results

- minimal
- minimal + trigram RMM (RM-A)
- composed Rule



Results

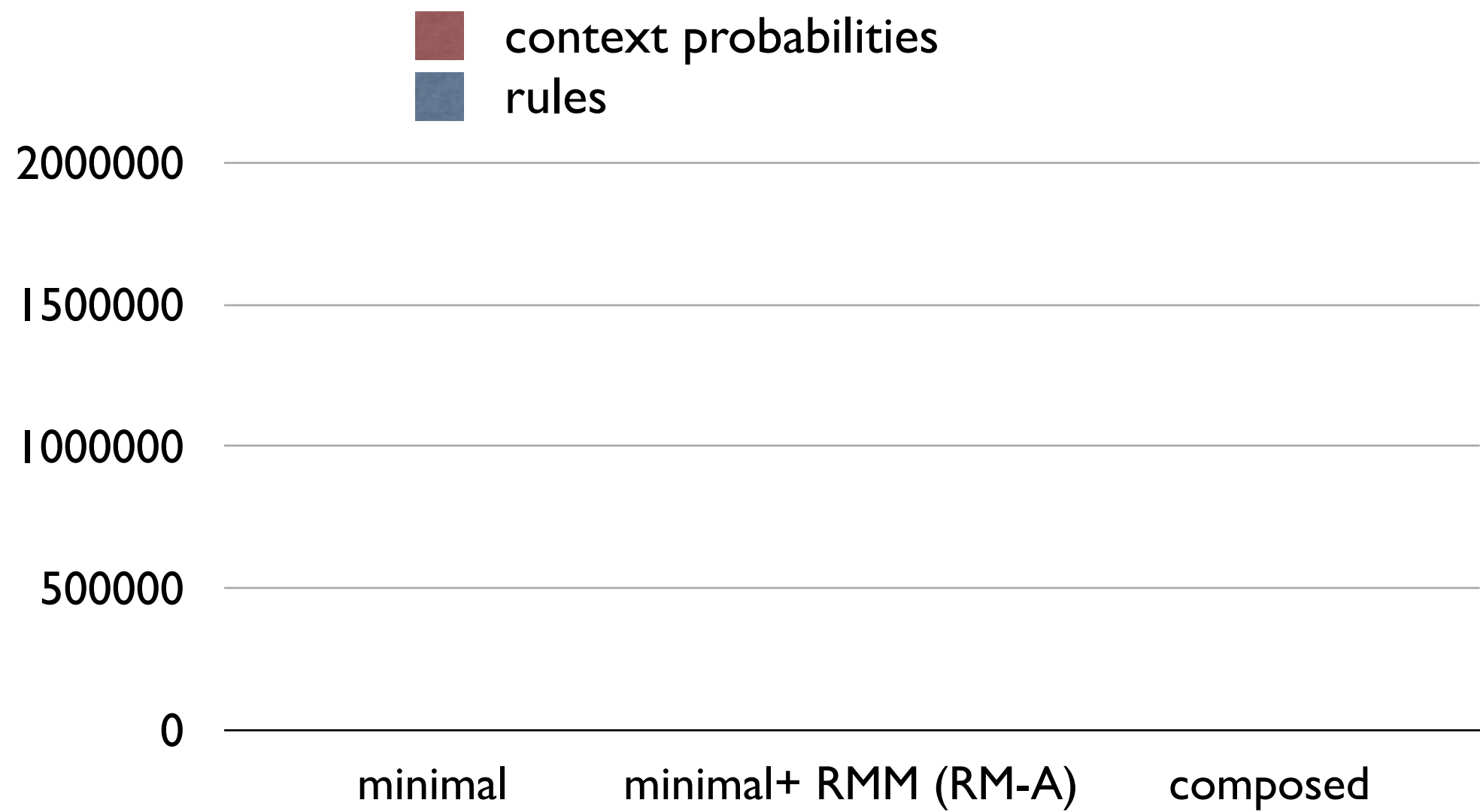
- minimal
- minimal + trigram RMM (RM-A)
- composed Rule



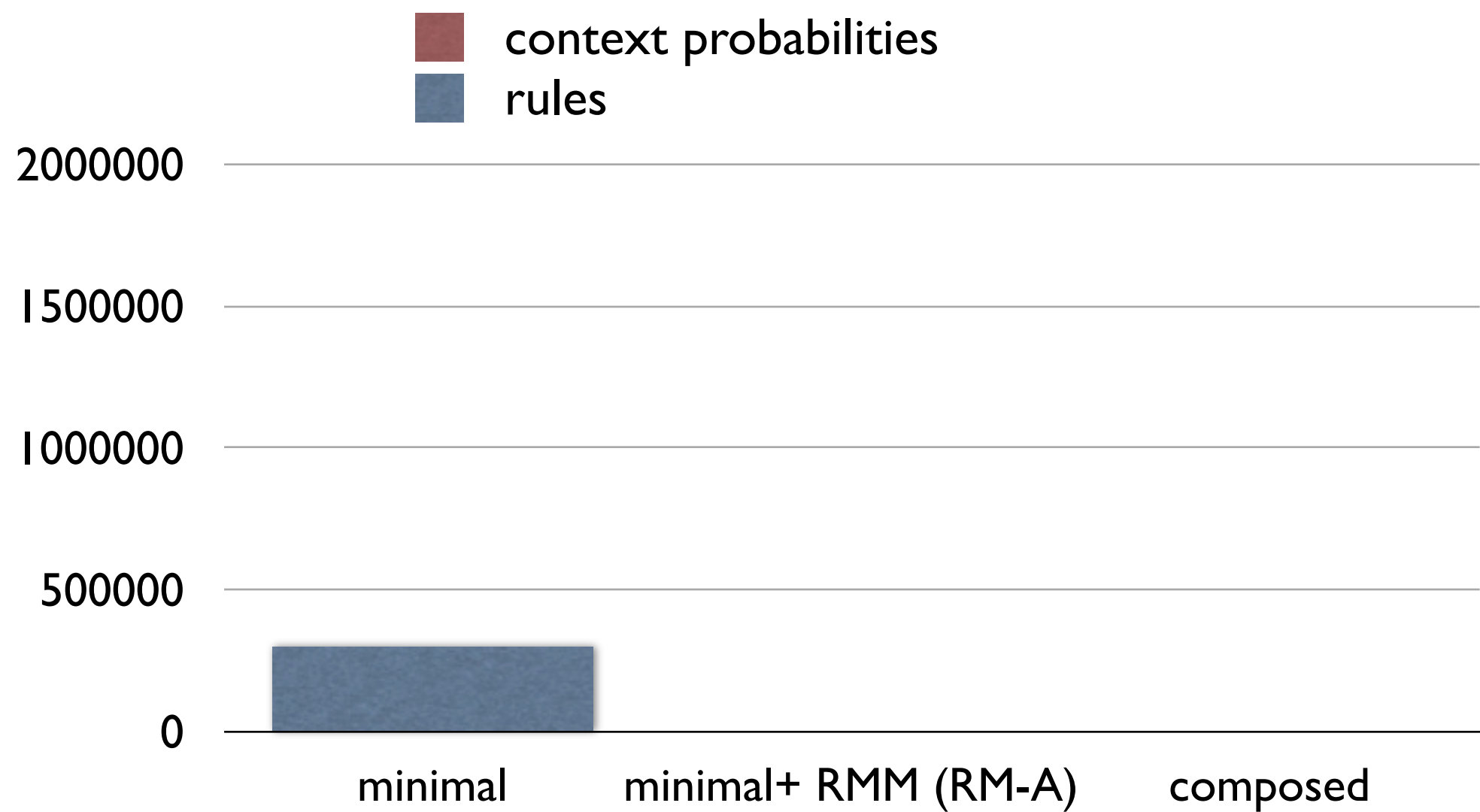
Results

- context probabilities
- rules

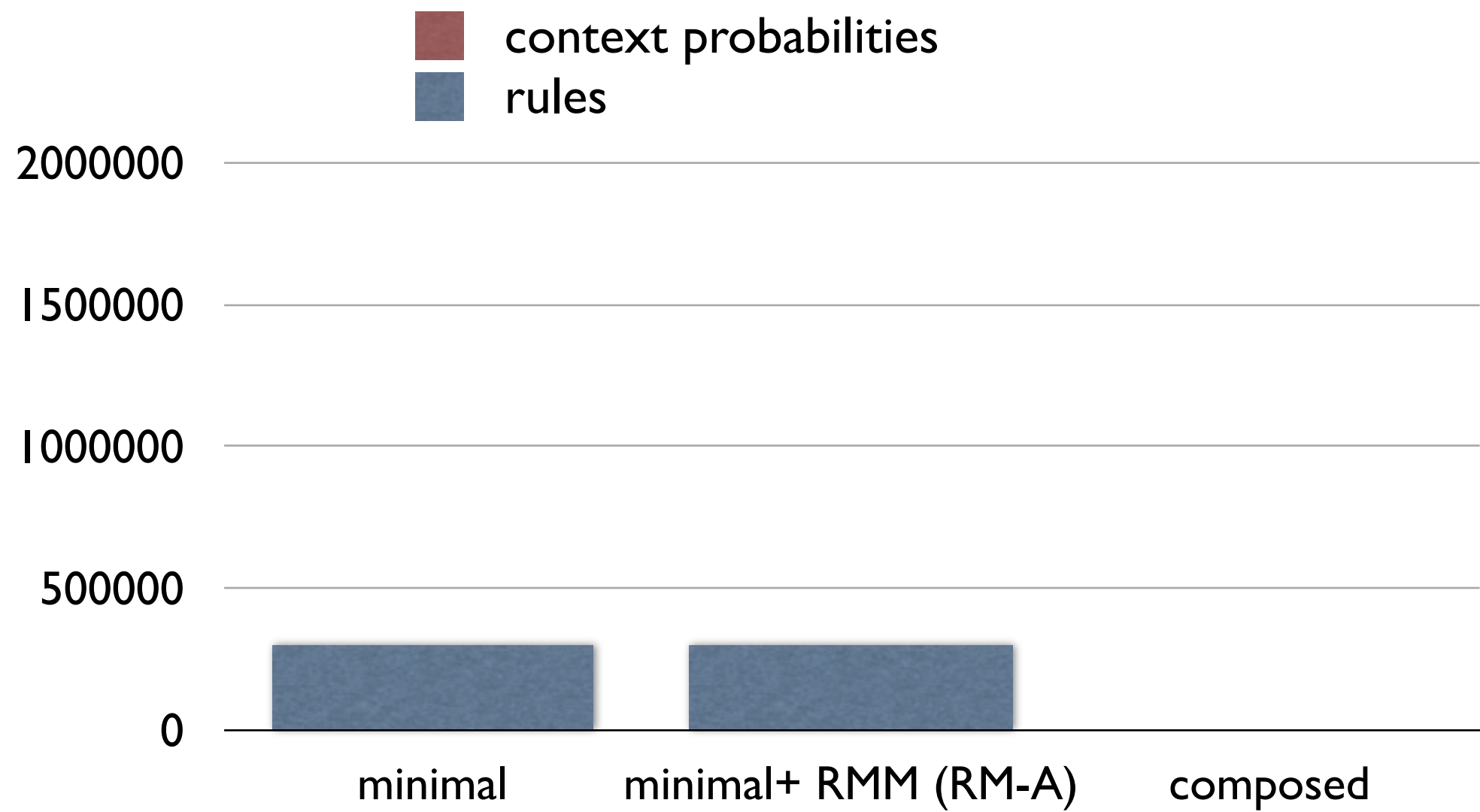
Results



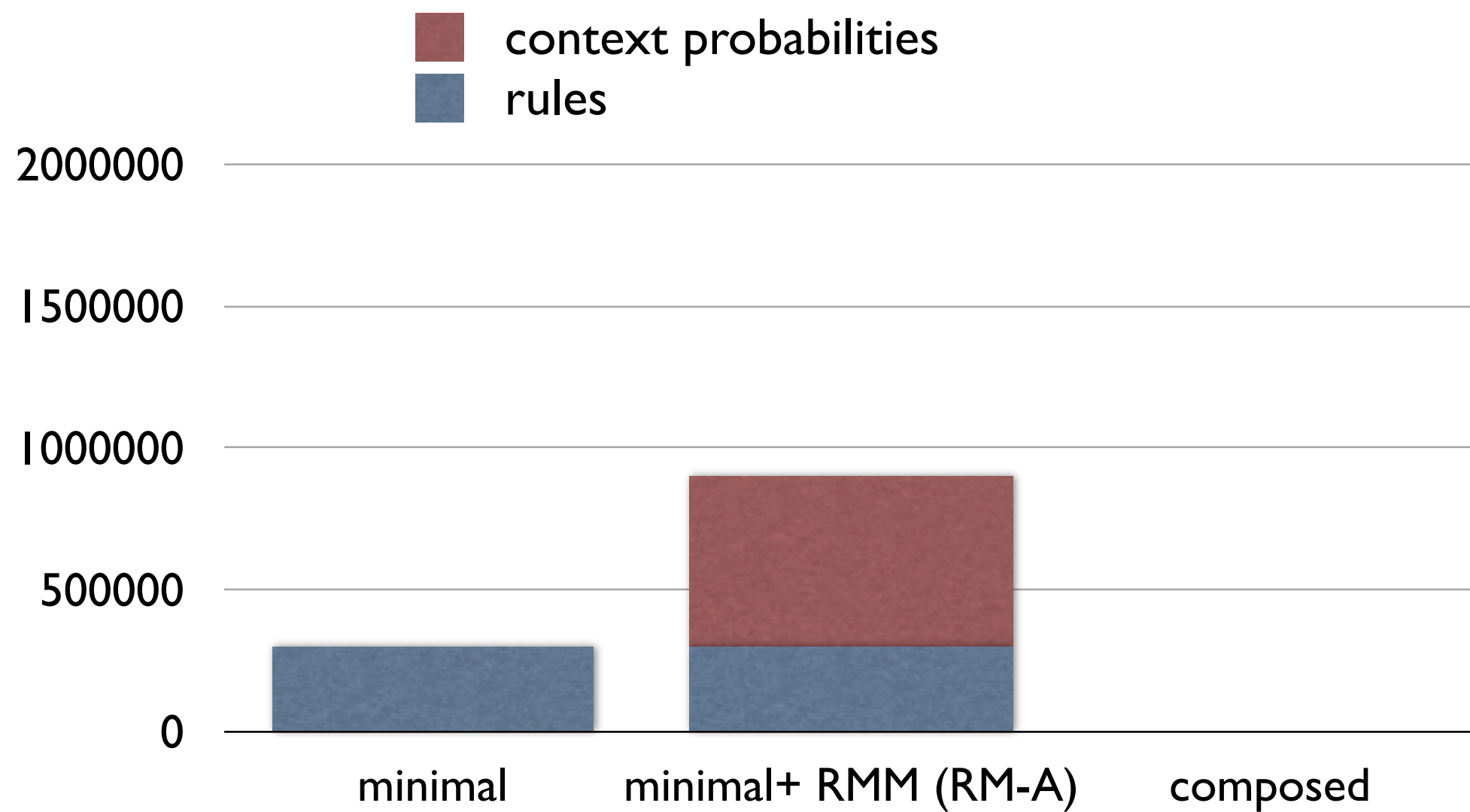
Results



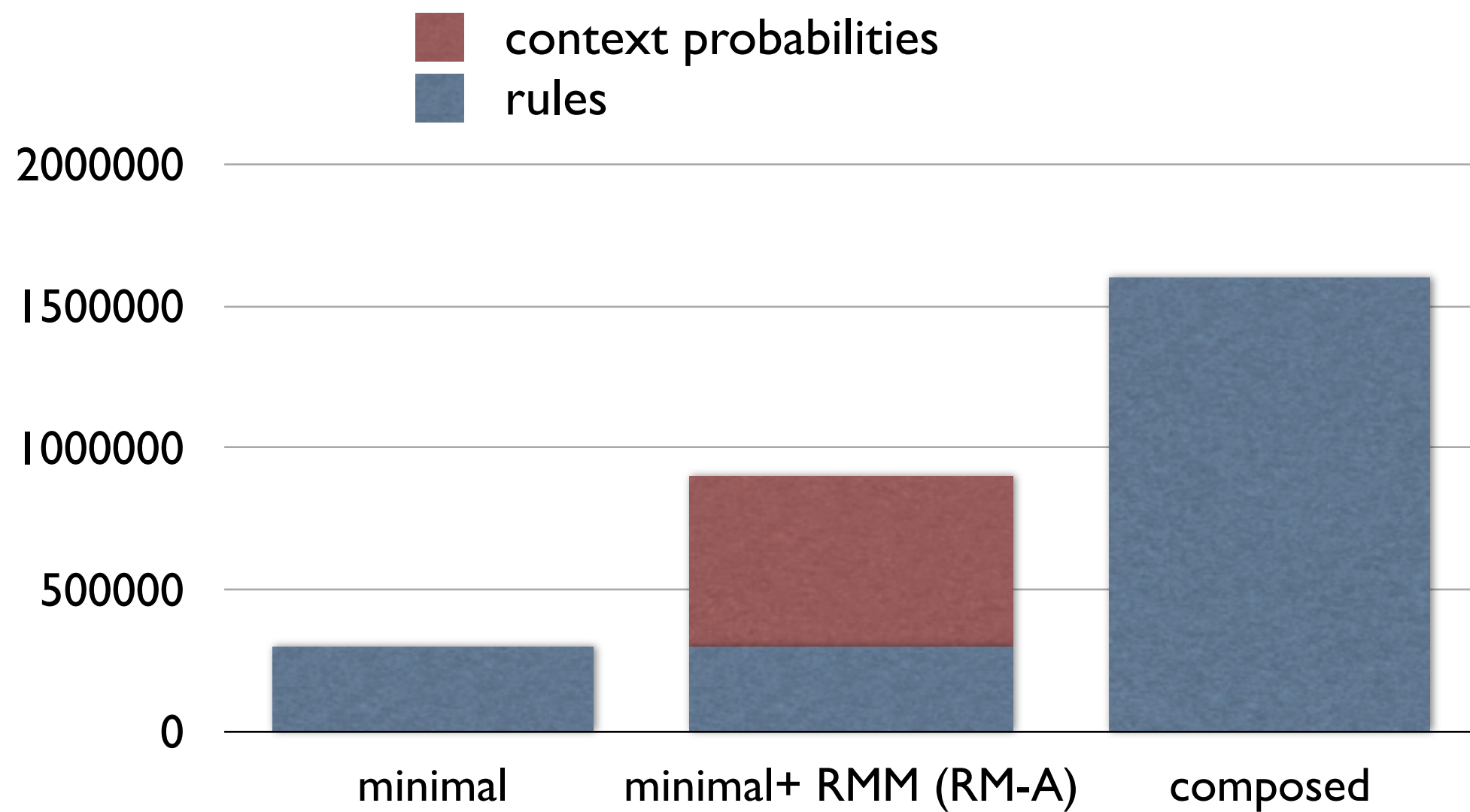
Results



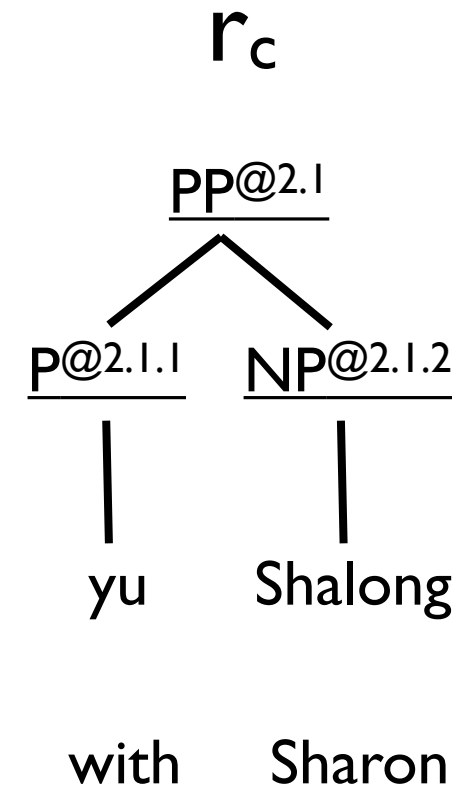
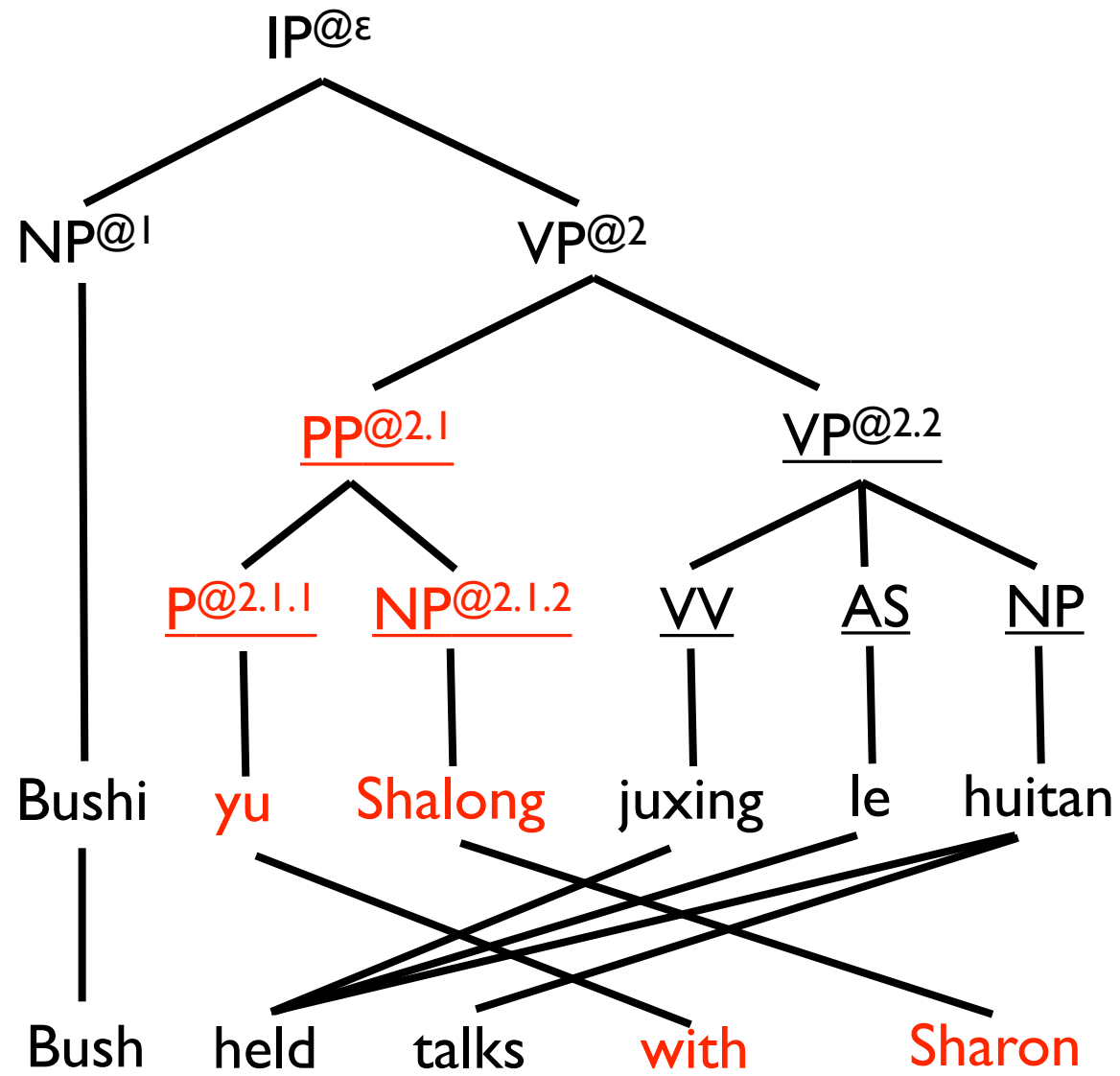
Results



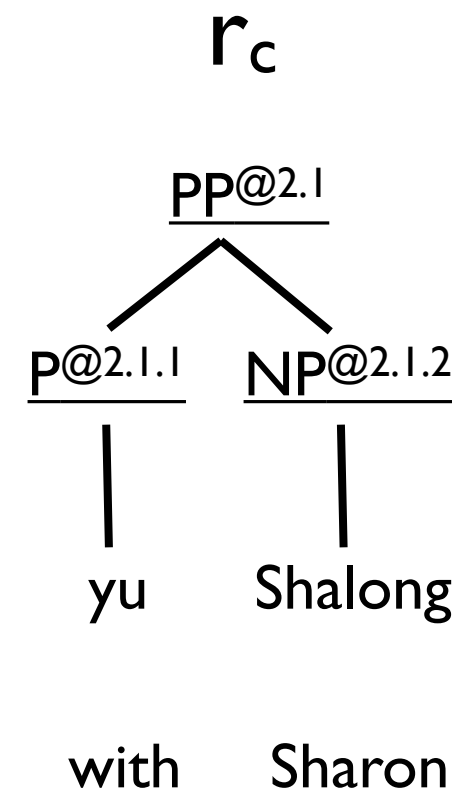
Results



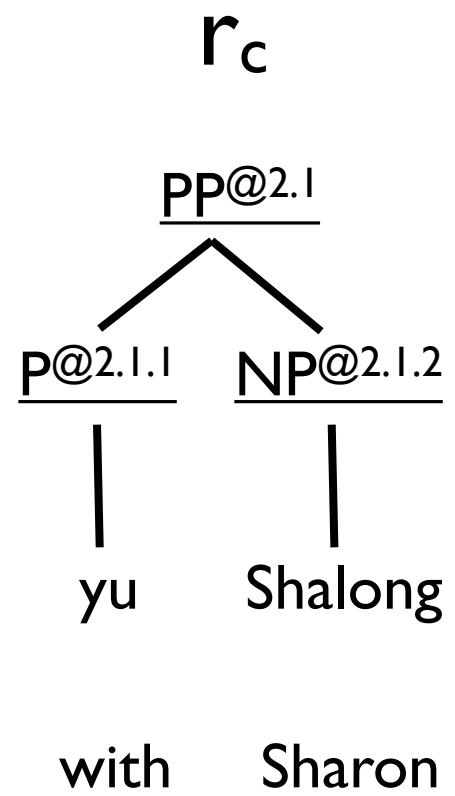
RMMs with composed Rules



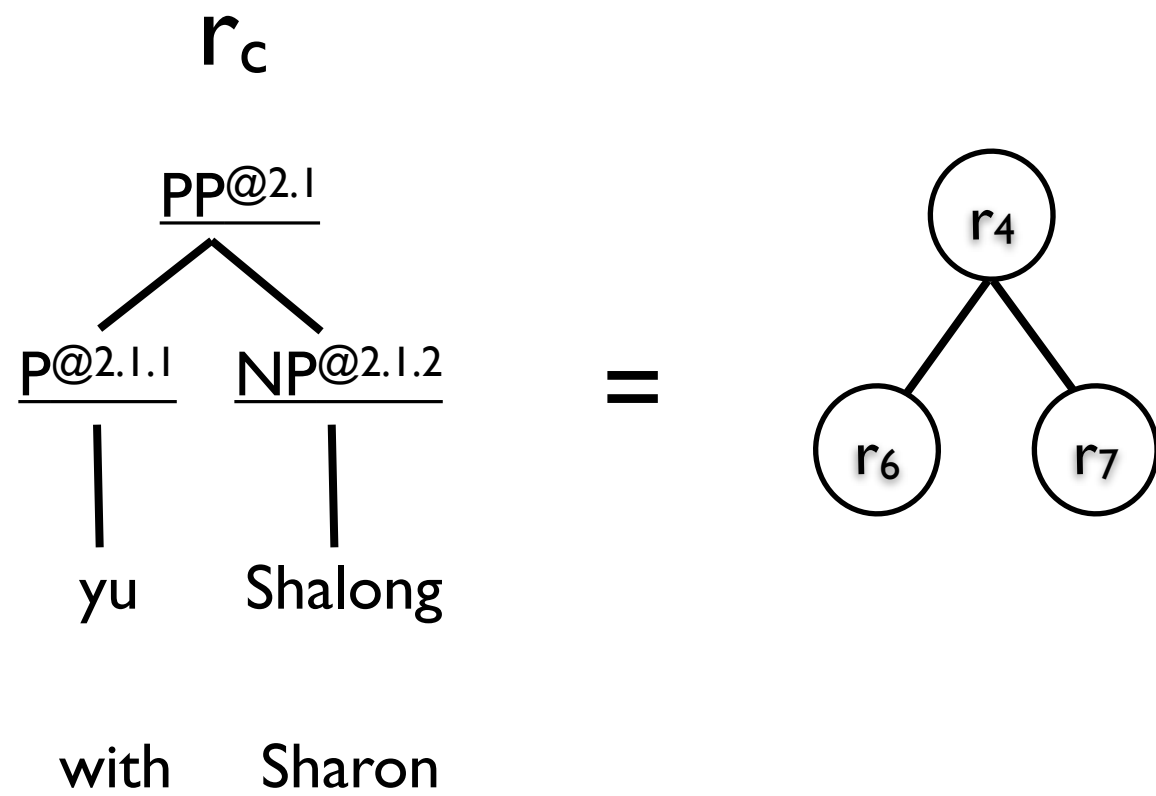
RMMs with composed Rules



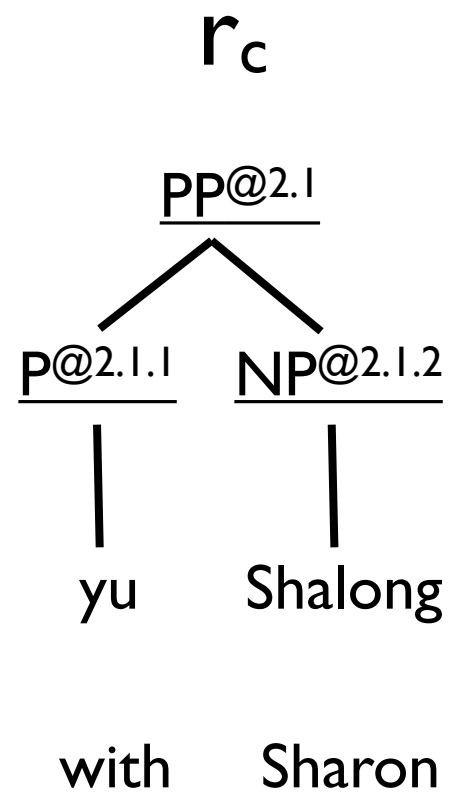
RMMs with composed Rules



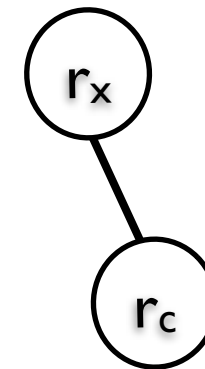
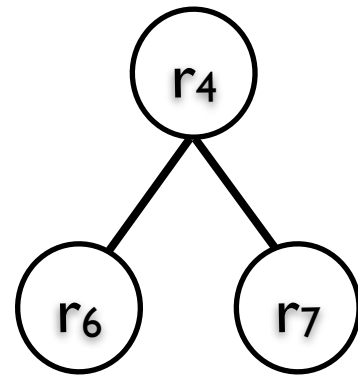
RMMs with composed Rules



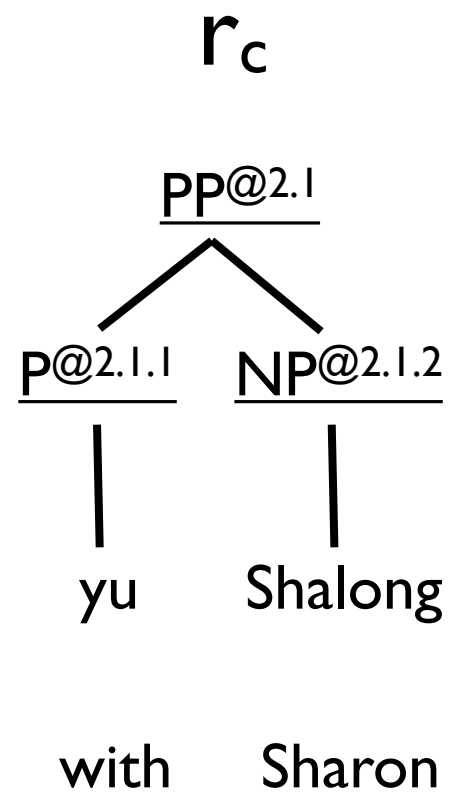
RMMs with composed Rules



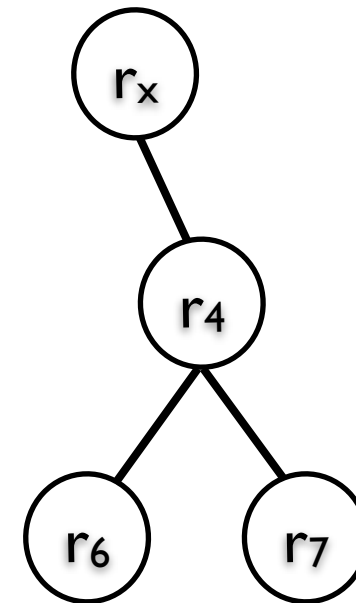
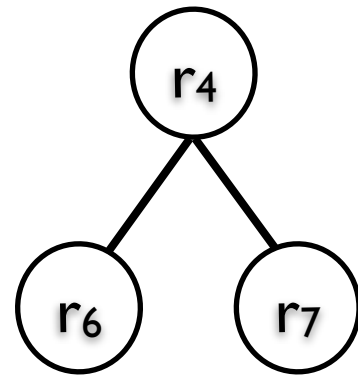
=



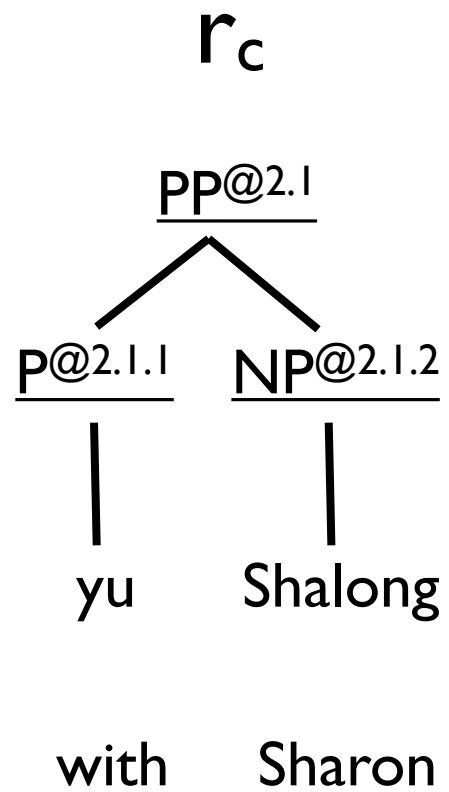
RMMs with composed Rules



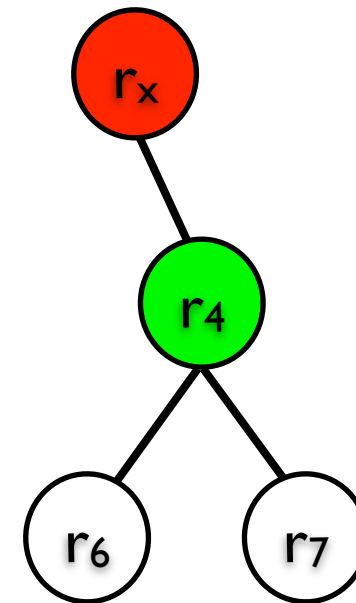
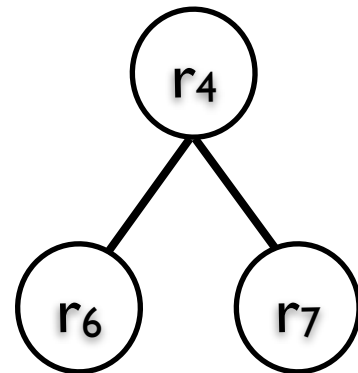
=



RMMs with composed Rules

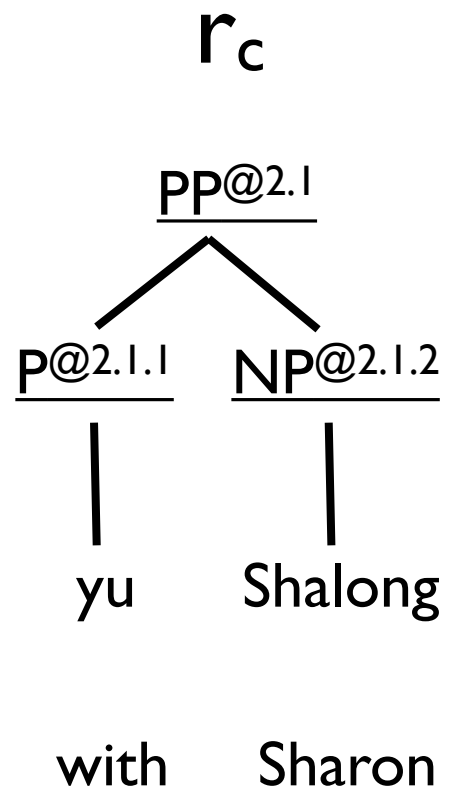


=

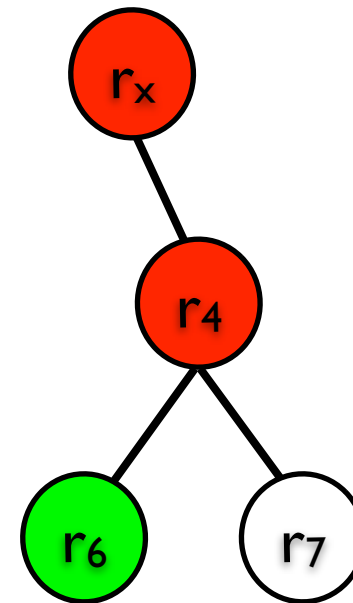
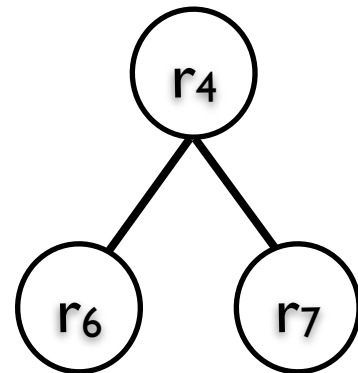


$$P(r_c|r_x) = P(r_4|r_x) \cdot P(r_6|r_4,r_x) \cdot P(r_7|r_4,r_x)$$

RMMs with composed Rules

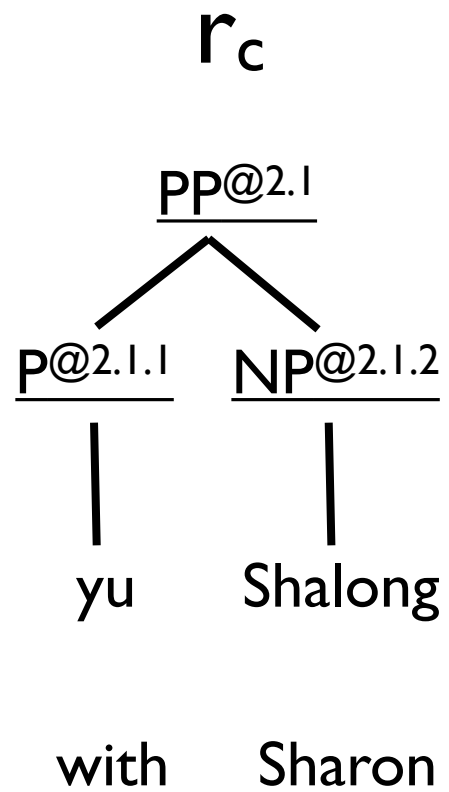


=

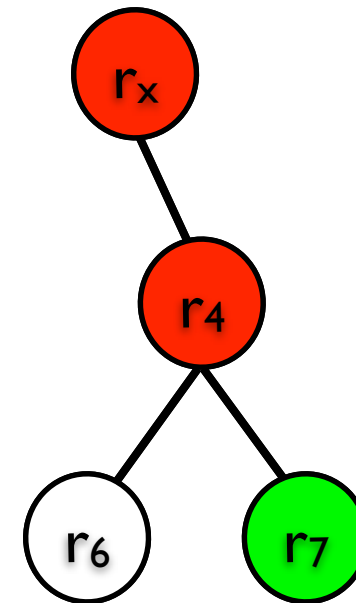
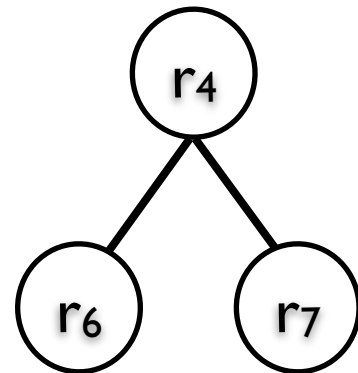


$$P(r_c|r_x) = P(r_4|r_x) \cdot P(r_6|r_4, r_x) \cdot P(r_7|r_4, r_x)$$

RMMs with composed Rules



=



$$P(r_c|r_x) = P(r_4|r_x) \cdot P(r_6|r_4,r_x) \cdot P(r_7|r_4,r_x)$$

Results: RMMs improve over composed rules

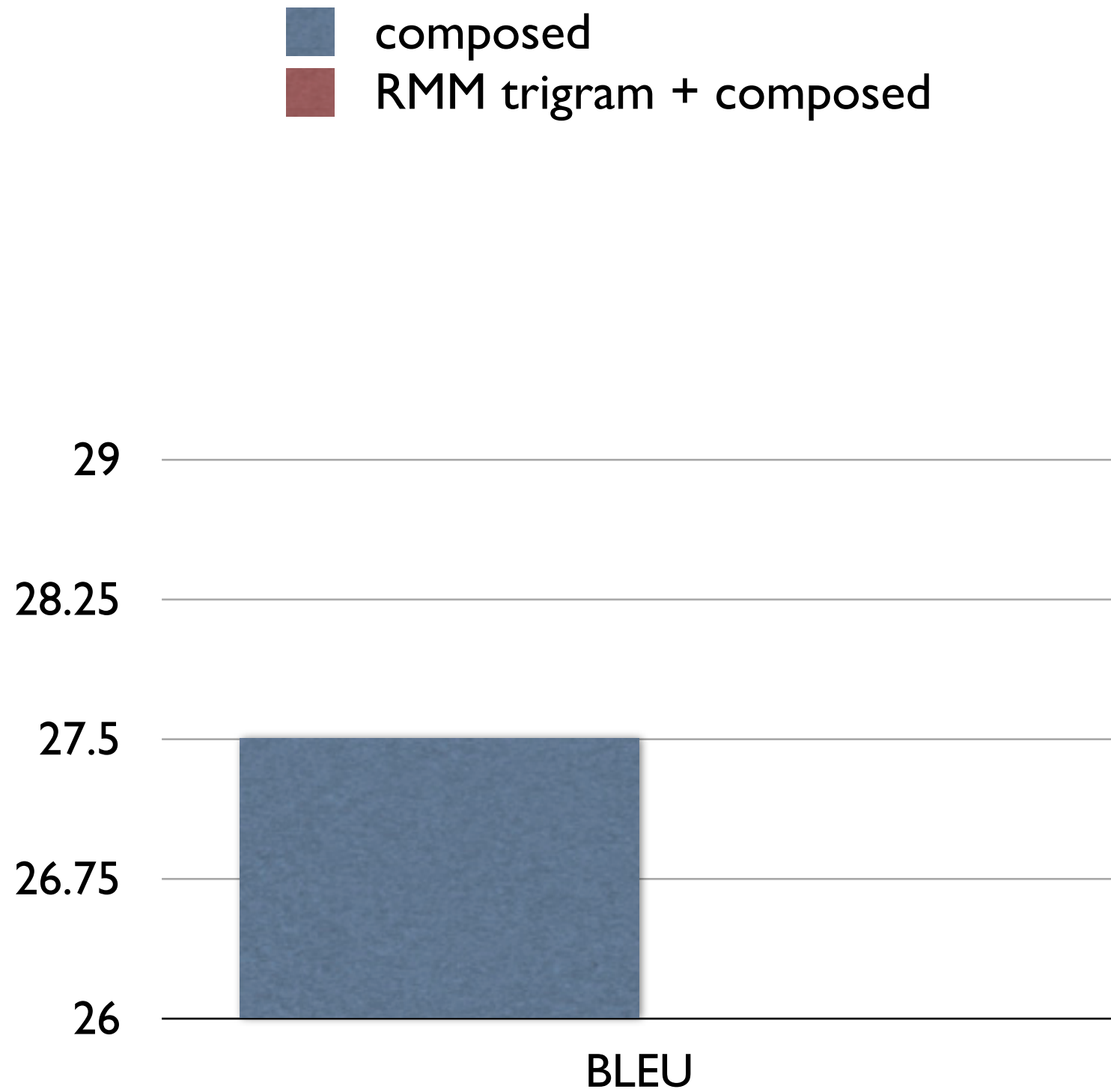
- composed
- RMM trigram + composed

Results: RMMs improve over composed rules

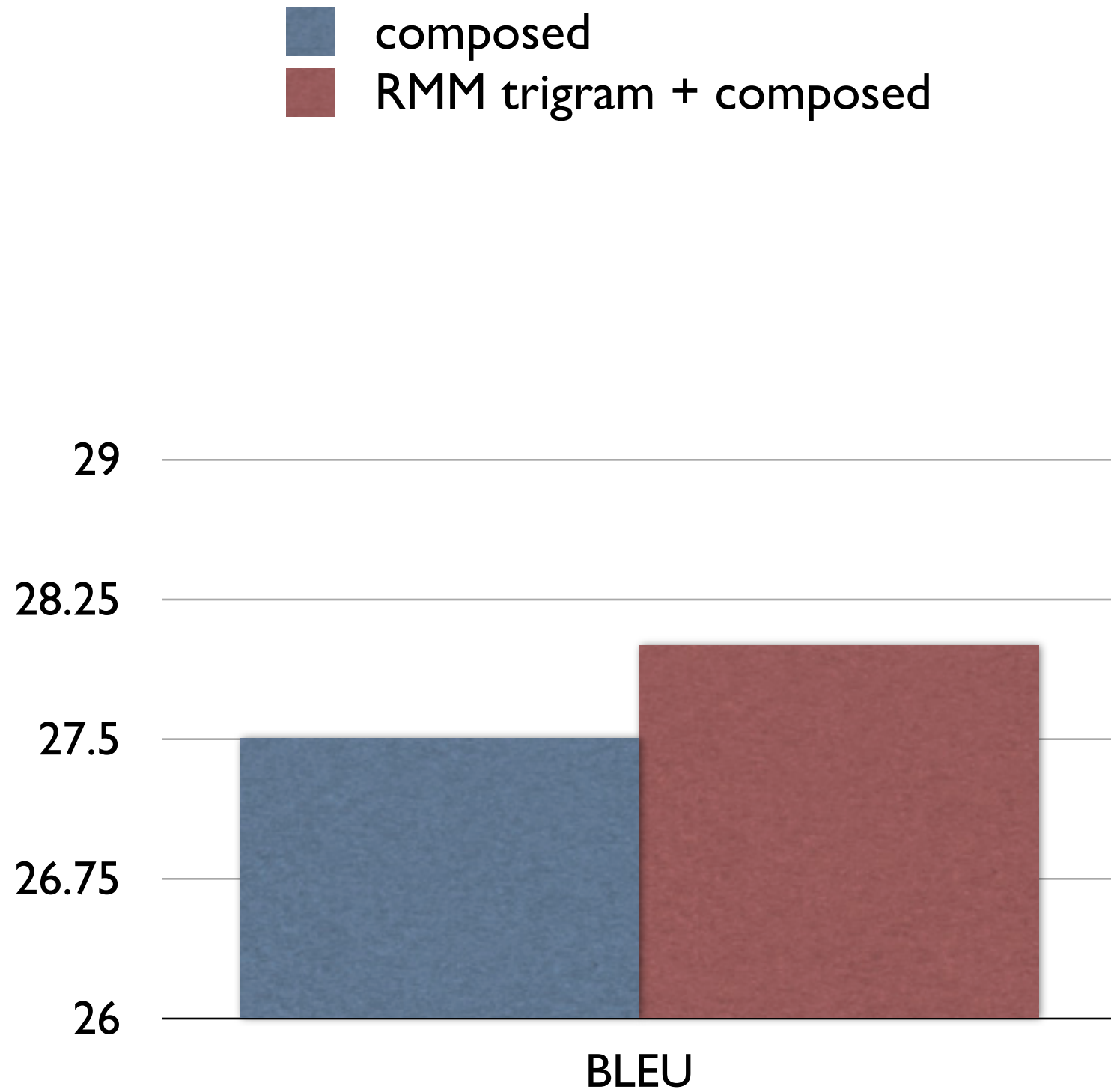
- composed
- RMM trigram + composed



Results: RMMs improve over composed rules



Results: RMMs improve over composed rules



Related Work

- Quirk and Menezes (2006)
- Ding and Palmer (2005)
- Liu and Gildea (2008)

Conclusion

- Using rule Markov models, we are able to get significant improvements in BLEU score.
- The grammar size and decoding time is less than the composed rule grammar
- Using rule Markov models with composed rule grammars further improves the BLEU score.

THANKS