

# A Generative Model for Punctuation in Dependency Parsing

Xiang Lisa Li\*, Dingquan Wang\*, and Jason Eisner

**\* Equal Contribution**

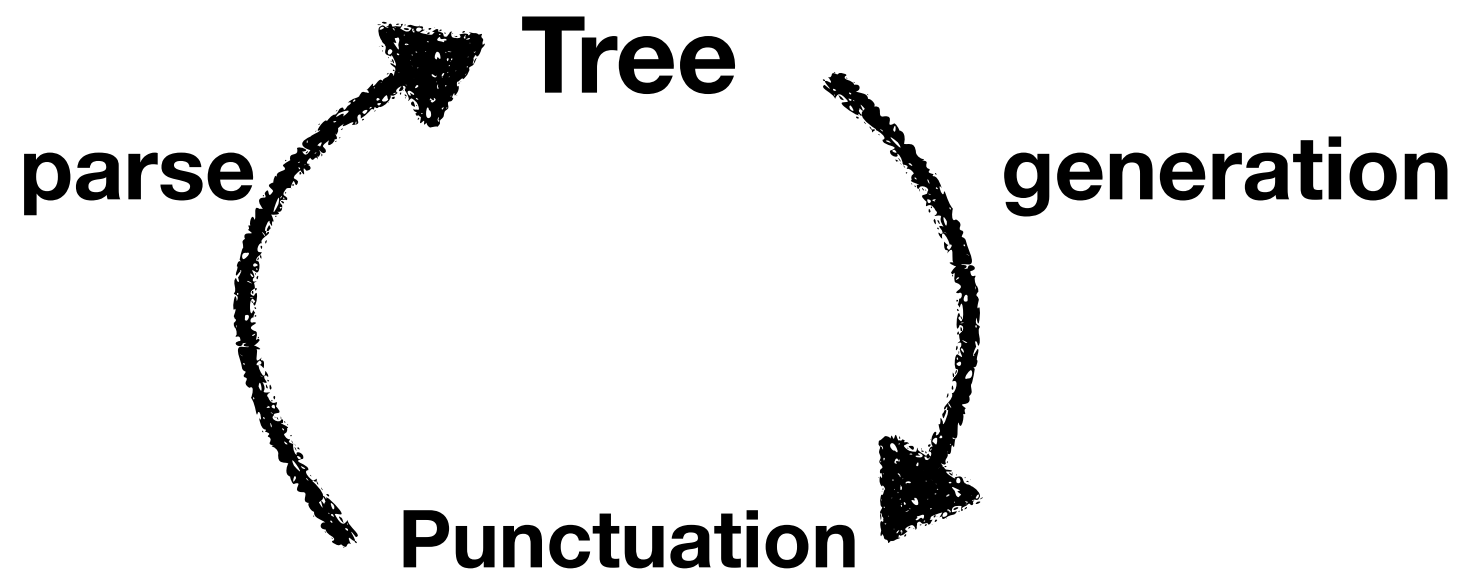


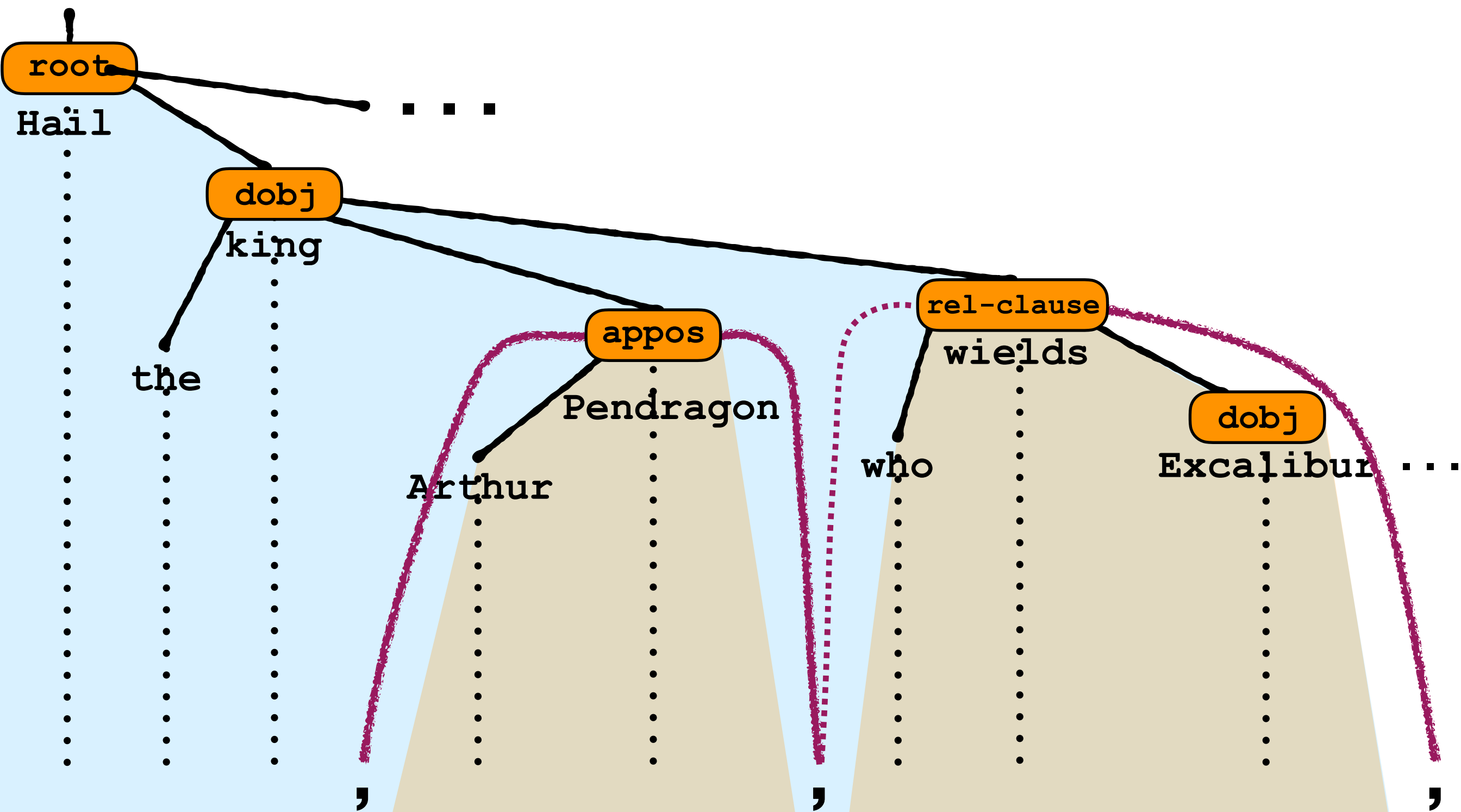
# NLP has neglected punctuation

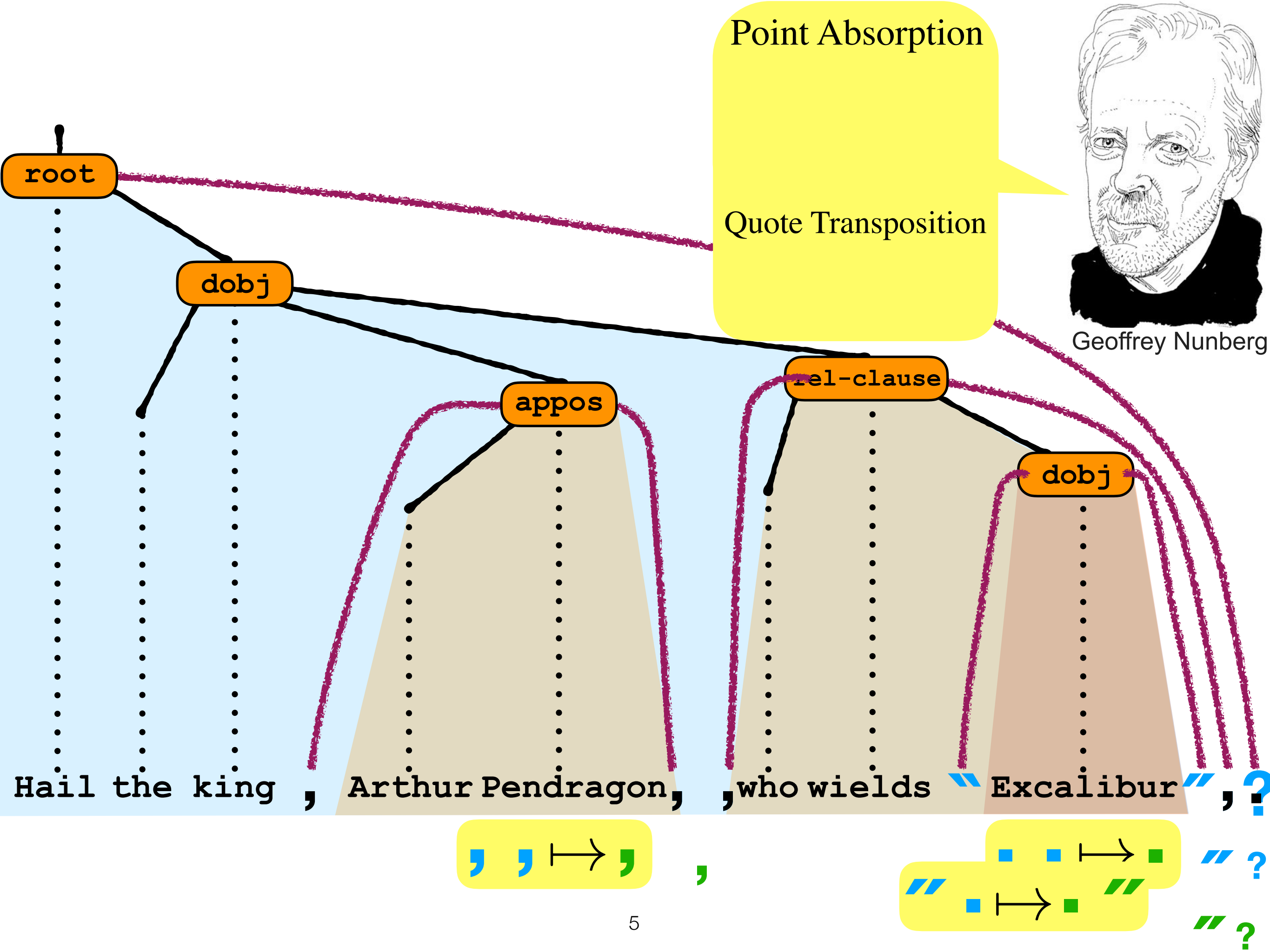
- ◆ Treebanks treat punctuation marks as ordinary tokens, but ARE THEY?
- ◆ *The Linguistics of Punctuation*, Nunberg (1990)

# Punctuation is useful

Punctuation marks are correlated with prosody and the syntactic tree structure.

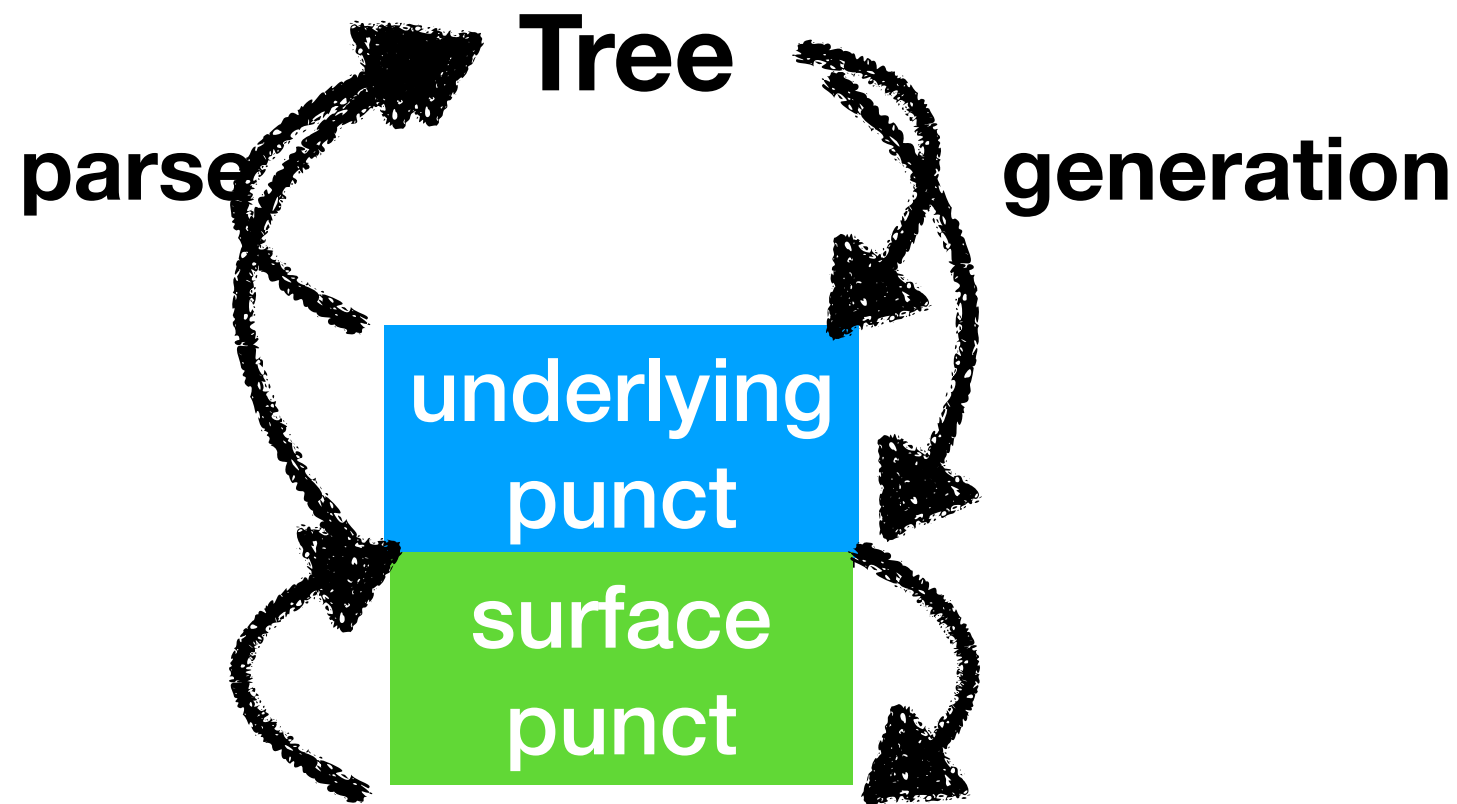






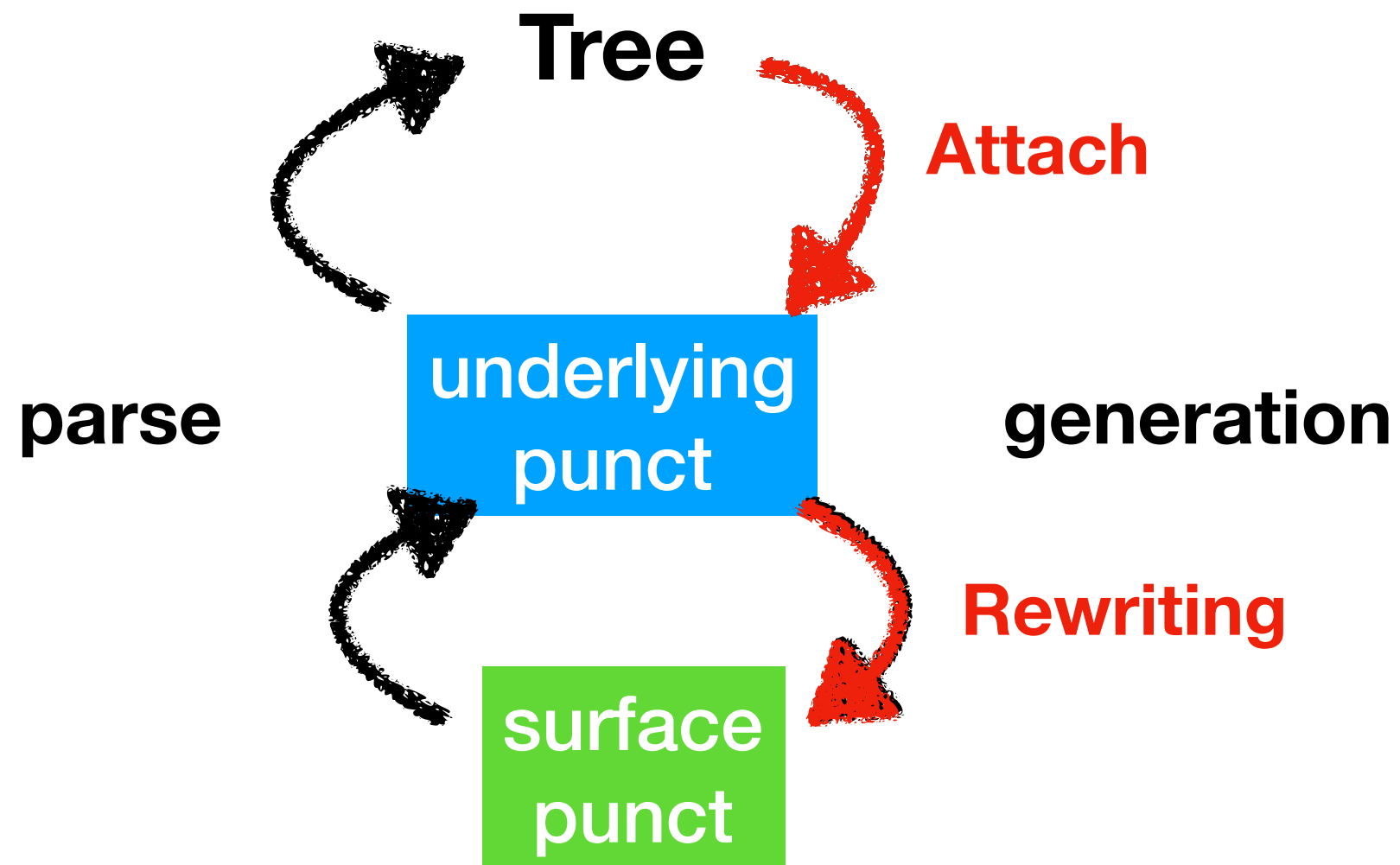
# Summary...

- Punctuation marks are not words.
  - Just as prosody is not words.
- They do not belong in the tree.
- Only underlying punctuation marks are in the tree, where they surround certain phrases.



The **surface** punctuation has a non-obvious correlation with the tree.

The **underlying** punctuation has a more direct correlation with the tree.



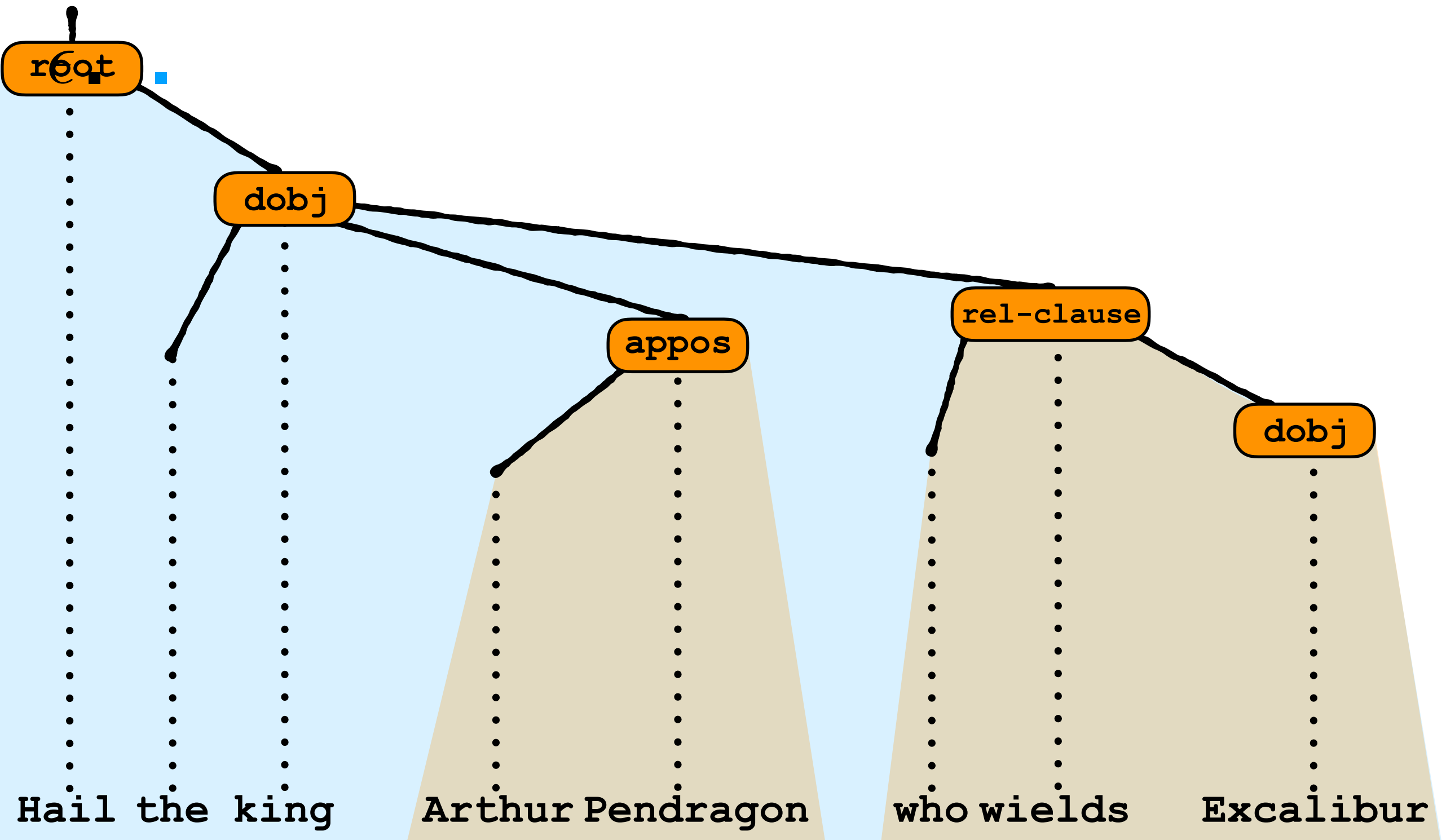
**Let's exploit the underlying punctuation  
under a generative model !**



$$P(\text{punctuated tree} \mid \text{unpunctuated tree})$$

$$= P(\epsilon \mid \text{root}, \text{unpunctuated tree})$$

**Attach**

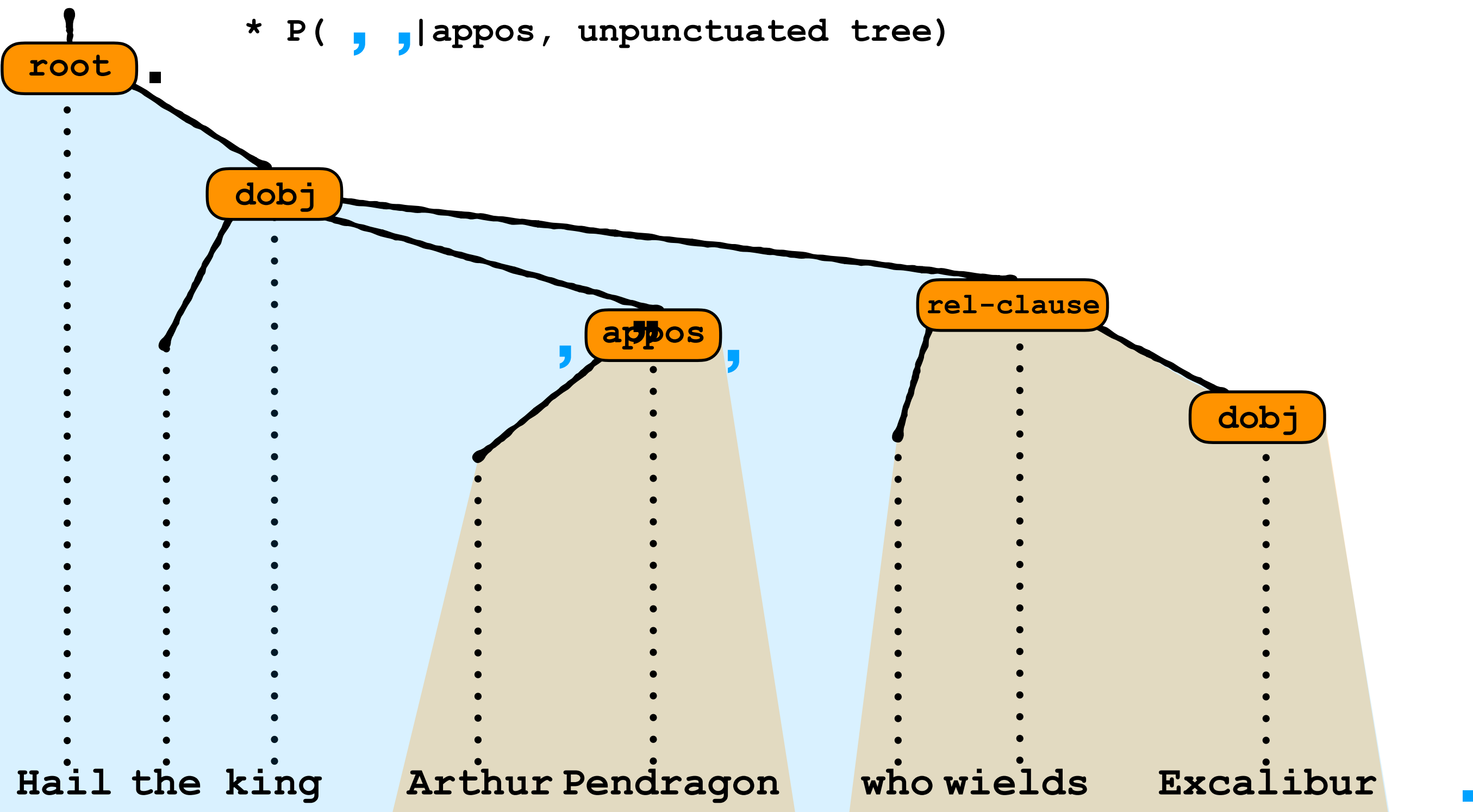


$P(\text{punctuated tree} \mid \text{unpunctuated tree})$

$= P(\epsilon \mid \text{root}, \text{unpunctuated tree})$

$* P(, \mid \text{appos}, \text{unpunctuated tree})$

**Attach**



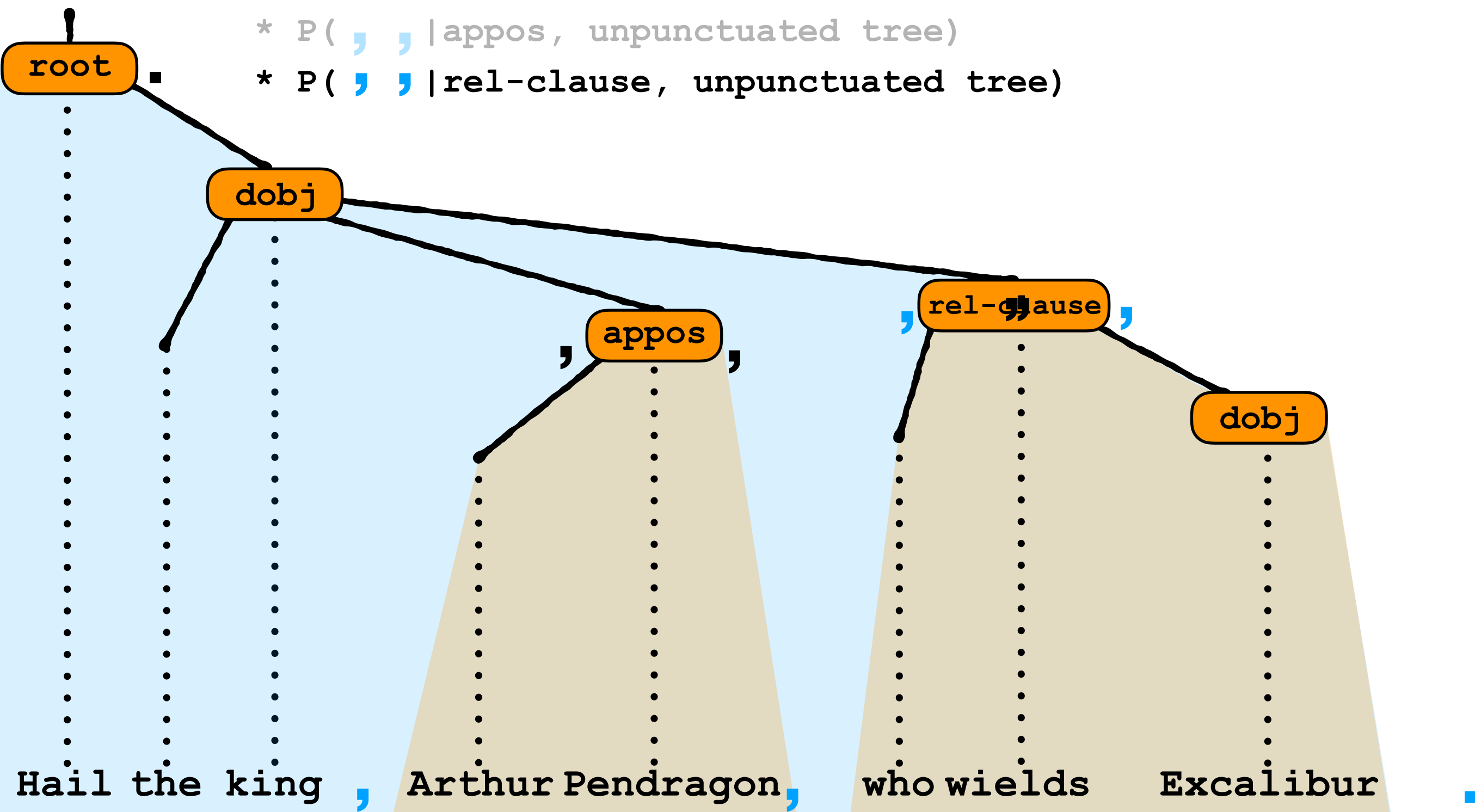
$P(\text{punctuated tree} \mid \text{unpunctuated tree})$

$= P(\epsilon \mid \text{root}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{appos}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{rel-clause}, \text{unpunctuated tree})$

**Attach**



$P(\text{punctuated tree} \mid \text{unpunctuated tree})$

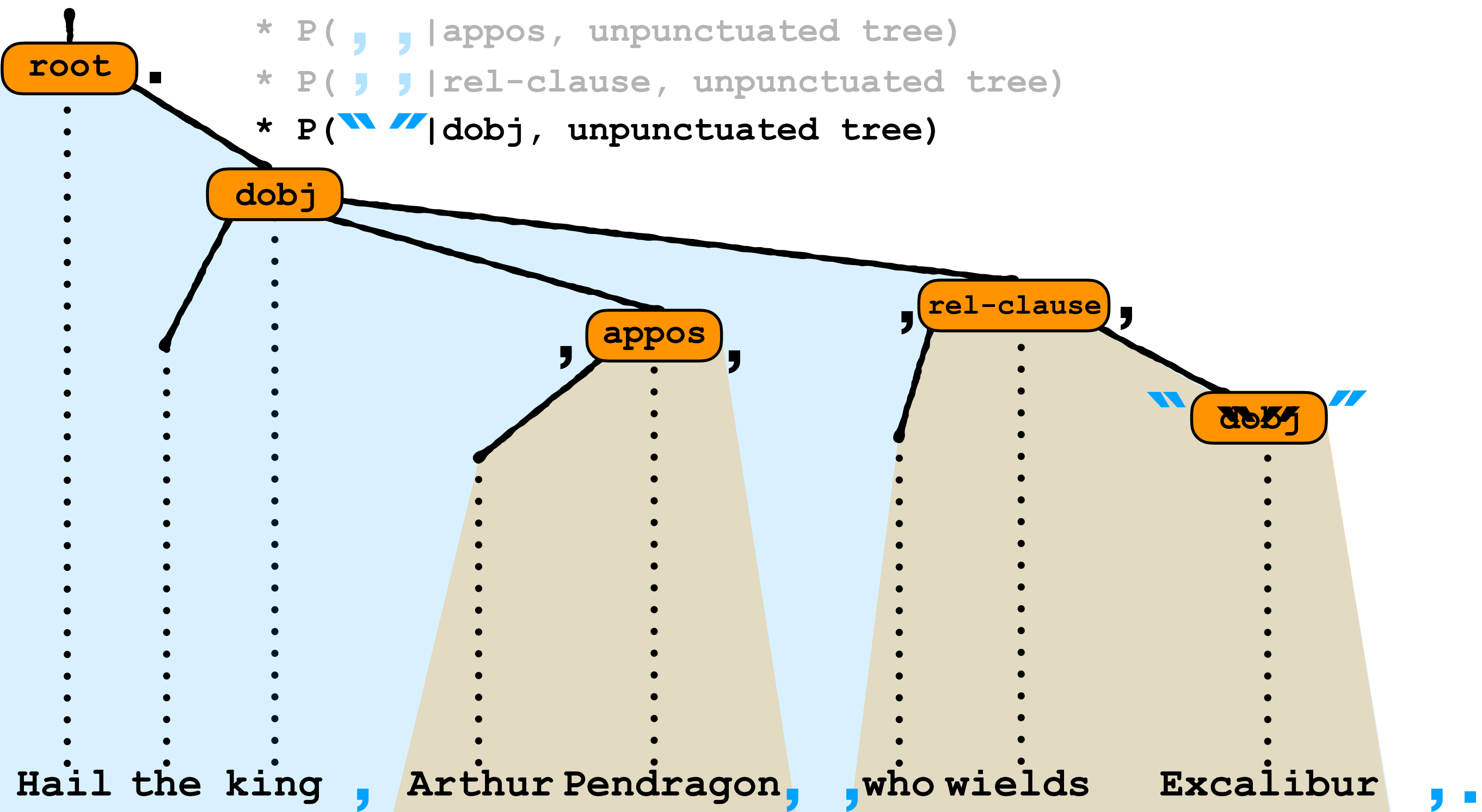
$= P(\epsilon \mid \text{root}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{appos}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{rel-clause}, \text{unpunctuated tree})$

$\quad * P(// \mid \text{dobj}, \text{unpunctuated tree})$

**Attach**



$P(\text{punctuated tree} \mid \text{unpunctuated tree})$

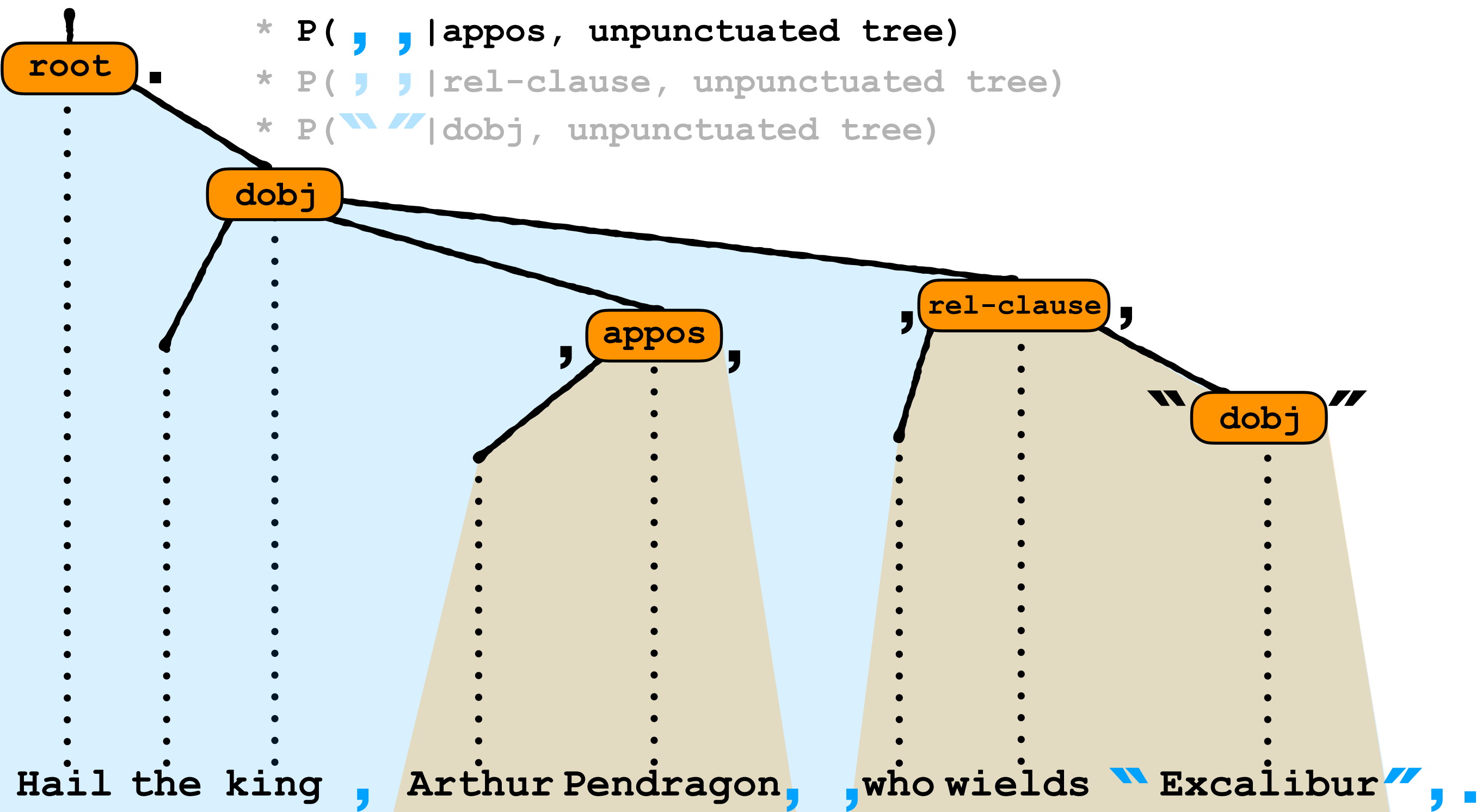
$= P(\epsilon \mid \text{root}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{appos}, \text{unpunctuated tree})$

$\quad * P(, \mid \text{rel-clause}, \text{unpunctuated tree})$

$\quad * P(\text{"/"} \mid \text{dobj}, \text{unpunctuated tree})$

**Attach**



$$P( \text{ , } \text{ , } \mid \text{ appos , unpunctuated tree } ) = 0.4$$

**appos**

left punct	right punct
,	,
“	”
€	,
,	€
€	€

**Prob**

0.4

0.3

0.05

0.05

0.2



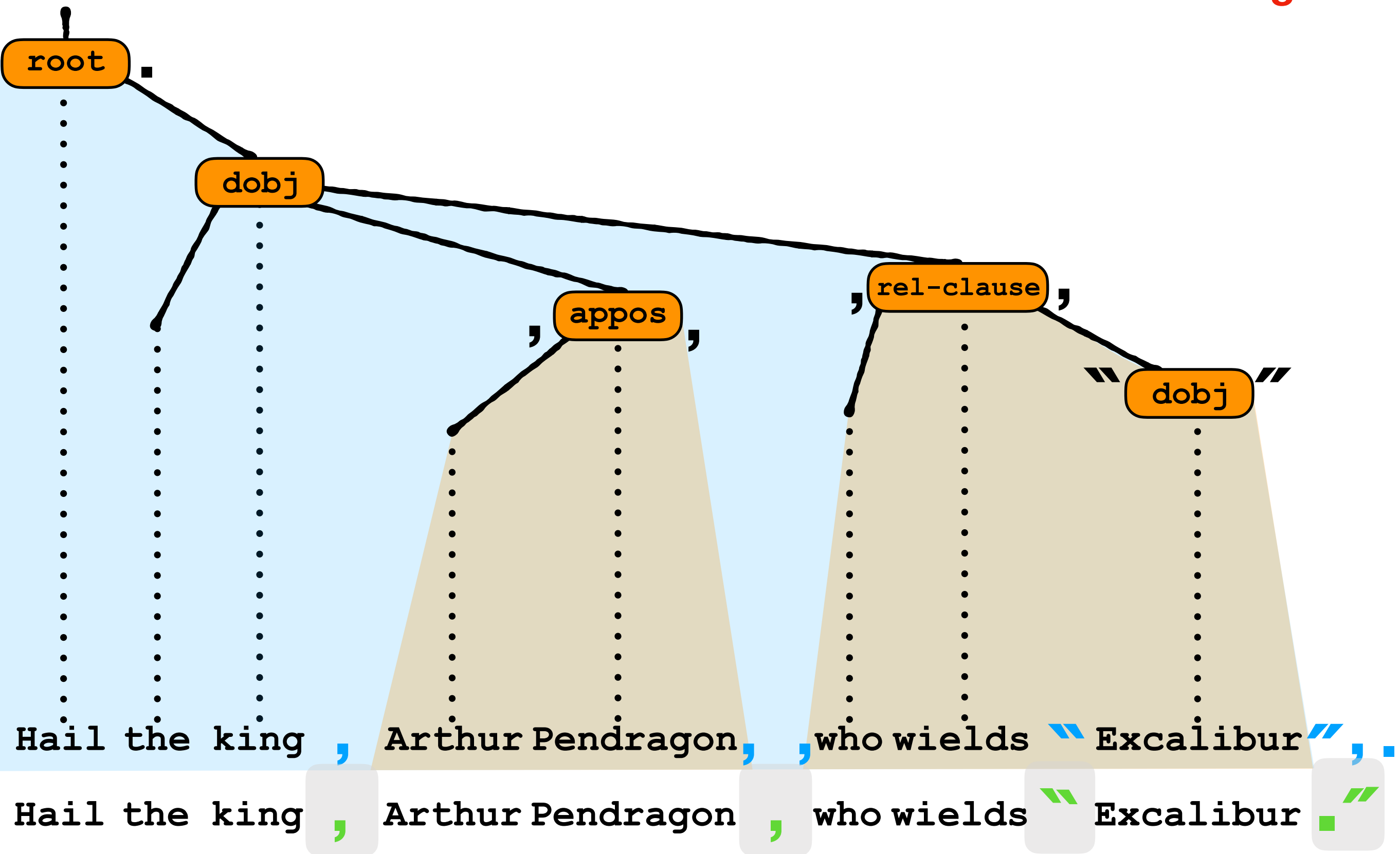
f1	f2	f3	f4
0	1	2	1
0	1	0	0
0	0	0	1
0	0	0	1
1	0	0	0

...

- **f1: indicator feature for no punctuation ( € € ).**
- **f2: indicator feature that checks for symmetry.**
- **f3: span length feature for attaching commas.**
- **f4: indicator feature that checks for internal comma.**

P(surface punct | punctuated tree)

Rewriting



# Rewriting — A Sliding Window Model



Left-absorb  $P ( \text{, ,} \mapsto \text{,} )$

Left-absorb  $P ( \text{, ,} \mapsto \text{,} )$

Left-absorb  $P ( \text{, )} \mapsto \text{,} )$

identity  $P ( \text{)"} \mapsto \text{,} " )$

swap  $P ( \text{"} \square \mapsto \square " )$

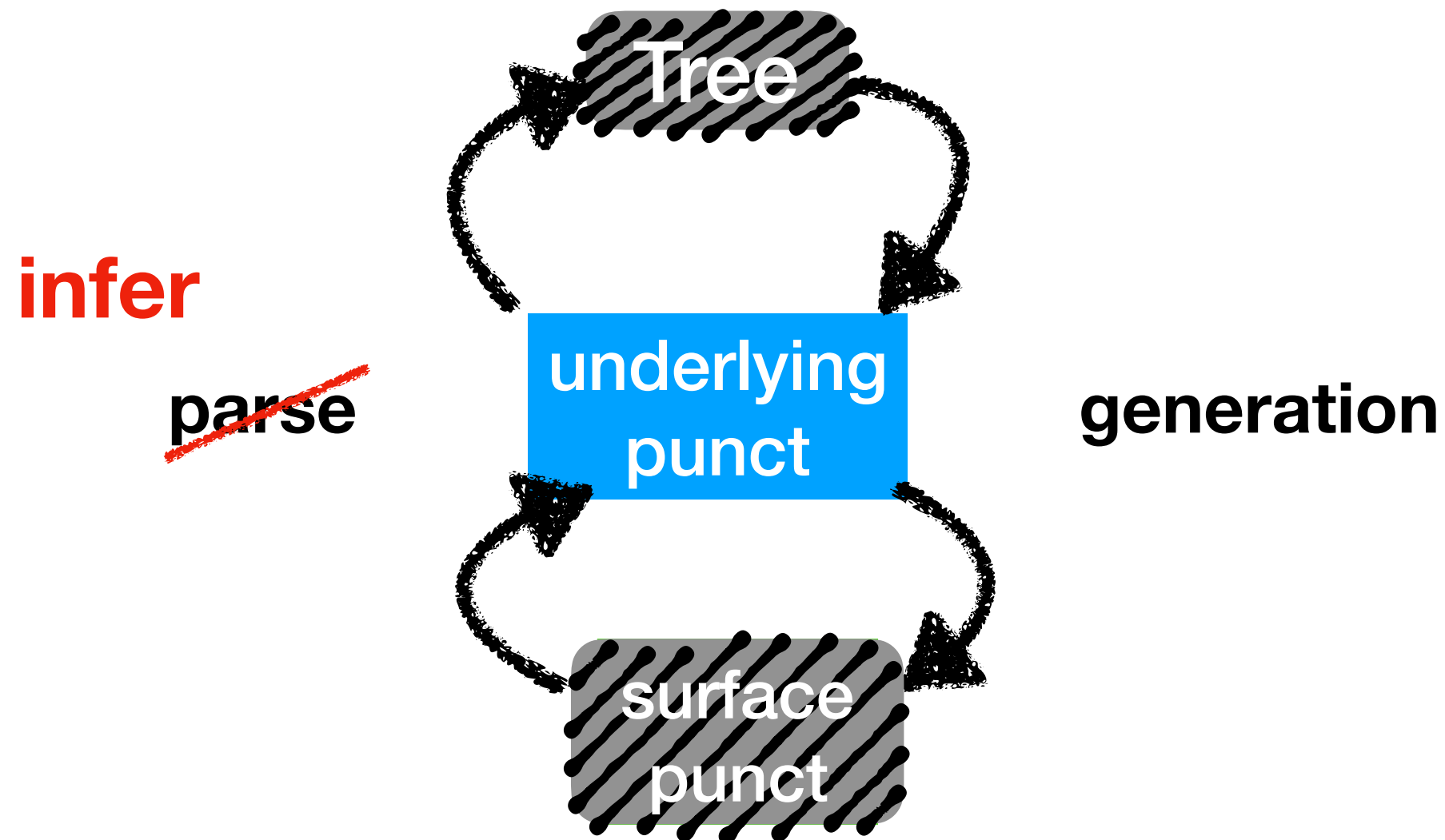


Geoffrey Nunberg



# Summary generative story

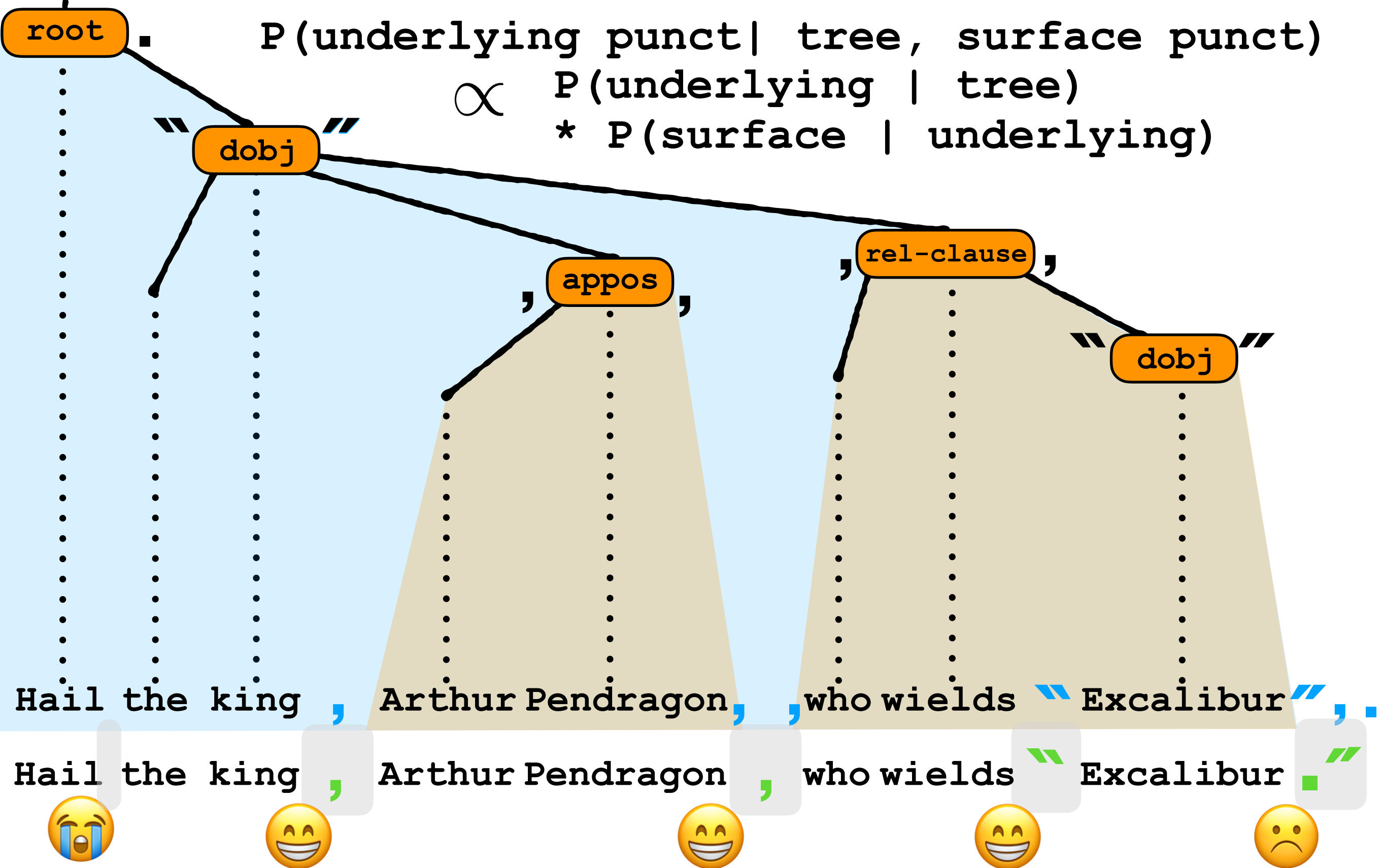
- Underlying punctuation is generated at each tree node (not quite independently).
- Total underlying punctuation at each slot between words is rewritten into surface punctuation (independently).



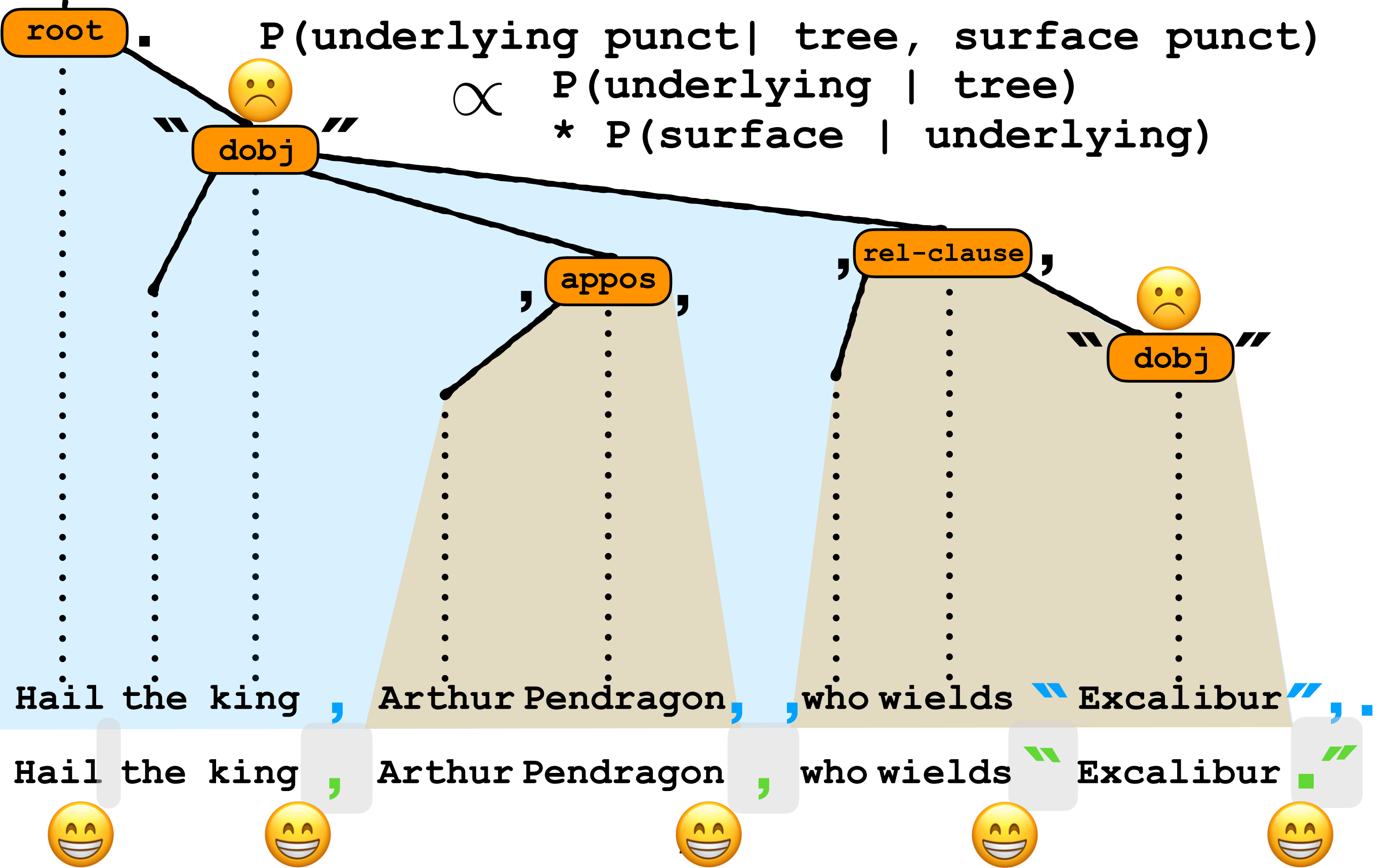
**In training, we observe the tree and surface punct.  
Want to recover the underlying punctuation.**

$P(\text{underlying punct} \mid \text{tree}, \text{surface punct})$

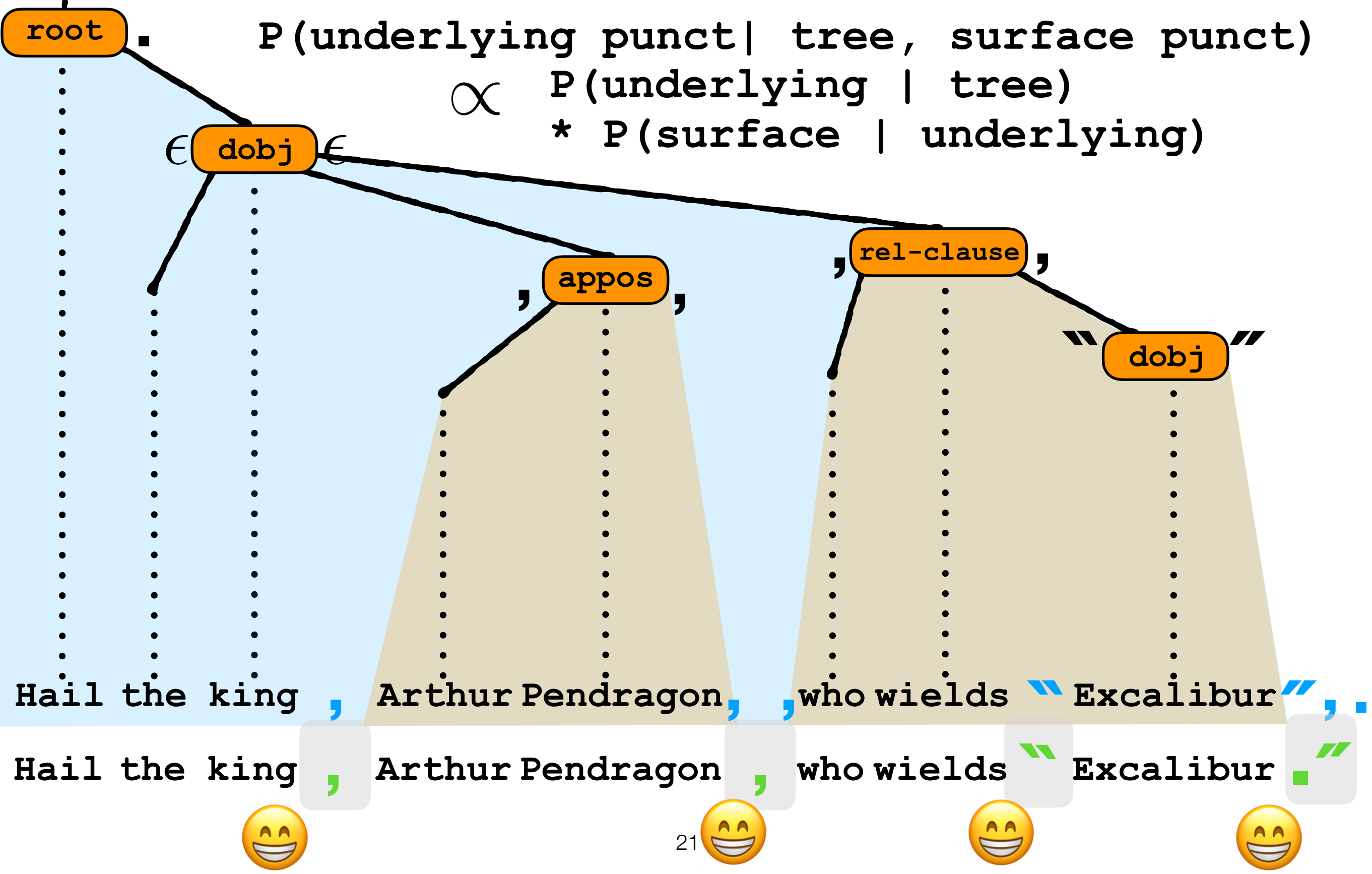
# Method 1. A Sampling Approach



# Method 1. A Sampling Approach



# Method 1. A Sampling Approach



# Method 2. $O(n)$ Dynamic Programming

**Idea : Make a (weighted) CFG that generates just the possible underlyingly punctuated trees for *this* sentence**

**Soft constraint ① : Underlying punctuation fits the tree structure**

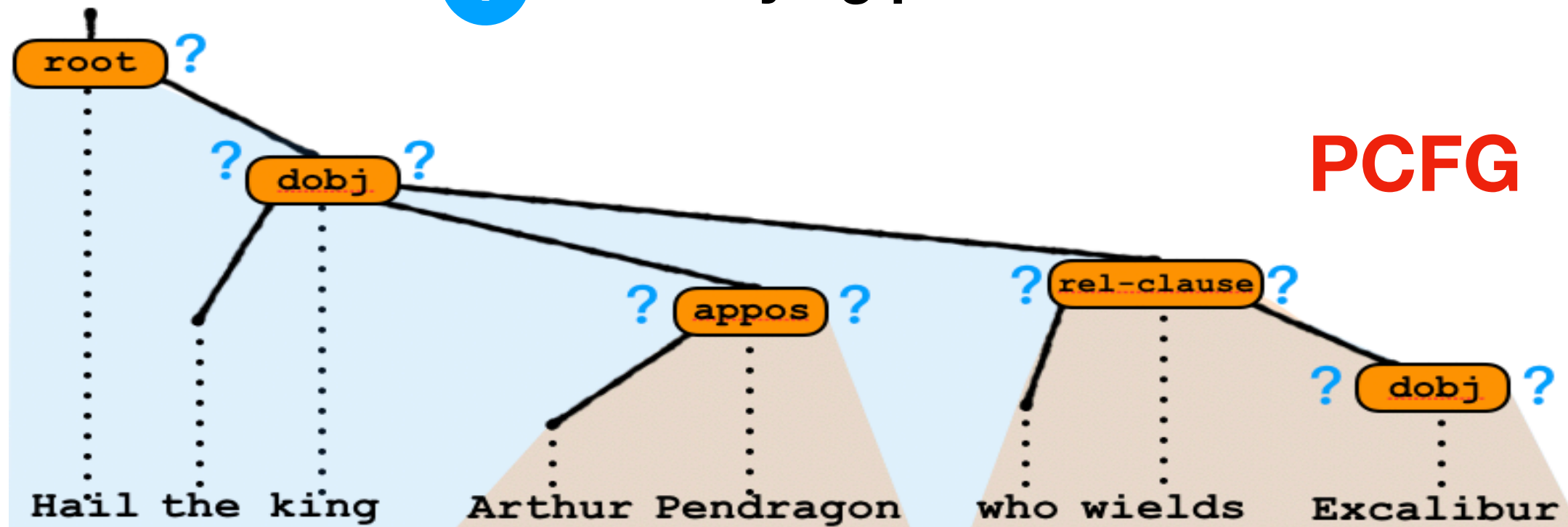
**Soft constraint ② : Underlying punctuation fits surface punctuation**

**Intersect ① ②  $\rightarrow$  our CFG**

# Method 2. $O(n)$ Dynamic Programming

Idea : Make a (weighted) CFG that generates just the possible underlyingly punctuated trees for *this* sentence

Soft constraint ① : Underlying punctuation fits the tree structure



Soft constraint ② : Underlying punctuation fits surface punctuation

Intersect ① ② → our CFG

# Method 2. $O(n)$ Dynamic Programming

Idea : Make a (weighted) CFG that generates just the possible underlyingly punctuated trees for *this* sentence

Soft constraint ① : Underlying punctuation fits the tree structure  
**PCFG**

Soft constraint ② : Underlying punctuation fits surface punctuation  
**weighted FSA**

Hail the king  Arthur Pendragon  who wields  Excalibur 

Intersect ① ②  **our CFG**



# Method 2. O(n) Dynamic Programming

Idea : Make a (weighted) CFG that generates just the possible underlyingly punctuated trees for *this* sentence

Soft constraint ① : Underlying punctuation fits the tree structure  
**PCFG**

Soft constraint ② : Underlying punctuation fits surface punctuation

Intersect ① ②  $\rightarrow$  our CFG **weighted FSA**

$$\text{PCFG} \cap \text{weighted FSA} = \text{weighted CFG}$$

# Method 2. $O(n)$ Dynamic Programming

**Idea : Make a (weighted) CFG that generates just the possible underlyingly punctuated trees for *this* sentence**

**Soft constraint ① : Underlying punctuation fits the tree structure**

**Soft constraint ② : Underlying punctuation fits surface punctuation**

**Intersect ① ②  $\rightarrow$  our CFG**

**Now we can find the best underlyingly punctuated tree**

- **or sum over all of them for the model likelihood in training**

**For our CFG, this can be done in  $O(n)$  time**

# Results

**Rules Learned from the Noisy Channel**

**Analysis of Attachment model**

**Punctuation Restoration**

**Syntactic Rephrasing**

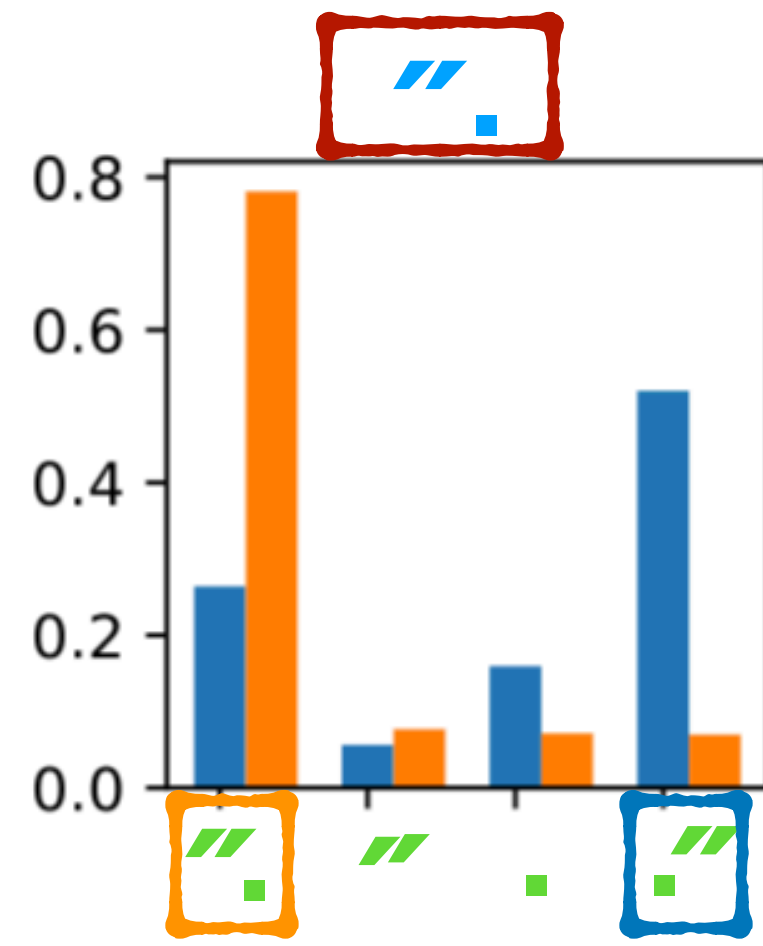
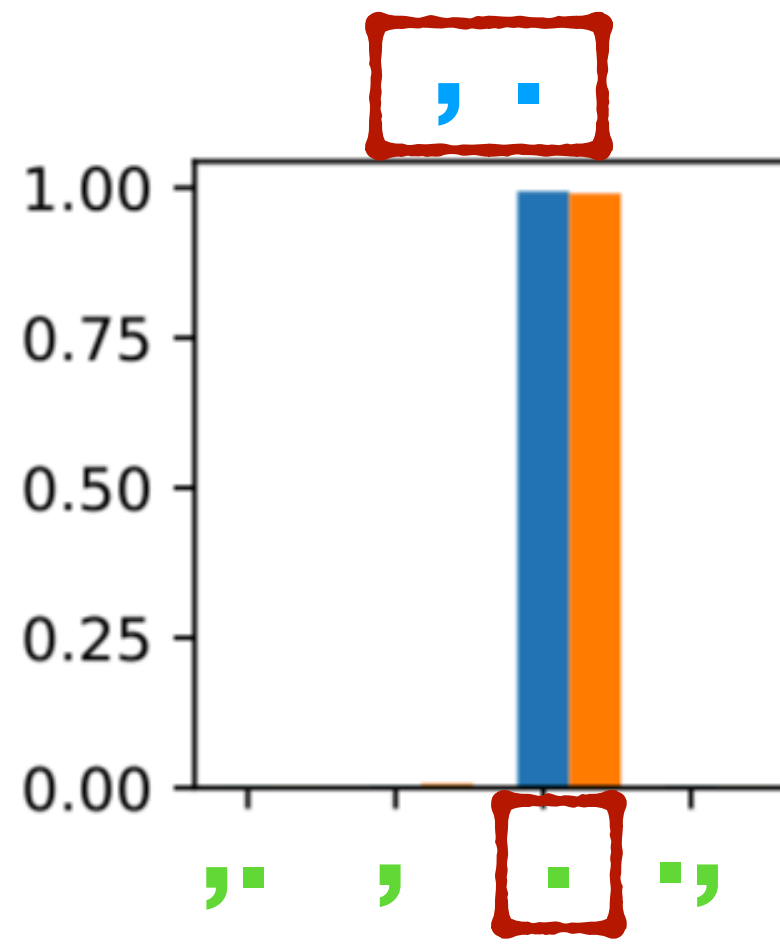
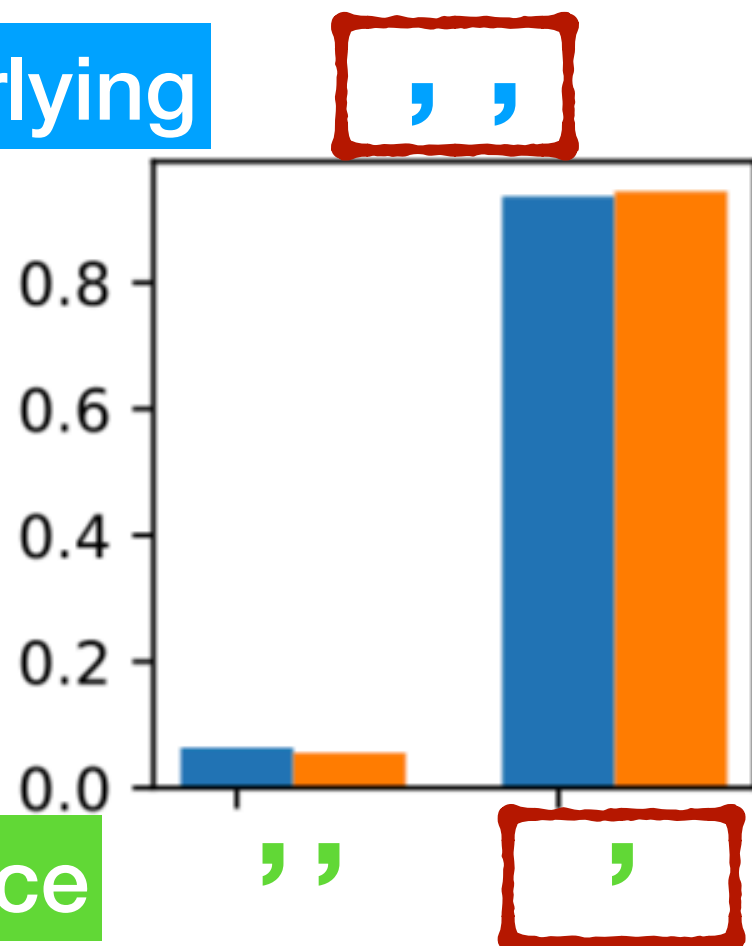
# We learn high probability for Nunberg's English Rules

American English

British English

underlying

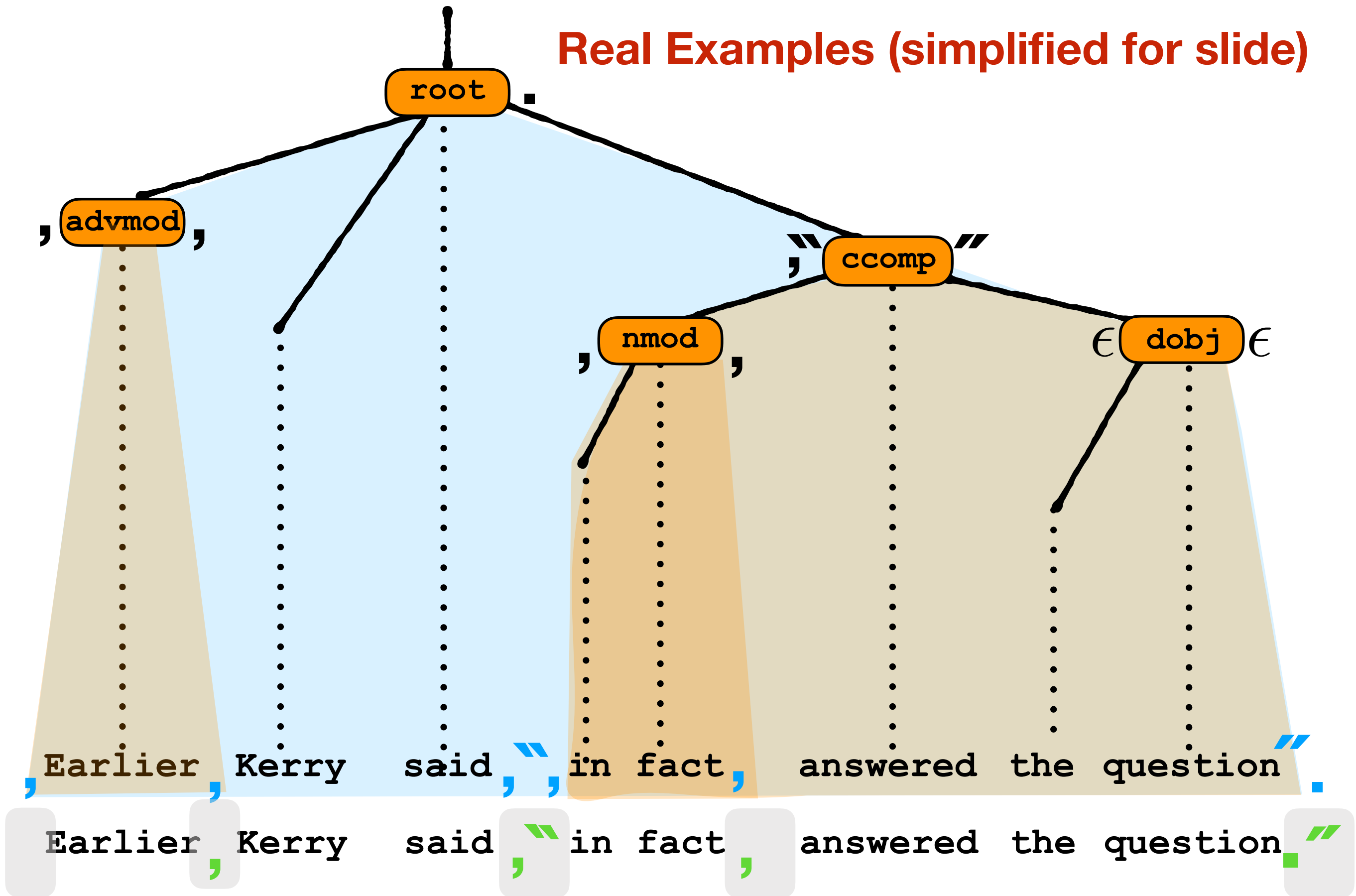
surface



Distribution over rewrite rules

# Viterbi recovers good underlying punctuation

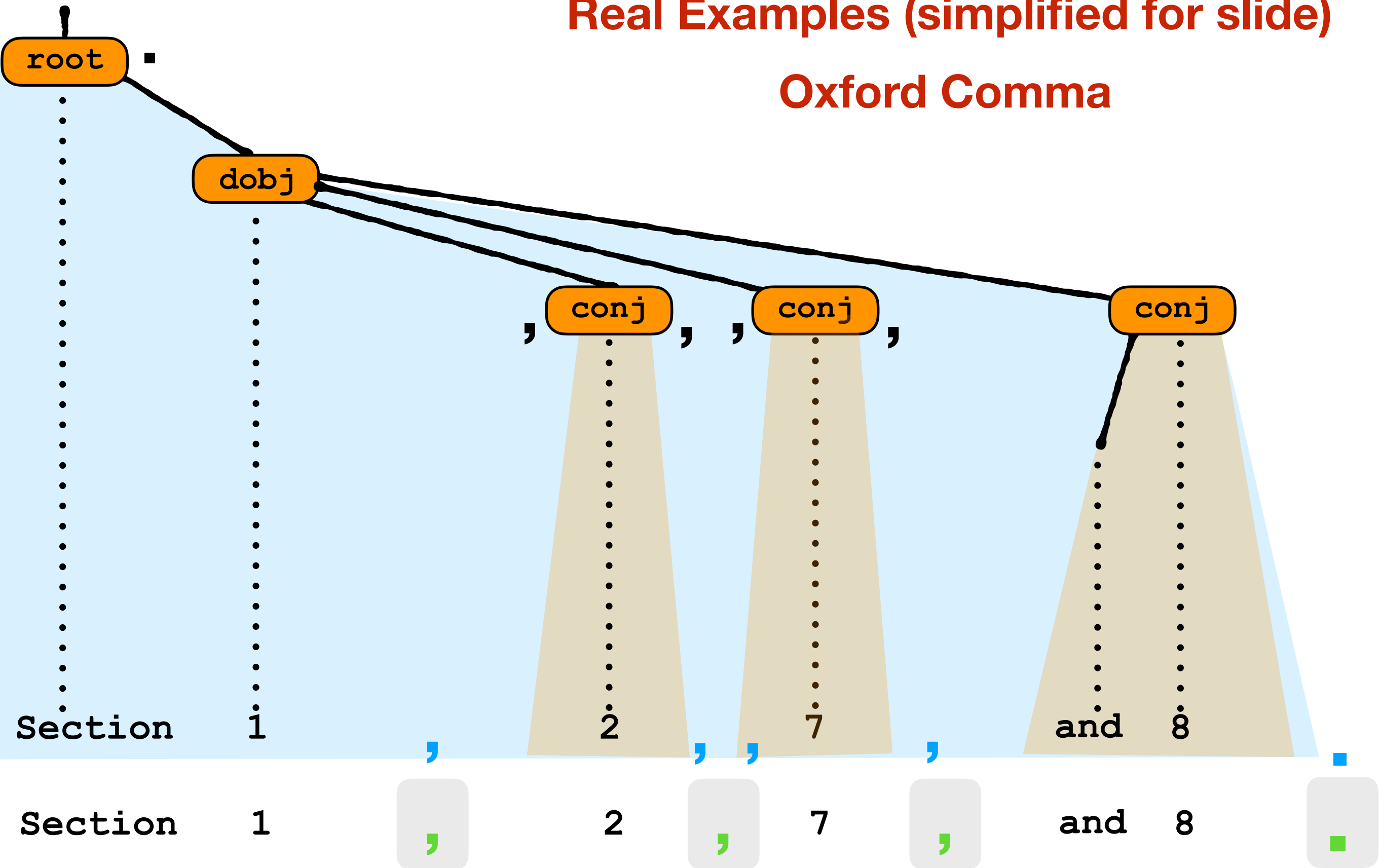
Real Examples (simplified for slide)

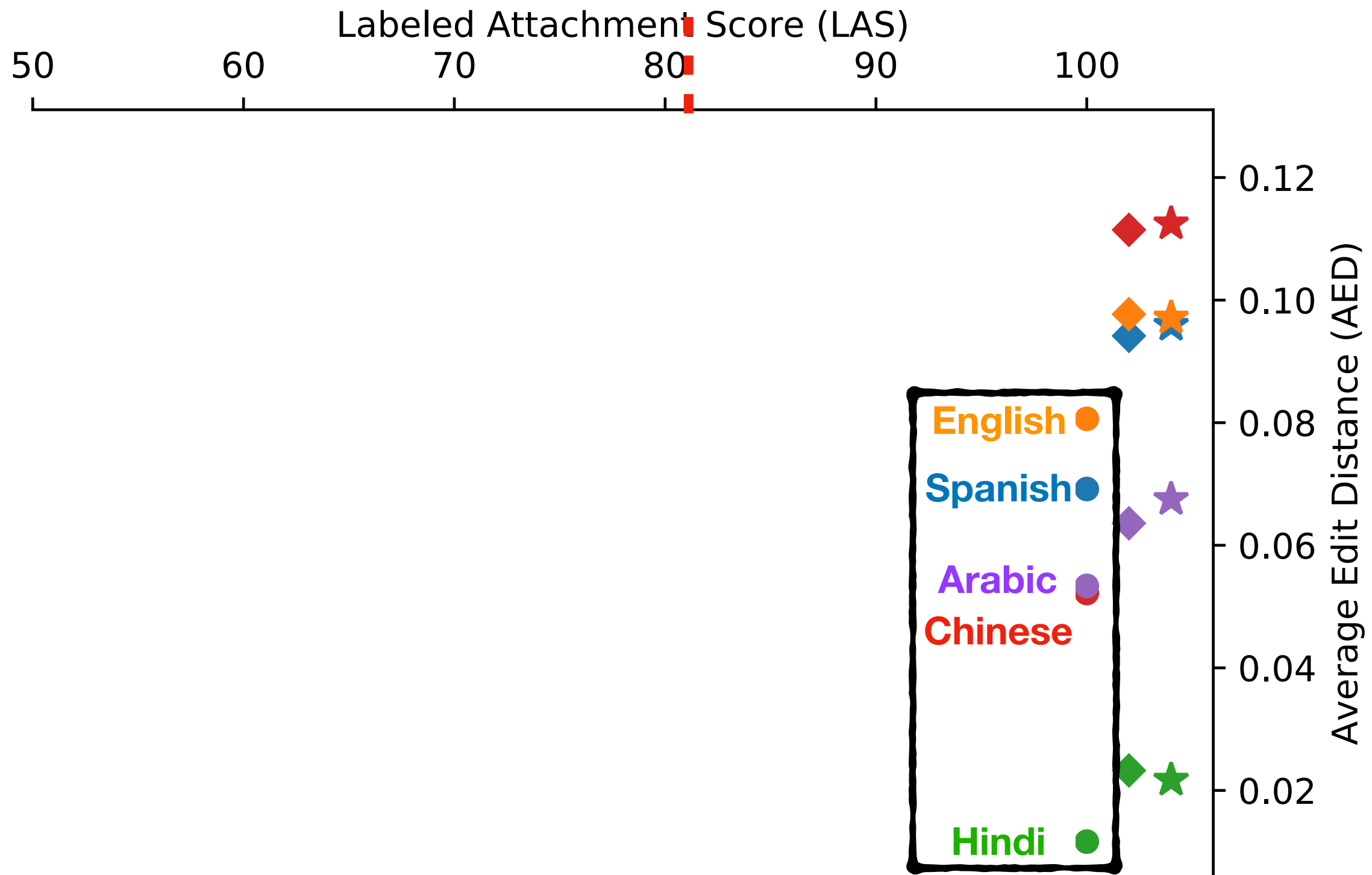


# Viterbi recovers good underlying punctuation

## Real Examples (simplified for slide)

# Oxford Comma





What if we don't have gold parses?

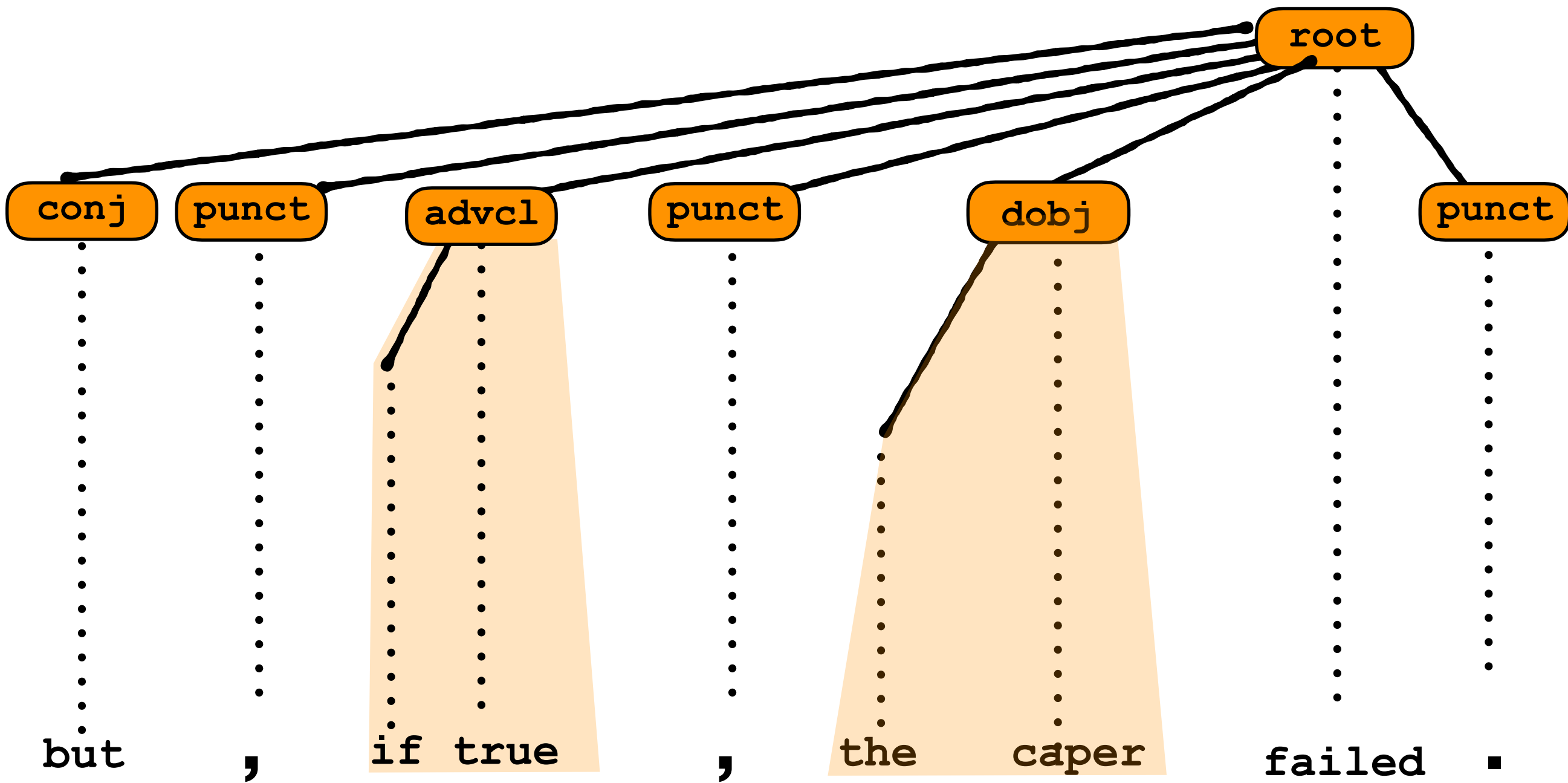
## Punctuation Restoration

average edit distance  
[the lower the better]

★ BiLSTM-CRF tagger

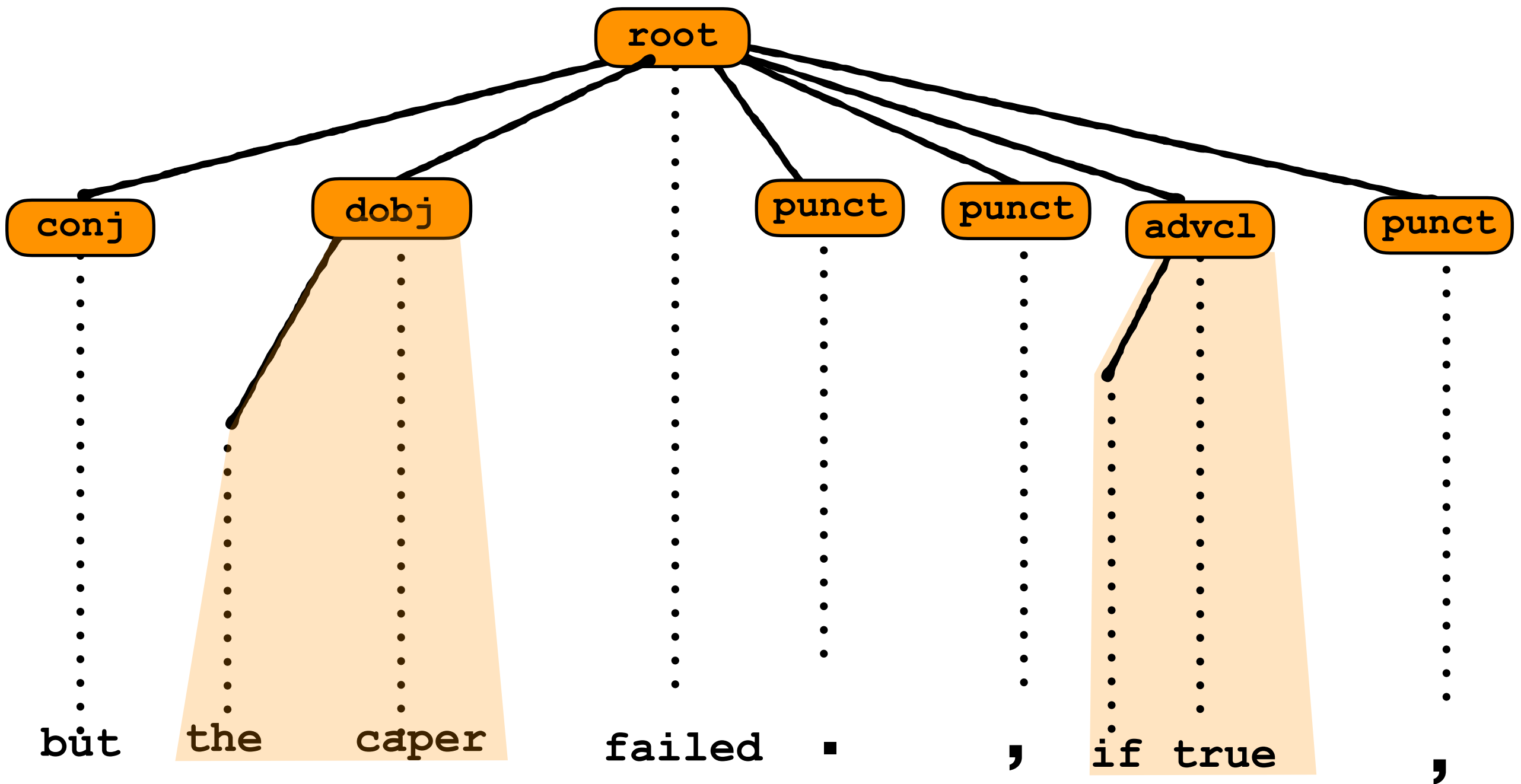
◆ Always Period tagger

● Ours



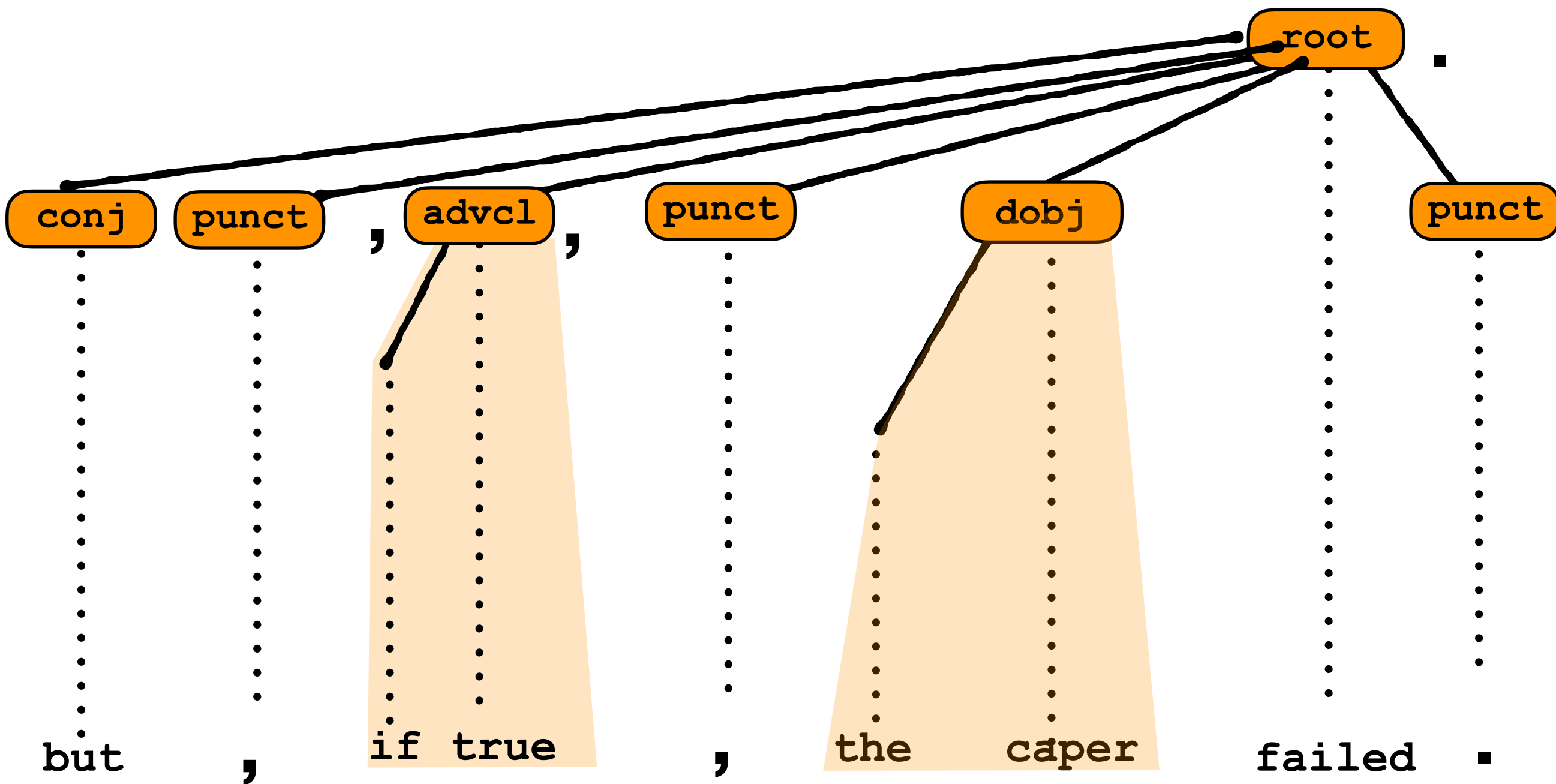
**Syntactically Transform Sentence from Treebank**



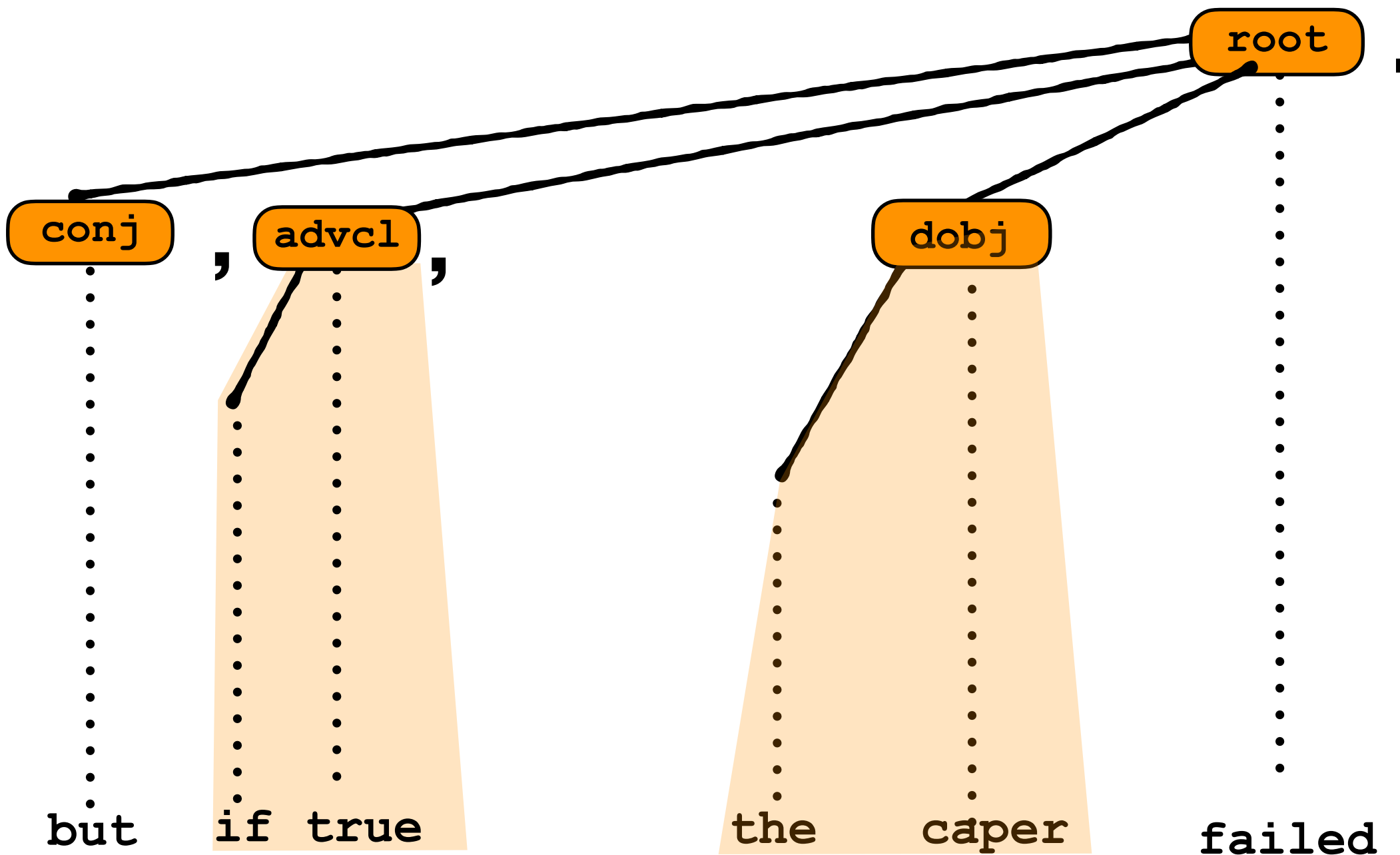


obviously ugly, and high perplexity under an LM

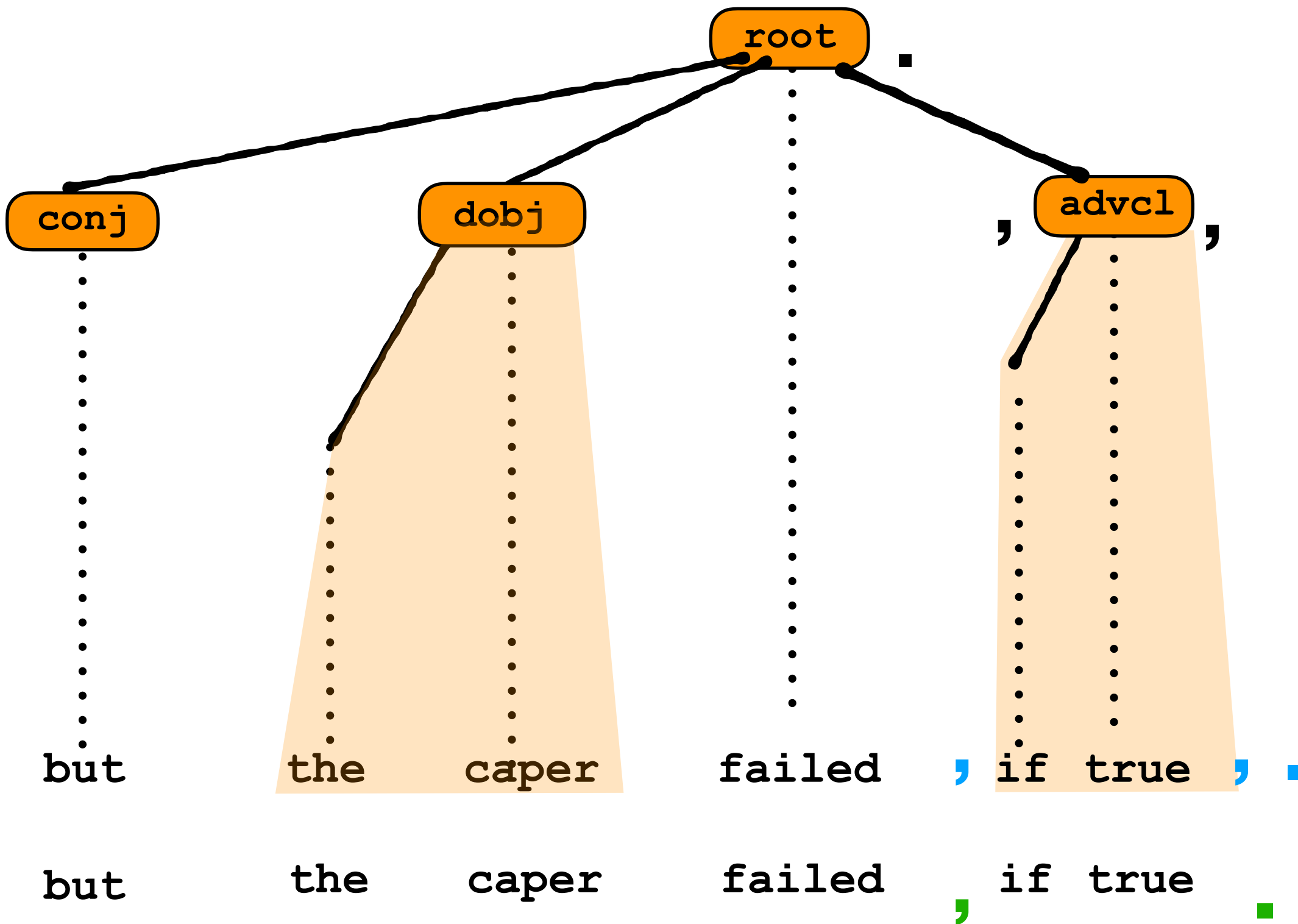
**Syntactically Transform Sentence from Treebank**



**Syntactically Transform Sentence with our Annotation**



**Syntactically Transform Sentence with our Annotation**



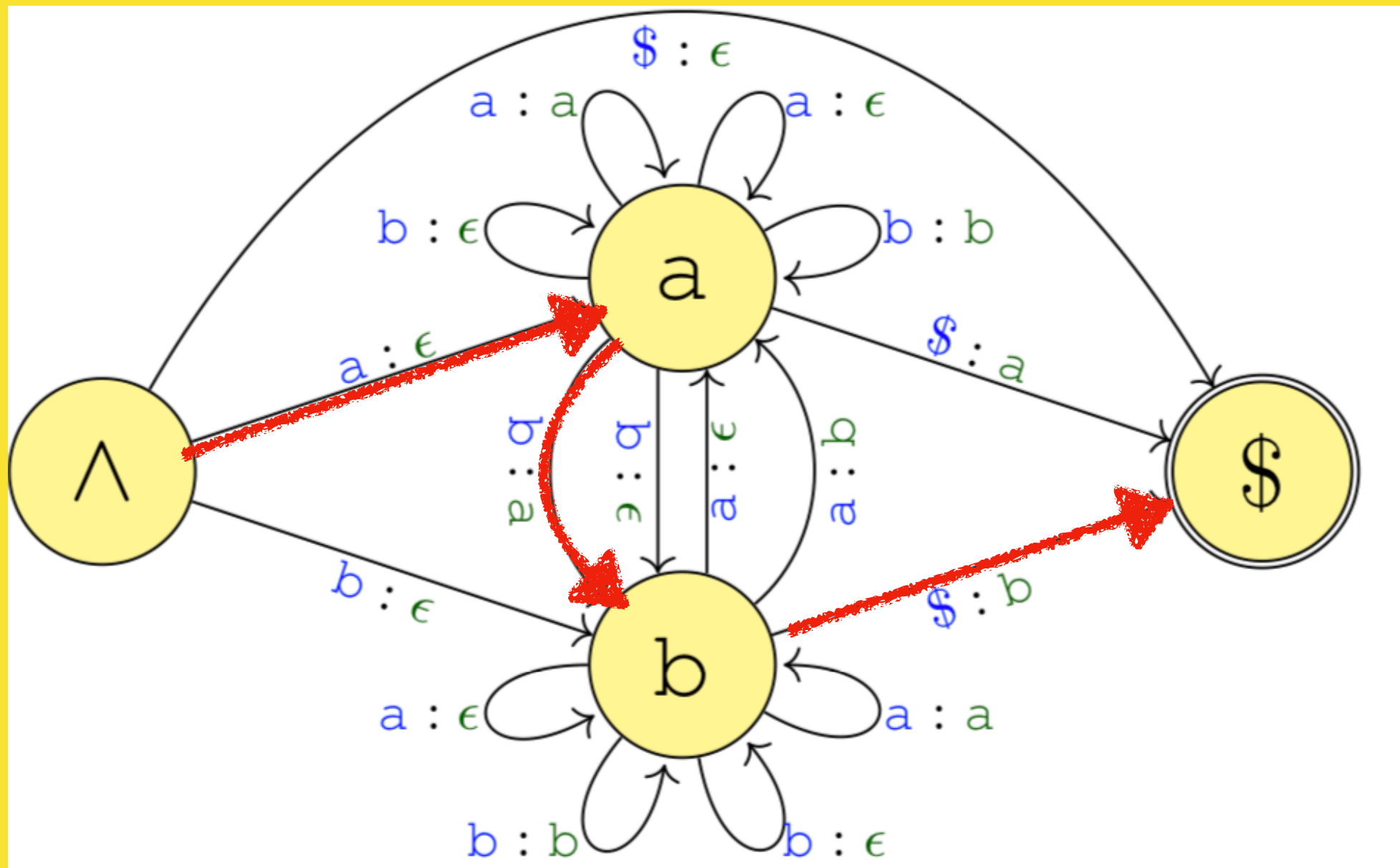
**Syntactically Transform Sentence with our Annotation**

Let's slow down  
for a second here.

We better just  
stop right now.

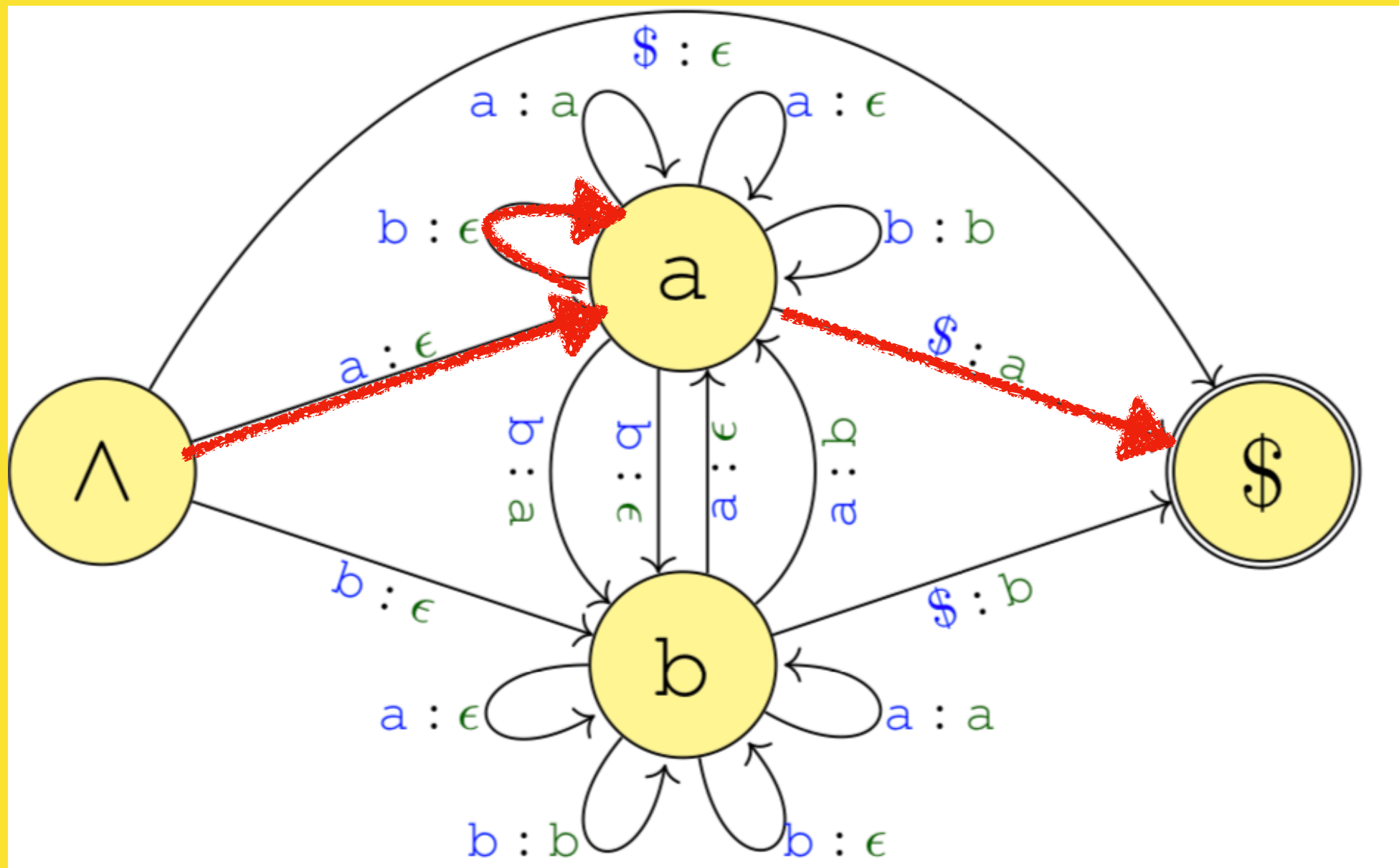


# Implemented by Finite State Transducer



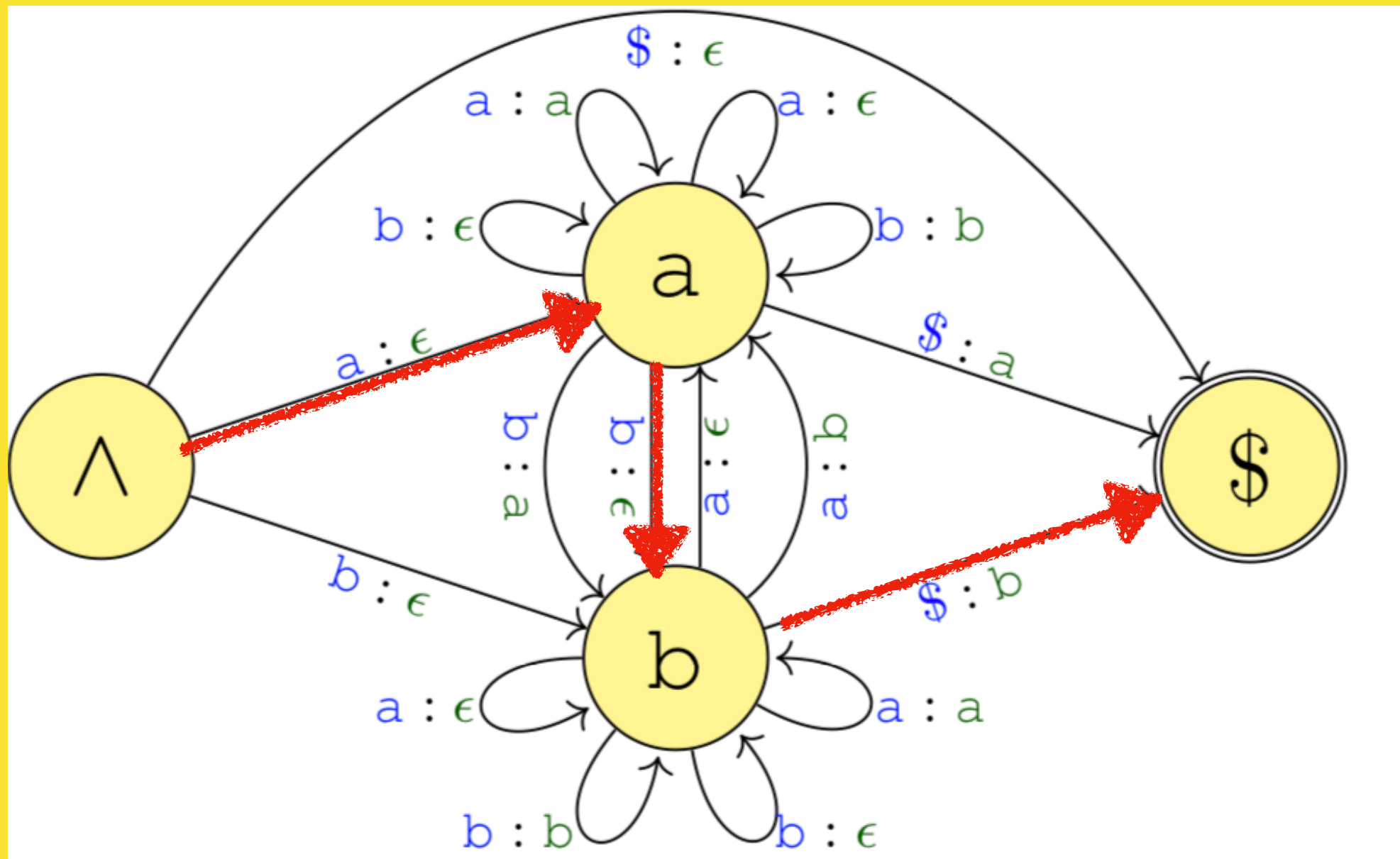
$ab \rightarrow ab$   
 $ab \rightarrow a$   
 $ab \rightarrow b$   
 $ab \rightarrow ba$

# Implemented by Finite State Transducer



$ab \rightarrow ab$   
 $ab \rightarrow a$   
 $ab \rightarrow b$   
 $ab \rightarrow ba$

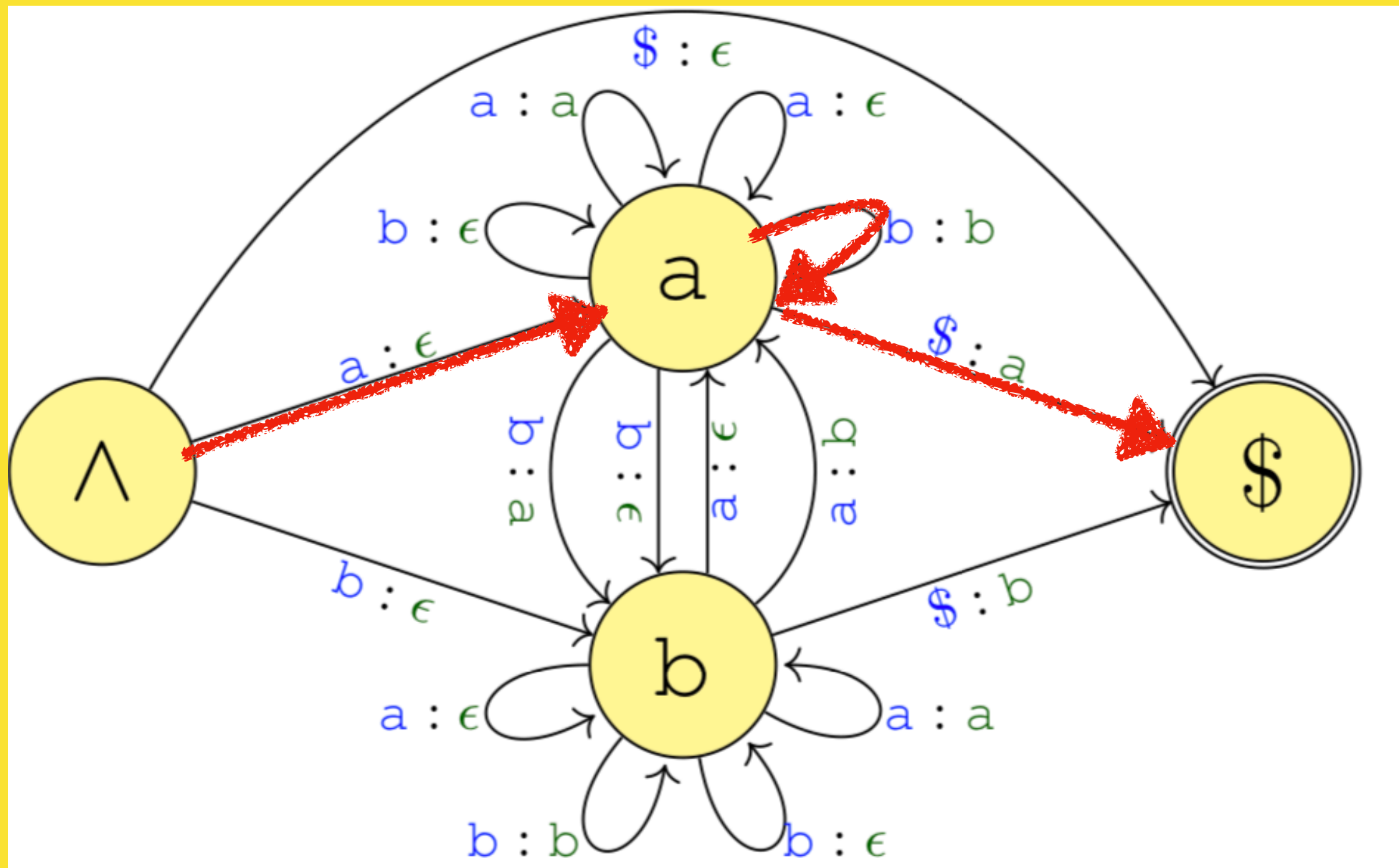




$ab \rightarrow ab$   
 $ab \rightarrow a$   
 $ab \rightarrow b$   
 $ab \rightarrow ba$



# Implemented by Finite State Transducer



$ab \rightarrow ab$   
 $ab \rightarrow a$   
 $ab \rightarrow b$   
 $ab \rightarrow ba$