# Lecture: Syntax Part 1

## CS 585, Fall 2016

Introduction to Natural Language Processing
http://people.cs.umass.edu/~brenocon/inlp2016
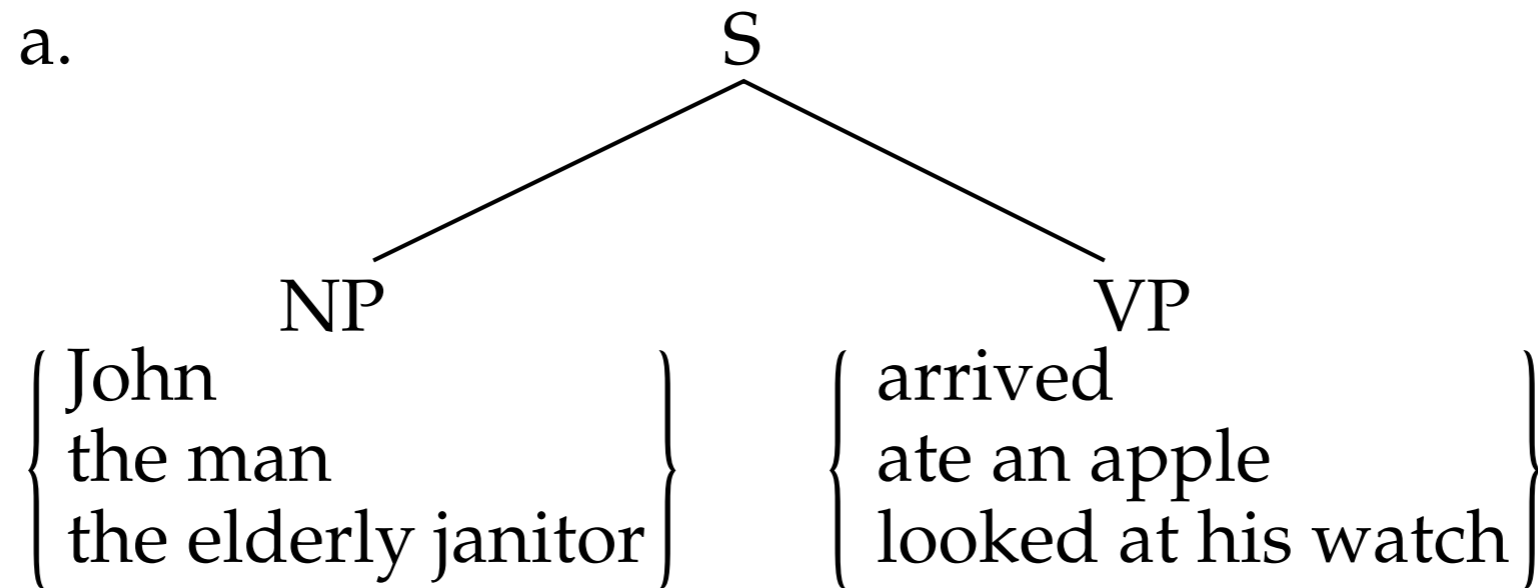
## Brendan O'Connor

College of Information and Computer Sciences
University of Massachusetts Amherst

# Midterm

- In-class next Thursday, 11/3. Review session on Tuesday 11/1.
- Closed book EXCEPT:
  One sheet of paper of any notes you want
  (front and back)
- Covers any material so far in the course; some subset of
  - Regular Expressions
  - Text normalization
  - Markov/N-gram Language Models
  - Naive Bayes
  - Classifiers
  - HMM
  - CRF
  - Perceptron
  - Parts of speech
  - Syntactic Parsing (this week)

- Syntax: how do words structurally combine to form sentences and meaning?

- Order
  - dogs chase cats  ..vs..  cats chase dogs
- Constituents
  - [the big dogs] chase cats
  - [colorless green clouds] chase cats
- Dependencies
  - The **dog chased** the cat.
  - My **dog**, a big old one, **chased** the cat.

- Idea of a *grammar*:  global template for how sentences / utterances / phrases are formed
  - Linguistics
  - Generation
  - Parsing (structured prediction)

- "a Sentence made of Noun Phrase followed by a Verb Phrase"

a.

$$S$$

NP

$\left\{\begin{array}{l}\text{John}\\\text{the man}\\\text{the elderly janitor}\end{array}\right\}$

VP

$\left\{\begin{array}{l}\text{arrived}\\\text{ate an apple}\\\text{looked at his watch}\end{array}\right\}$

b. S → NP VP                                                                 (1)

# Context-Free Grammar

- CFG describes a generative process for an (infinite) set of strings
  - 1. Nonterminal symbols
    - "S": START symbol / "Sentence" symbol
  - 2. Terminal symbols: word vocabulary
  - 3. Rules (a.k.a. Productions). Practically, two types:

"Grammar": *one* NT expands to >=1 NT
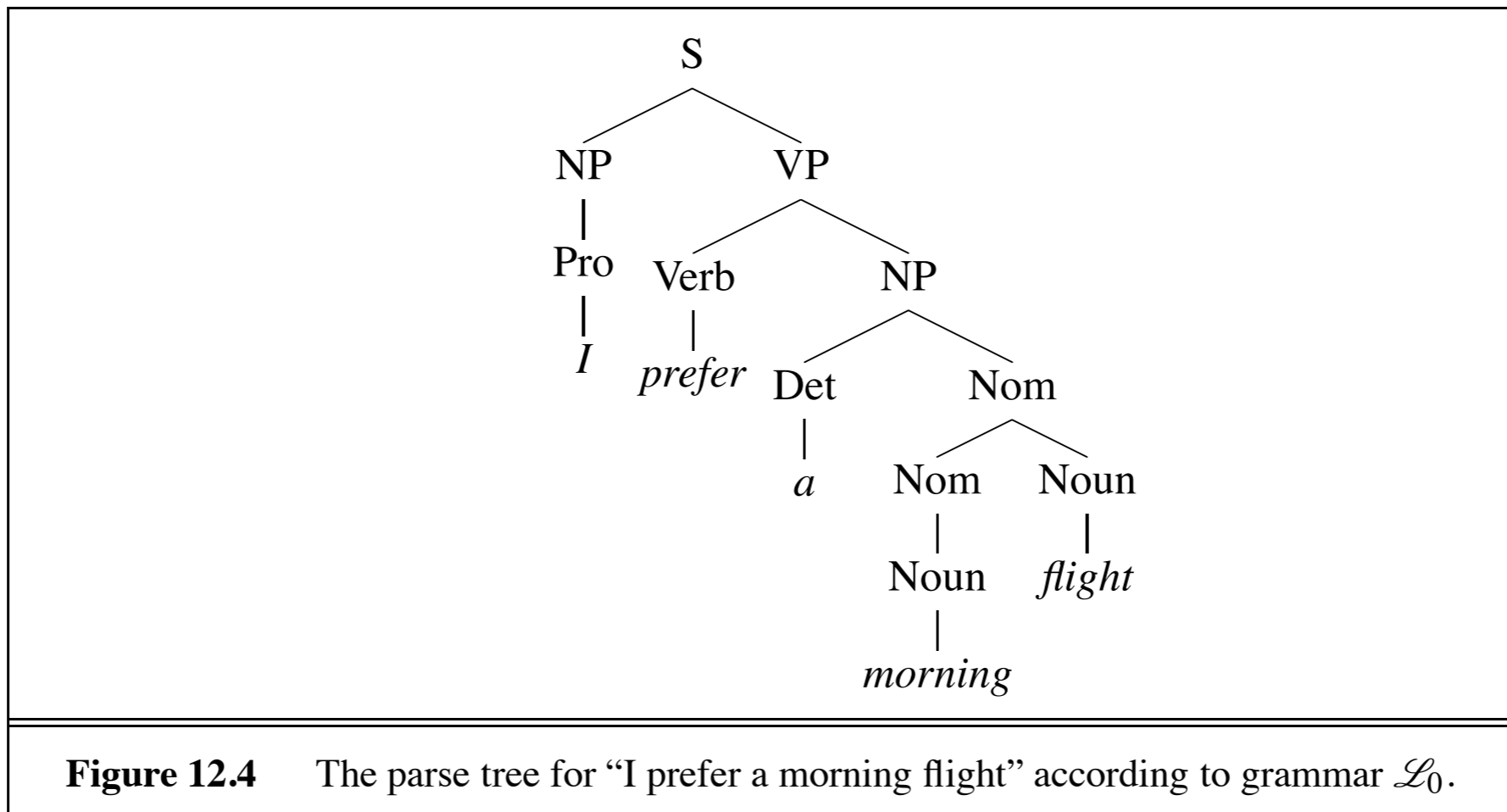always one NT on left side of rulep

Lexicon: NT expands to a terminal

| | | |
|---|---|---|
| $S \rightarrow$ | $NP\ VP$ | I + want a morning flight |
| $NP \rightarrow$ | $Pronoun$ | I |
| $\mid$ | $Proper\text{-}Noun$ | Los Angeles |
| $\mid$ | $Det\ Nominal$ | a + flight |
| $Nominal \rightarrow$ | $Nominal\ Noun$ | morning + flight |
| $\mid$ | $Noun$ | flights |
| $VP \rightarrow$ | $Verb$ | do |
| $\mid$ | $Verb\ NP$ | want + a flight |
| $\mid$ | $Verb\ NP\ PP$ | leave + Boston + in the morning |
| $\mid$ | $Verb\ PP$ | leaving + on Thursday |
| $PP \rightarrow$ | $Preposition\ NP$ | from + Los Angeles |

$$Noun \rightarrow flights \mid breeze \mid trip \mid morning \mid \ldots$$
$$Verb \rightarrow is \mid prefer \mid like \mid need \mid want \mid fly$$
$$Adjective \rightarrow cheapest \mid non-stop \mid first \mid latest$$
$$\mid other \mid direct \mid \ldots$$
$$Pronoun \rightarrow me \mid I \mid you \mid it \mid \ldots$$
$$Proper\text{-}Noun \rightarrow Alaska \mid Baltimore \mid Los\ Angeles$$
$$\mid Chicago \mid United \mid American \mid \ldots$$
$$Determiner \rightarrow the \mid a \mid an \mid this \mid these \mid that \mid \ldots$$
$$Preposition \rightarrow from \mid to \mid on \mid near \mid \ldots$$
$$Conjunction \rightarrow and \mid or \mid but \mid \ldots$$

*[only one token. ignore "L A"]*

5

# Constituent Parse Trees



**Figure 12.4**     The parse tree for "I prefer a morning flight" according to grammar $\mathscr{L}_0$.

Representations:

Bracket notation

(12.2)     $[_S [_{NP} [_{Pro} \text{I}]] [_{VP} [_V \text{prefer}] [_{NP} [_{Det} \text{a}] [_{Nom} [_N \text{morning}] [_{Nom} [_N \text{flight}]]]]]]$

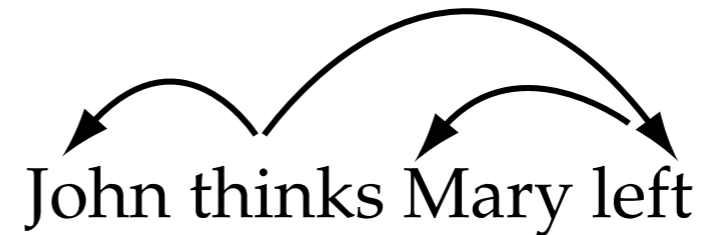Non-terminal positional spans

e.g. (NP, 0, 1), (VP, 1, 5), (NP, 2, 5), etc.

# Dependencies
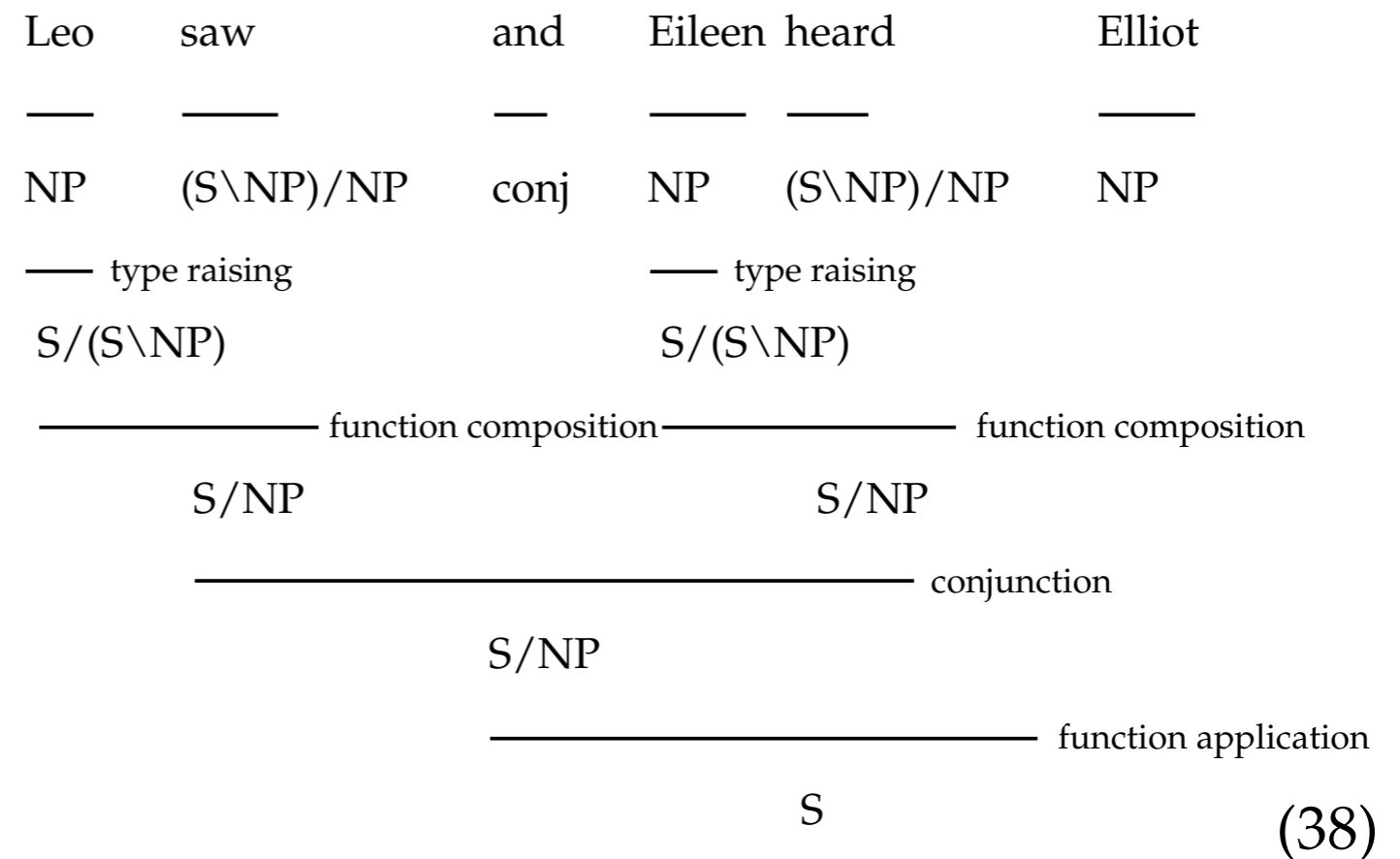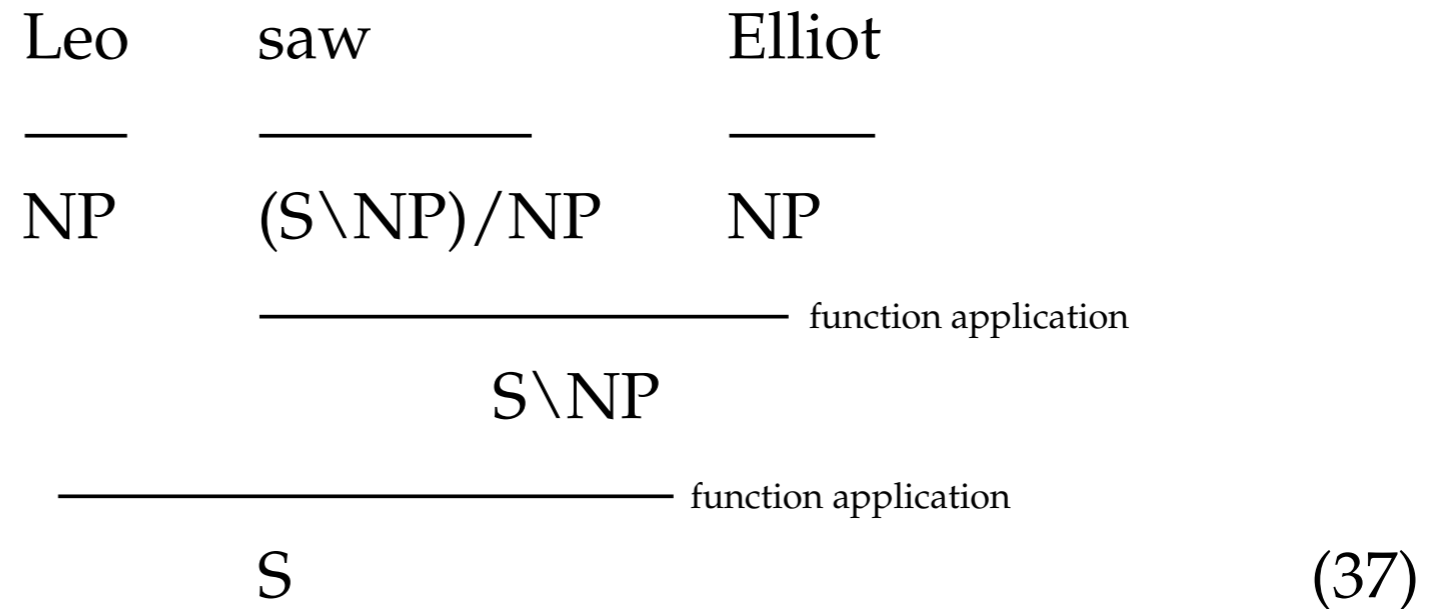
John saw Mary          John thinks Mary left

- Grammatical relationships between words in a sentence
- Most typically used type of parsing in modern NLP:
  e.g. "dependency bigrams" features
- e.g. http://nlp.stanford.edu:8080/corenlp/process
- Creeping toward semantics

# CCG

- Combinatory Categorial Grammar

- Syntactic theory based on constituent types that can combine to left and right

$$
\begin{array}{ccc}
\text{Leo} & \text{saw} & \text{Elliot} \\
\hline
\text{NP} & \text{(S\textbackslash NP)/NP} & \text{NP}
\end{array}
$$

Leo        saw              Elliot
——        ————            ——
NP        (S\NP)/NP        NP
                   ———————————— function application
                          S\NP
          ———————————————— function application
                   S                                (37)

Leo    saw              and     Eileen  heard        Elliot
——    ———            —      ——    ——          ——
NP    (S\NP)/NP        conj    NP     (S\NP)/NP    NP
—— type raising                      —— type raising
 S/(S\NP)                              S/(S\NP)
————————— function composition —————————— function composition
     S/NP                                   S/NP
     ———————————————————————— conjunction
                   S/NP
               ————————————————— function application
                          S                           (38)

# Modern parsing

- Structured prediction:  a tree structure
- Supervised training data: treebanks
  - Very labor intensive to create!
  - English: 1993 Penn Treebank / Wall Street Journal is still standard, more or less...
- Either constituents or dependencies
  - Dependencies are most common right now
  - Creeping towards semantics...

- Tomorrow: parsing with CFGs (most basic)