

# Bilexical Grammars and a $O(n^3)$ Probabilistic Parser

Jason Eisner  
University of Pennsylvania  
IWPT - 1997

## Soft Selection

**doff** a **cap**  
**hat**  
sombrero  
shirt  
sink  
clothe  
about  
...



*Adjuncts too:*  
doffed his cap **to her**  
**at her**  
**for her**

monkeys doffing their hats

## Lexicalized Grammars

*doff*: \_\_\_ NP      *doff*: (S\NP)/NP

$S \rightarrow NP \text{ doff } NP$

*doff*:  $\begin{bmatrix} s \\ L(np) \\ R(np) \end{bmatrix}$

Diagram showing a tree structure for the sentence "The cat in the hat wore a striped stovepipe to our house today." with labels for subject (subj) and object (obj) and the verb "doff".

Rules are specialized for individual words  
(or are implicit in lexical entries)

Jason Eisner (U. Penn)

3

## From lexical to bilexical

- 1 Lafferty et al. 92, Charniak 95, Alshawi 96, Collins 96, Eisner 96, Goodman 97
- 1 Also see Magerman 94, Ratnaparkhi 97, etc.

### 1 Rules mention two words

E.g., each verb can have its *own* distribution of arguments

### 1 Goal: No parsing performance penalty

Alas, with standard chart parser:

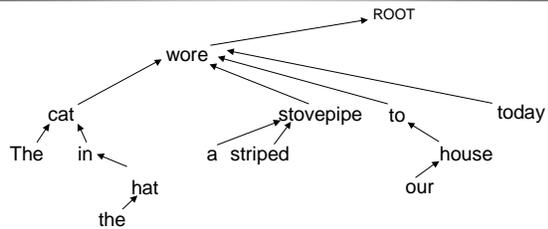
**nonlexical  $O(n^3)$**

**lexical  $O(n^3)$**

**bilexical  $O(n^2)$**

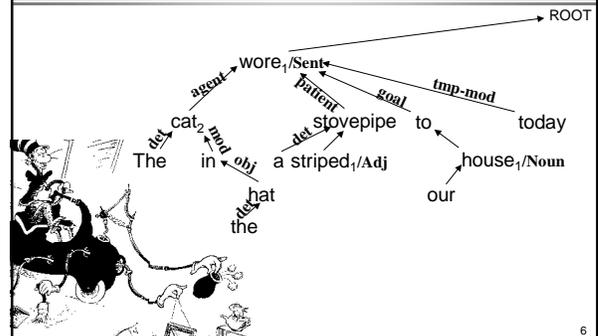
← other methods give  $O(n^4)$  or  $O(n^3)$

## Simplified Formalism (1)

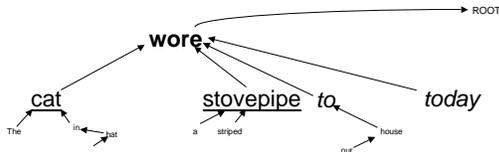


The cat in the hat wore a striped stovepipe to our house today.

(save these gewgaws for later)

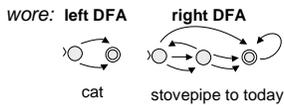


## Simplified Formalism (2)



Need a flexible mechanism to score the possible sequences of dependents.

every lexical entry lists 2 idiosyncratic DFAs. These accept dependent sequences the word likes.



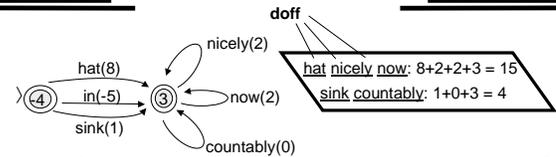
Jason Eisner (U. Penn)

7

## Weighting the Grammar

doff: right DFA Transitive verb.  
accepts: Noun Adv\*

likes: hat nicely now (e.g., "[Bentley] doffed [his hat] [nicely] [just now]")  
hates: sink countably (e.g., "#Bentley doffed [the sink] [countably]")



Jason Eisner (U. Penn)

8

## Why CKY is slow

1. visiting relatives is boring
2. visiting relatives wear funny hats
3. visiting relatives, we got bored and stole their funny hats

visiting relatives: NP(visiting), NP(relatives), AdvP, ...

CFG says that all NPs are interchangeable  
So we only have to use generic or best NP.

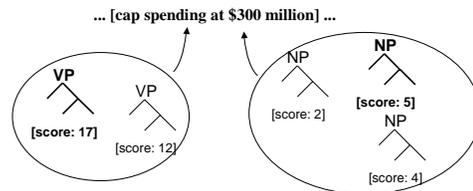
But billexical grammar disagrees:  
e.g., NP(visiting) is a poor subject for wear  
We must try combining each analysis w/ context

Jason Eisner (U. Penn)

9

## Generic Chart Parsing (1)

- 1 interchangeable analyses have same **signature**
- 1 "analysis" = tree or dotted tree or ...



- 1 if  $\leq S$  signatures, keep  $\leq S$  analyses per substring

Jason Eisner (U. Penn)

10

## Generic Chart Parsing (2)

for each of the  $O(n^2)$  substrings,  
for each of  $O(n)$  ways of splitting it,  
for each of  $\leq S$  analyses of first half  
for each of  $\leq S$  analyses of second half,  
for each of  $\leq c$  ways of combining them:  
**combine, & add result to chart if best**

$O(n^3 S^2 c)$

[cap spending] + [at \$300 million] = [[cap spending] [at \$300 million]]  
 $\leq S$  analyses       $\leq S$  analyses       $\leq c S^2$  analyses  
of which we keep  $\leq S$

Jason Eisner (U. Penn)

11

## Headed constituents ...

... have too many signatures.

**How bad is  $\Theta(n^3 S^2 c)$ ?**

For unheaded constituents, S is constant: NP, VP ...  
(similarly for dotted trees). So  $\Theta(n^3)$ .

But when different heads  $\Rightarrow$  different signatures,  
the average substring has  $\Theta(n)$  possible heads  
and  $S = \Theta(n)$  possible signatures. So  $\Theta(n^5)$ .

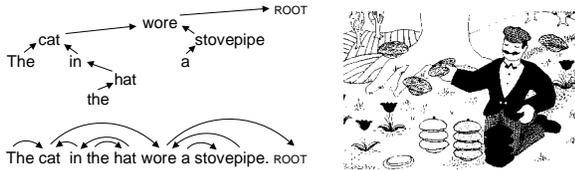
Jason Eisner (U. Penn)

12

## Forget heads - think hats!

Solution:

Don't assemble the parse from constituents.  
Assemble it from spans instead.



Jason Eisner (U. Penn)

13

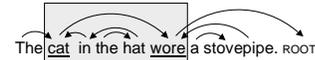
## Spans vs. constituents

Two kinds of substring.

» **Constituent** of the tree: links to the rest only through its head.



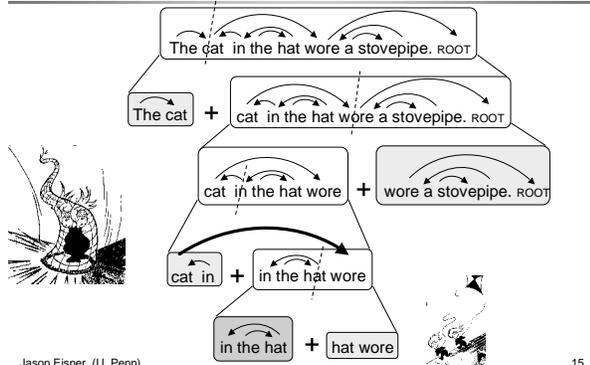
» **Span** of the tree: links to the rest only through its endwords.



Jason Eisner (U. Penn)

14

## Decomposing a tree into spans



Jason Eisner (U. Penn)

15

## Maintaining weights

Seed chart w/ word pairs  $\boxed{x \ y}$ ,  $\boxed{\overleftarrow{x} \ y}$ ,  $\boxed{x \ \overleftarrow{y}}$

Step of the algorithm:

$$\boxed{a \dots b} + \boxed{b \dots c} = \begin{cases} \boxed{a \dots b \dots c} \\ \boxed{a \dots b \dots c} \\ \boxed{a \dots b \dots c} \end{cases}$$

We can add an arc only if a, c are both parentless.

$$\text{weight}(\boxed{a \dots b \dots c}) = \text{weight}(\boxed{a \dots b}) + \text{weight}(\boxed{b \dots c}) + \text{weight of } c \text{ arc from } a\text{'s right DFA state} + \text{weights of stopping in } b\text{'s left and right states}$$

Jason Eisner (U. Penn)

16

## Analysis

Algorithm is  $O(n^3 S^2)$  time,  $O(n^2 S)$  space. What is S?

$$\boxed{a \dots b} + \boxed{b \dots c} = \boxed{a \dots b \dots c}$$

Where:

- » b gets a parent from exactly one side
- » Neither a nor c previously had a parent
- » a's right DFA accepts c; b's DFAs can halt

Signature of  $\boxed{a \dots b}$  has to specify parental status & DFA state of a and b

$\therefore S = O(t^2)$  where t is the maximum # states of any DFA

**S independent of n because all of a substring's analyses are headed in the same place - at the ends!**

Jason Eisner (U. Penn)

17

## Improvement

Can reduce S from  $O(t^2)$  to  $O(t)$

$$\boxed{a \dots b} + \boxed{b \dots c} = \boxed{a \dots b \dots c}$$

state of b's left automaton tells us weight of halting

likewise for b's right automaton

The halt-weight for each half is independent of the other half.

Add every span to both **left chart** & **right chart**

Above, we draw  $\boxed{a \dots b}$  from left chart,  $\boxed{b \dots c}$  from right chart

Copy of  $\boxed{a \dots b}$  in left chart has halt weight for b already added so its signature needn't mention the state of b's automaton

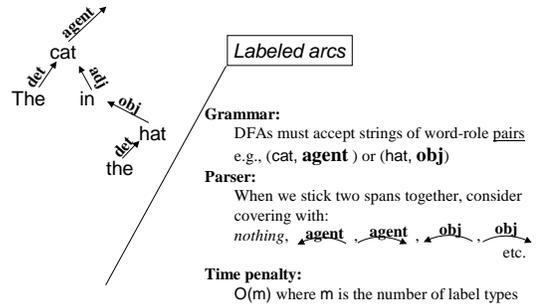
Jason Eisner (U. Penn)

18

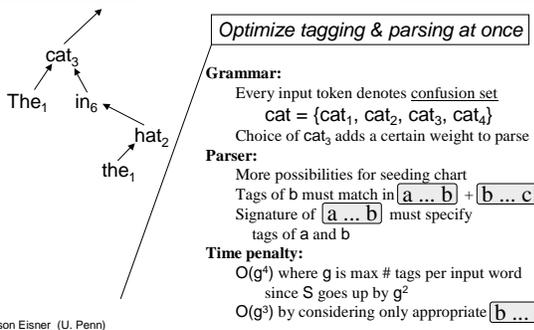
## Embellishments

- 1 More detailed parses
  - » Labeled edges
  - » Tags (part of speech, word sense, ...)
  - » Nonterminals
- 1 How to encode probability models

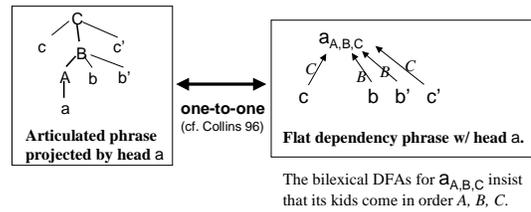
## More detailed parses (1)



## More detailed parses (2)

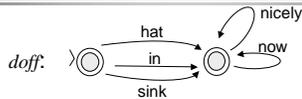


## Nonterminals



Use fast bilexical algorithm, then convert result to nonterminal tree.  
Want small (and finite) set of tags like  $a_{A,B,C}$ .  
(Guaranteed by X-bar theory:  $doff = \{doff_{V,VP}, doff_{V,VP,S}\}$ .)

## Using the weights



- 1 **Deterministic grammar:** All weights 0 or  $-\infty$
- 1 **Generative model:**  
 $\log \Pr(\text{next kid} = \text{nicely} \mid \text{doff in state 2})$
- 1 **Comprehension model:**  
 $\log \Pr(\text{next kid} = \text{nicely} \mid \text{doff in state 2, nicely present})$
- 1 **Eisner 1996** compared several models, found significant differences

## String-local constraints

Seed chart with word pairs like  $[x \ y]$   
We can choose to exclude some such pairs.

**Example:** k-gram tagging. (here  $k=3$ )

$[N \ P \ Det]$  tag with part-of-speech *trigrams*  
one cat in the hat weight =  $\log \Pr(\text{the} \mid \text{Det})\Pr(\text{Det} \mid N, P)$

$[Det \ V \ P]$   $[N \ P \ Det]$  excluded bigram:  
the the the 2 words disagree on tag for "cat"

## Conclusions

### 1 Bilexical grammar formalism

How much do 2 words want to relate?

Flexible: encode your favorite representation

Flexible: encode your favorite prob. model

### 1 Fast parsing algorithm

Assemble spans, not constituents

$O(n^3)$ , not  $O(n^5)$ . Precisely,  $O(n^3 t^2 g^3 m)$ .

$t$ =max DFA size,  $g$ =max senses/word,  $m$ =# label types

These grammar factors are typically small