# Computational Linguistics



Today, we explore this interstice

# Computational Linguistics
## (aka Natural Language Processing)



Christopher Manning

Ling 1

November 4, 2011

# What is Computational Linguistics?

- Getting computers to perform useful tasks involving human languages whether for:
  - Enabling human-machine communication
  - Improving human-human communication
  - Doing stuff with language data … email, blogs, etc.

- Examples:
  - Machine Translation
  - Automatic Question Answering
  - Speech Recognition
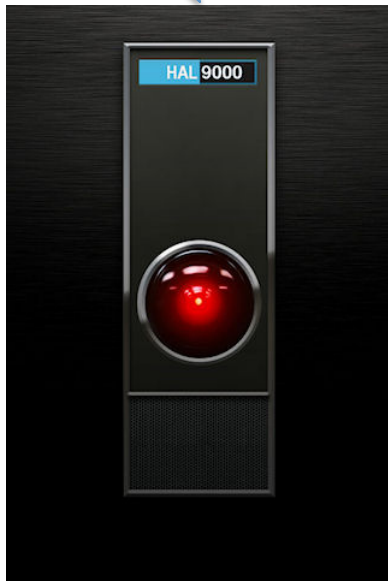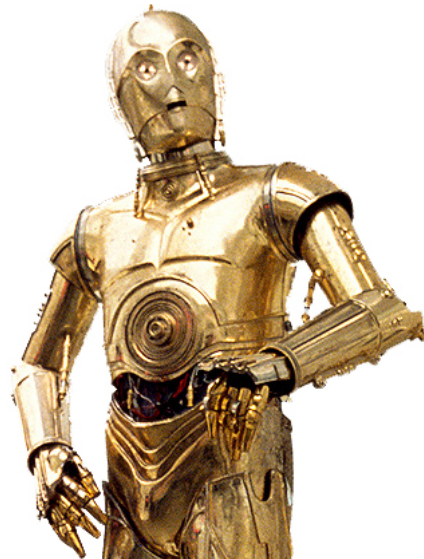  - Text-to-Speech Synthesis
  - Text Understanding

# CL vs. NLP

- Why say "Computational Linguistics (CL)" versus "Natural Language Processing" (NLP)?

- Either choose either freely or …

- Computational Linguistics
  - The science of computers dealing with language
  - Some interest in modeling what people do

- Natural Language Processing
  - Developing computer systems for processing and understanding human language text

# The Vision

# Language: the ultimate UI



Where is A Bug's Life playing in Mountain View?

A Bug's Life is playing at the Century 16 Theater.

When is it playing there?

It's playing at 2pm, 5pm, and 8pm.

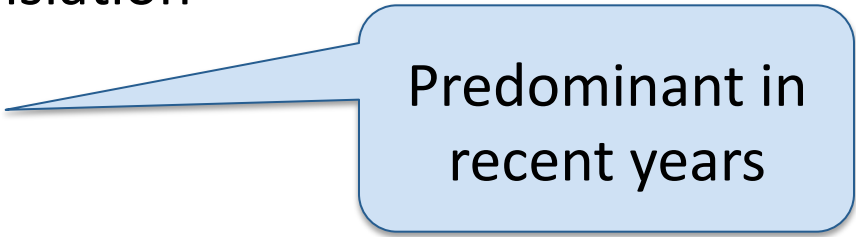OK. I'd like 1 adult and 2 children for the first show. How much would that cost?

But we need domain knowledge, discourse knowledge, world knowledge
(Not to mention linguistic knowledge!)

# NLP: Goals of the field

- ## From the lofty …

  - full-on natural language understanding
  - participation in spoken dialogues
  - open-domain question answering
  - real-time bi-directional translation

- ## … to the mundane

  - identifying spam
  - categorizing news stories
  - finding & comparing product information on the web
  - assessing sentiment toward products, brands, stocks, …

Predominant in recent years

# NLP in the commercial world

# Current motivations for NLP

## What's driving NLP?  Three trends:

- The explosion of machine-readable natural language text
  - Exabytes ($10^{18}$ bytes) of text, doubling every year or two
  - Web pages, emails, IMs, SMSs, tweets, docs, PDFs, …
  - Opportunity — and increasing necessity — to extract meaning

- Mediation of human interactions by computers
  - Opportunity for the computer in the loop to do much more

- Growing role of language in human-computer interaction

# Further motivation for CL

*One reason for studying language — and for me personally the most compelling reason — is that it is tempting to regard language, in the traditional phrase, as a "mirror of mind".*
Chomsky, 1975

For the same reason, computational linguistics is a compelling way to study human language acquisition and processing.

Sometimes, the best way to understand something is to build a model of it.

*What I cannot create, I do not understand.*  Feynman, 1988

# Subfields and tasks

| mostly solved | making good progress | still really hard |
|---|---|---|
| **Spam detection**<br>OK, let's meet by the big … ✔<br>D1ck too small? Buy V1AGRA … ✘ | **Sentiment analysis**<br>The pho was authentic and yummy. 👍<br>Waiter ignored us for 20 minutes. 👎 | **Semantic search**<br>people protesting globalization [Search]<br>➡ …demonstrators stormed IMF offices… |
| **Text categorization**<br>Phillies shut down Rangers 2-0  SPORTS<br>Jobless rate hits two-year low  BUSINESS | **Coreference resolution**<br>Obama told Mubarak he shouldn't run again. | **Question answering (QA)**<br>Q. What currency is used in China?<br>A. The yuan |
| **Part-of-speech (POS) tagging**<br>ADJ   ADJ  NOUN VERB   ADV<br>Colorless  green  ideas  sleep  furiously. | **Word sense disambiguation (WSD)**<br>I need new batteries for my *mouse*. | **Textual inference & paraphrase**<br>T. Thirteen soldiers lost their lives …<br>H. Several troops were killed in the …   YES |
| **Named entity recognition (NER)**<br>PERSON        ORG        LOC<br>Obama met with UAW leaders in Detroit … | **Syntactic parsing**<br>I can see Russia from my house! | **Summarization**<br>Sheen continues rant against … ➡ Sheen is nuts |
| **Information extraction (IE)**<br>You're invited to our bunga bunga party, Friday May 27 at 8:30pm in Cordura Hall<br>Party May 27 add | **Machine translation (MT)**<br>Our specialty is panda fried rice. ➡<br>我们的专长是熊猫炒饭 | **Discourse & dialog**<br>Where is Thor playing in SF?<br>Metreon at 4:30 and 7:30 |

# Why is computational linguistics hard?

Human languages:

- are highly ambiguous at all levels
- are complex, with recursive structures and coreference
- subtly exploit context to convey meaning
- are fuzzy and vague
- require reasoning about the world for understanding
- are part of a social system: persuading, insulting, amusing, …

(Nevertheless, simple features often do half the job!)

# OK, why *else* is NLP hard?

Oh so many reasons!

**non-standard English**

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either♥

**segmentation issues**

the New York-New Haven Railroad
the New York-New Haven Railroad

**idioms**

dark horse
get cold feet
lose face
throw in the towel

**neologisms**

unfriend
retweet
bromance
teabagger

**garden path sentences**

The man who hunts ducks out on weekends.
The cotton shirts are made from grows here.

**tricky entity names**

… a mutation on the *for* gene …
Where is *A Bug's Life* playing …
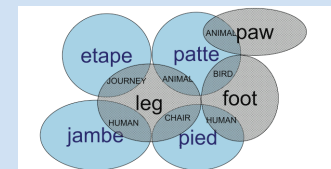Most of *Let It Be* was recorded …

**world knowledge**

Mary and Sue are sisters.
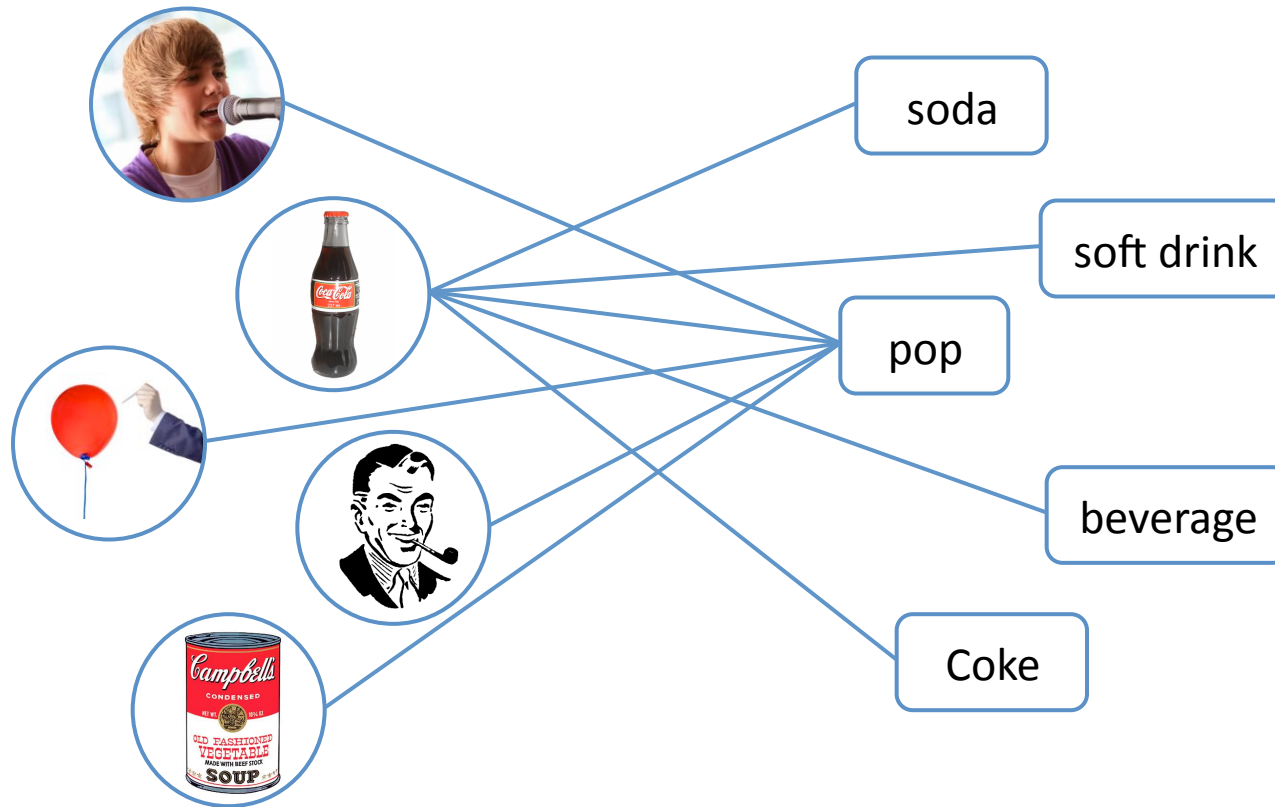Mary and Sue are mothers.

**prosody**

I never said *she* stole my money.
I never said she *stole* my money.
I never said she stole *my* money.

**lexical specificity**



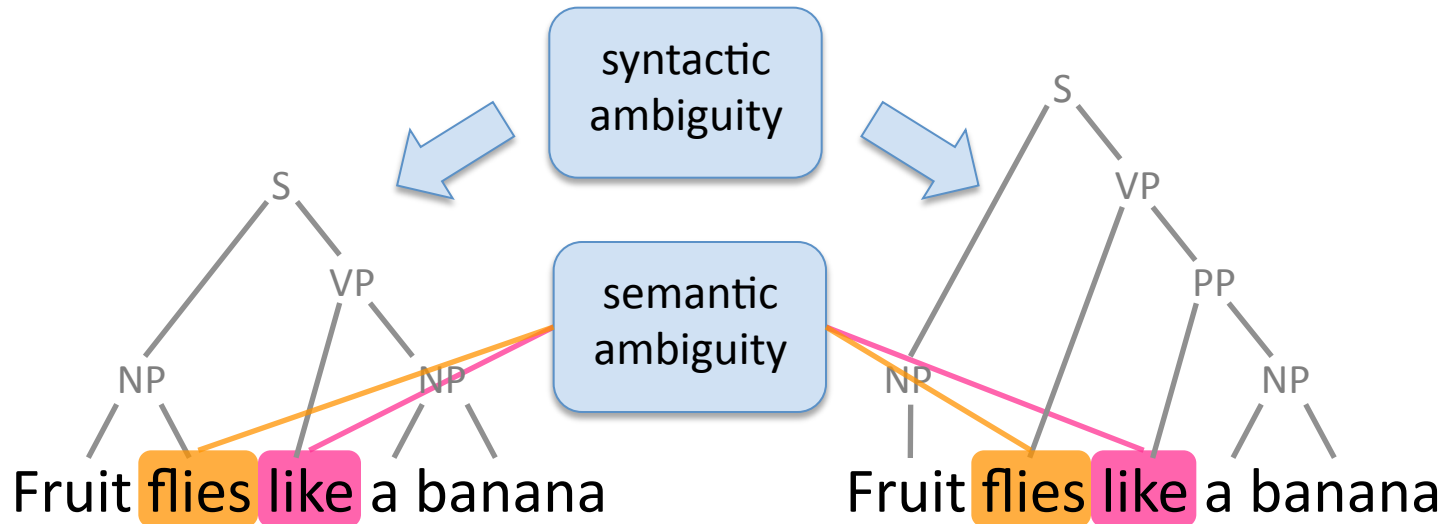But that's what makes it fun!

# Meanings and expressions

# One meaning, many expressions

Consider a semantic search application:

| Russia increasing price of gas for Georgia | Search |
|---|---|

Russia hits Georgia with huge rise in its gas bill

Russia plans to double Georgian gas price

Russia gas monopoly to double price of gas

Gazprom confirms two-fold increase in gas price for Georgia

Russia doubles gas bill to "punish" neighbour Georgia

Gazprom doubles Georgia's gas bill

# One expression, many meanings
# Syntactic & semantic ambiguity



syntactic ambiguity

semantic ambiguity

S
VP
NP
NP

Fruit flies like a banana

S
VP
PP
NP
NP

Fruit flies like a banana

# Ambiguous headlines

**Minister Accused Of Having 8 Wives In Jail**

May 21, 2007 06:49 AM



**ATLANTA (AP) --** A tra...
served two years in pri...
has been jailed again fo...
marry more women.

Bishop Anthony Owens,...
Ga., is in a Gwinnett Co...
four women claimed he...
after being released fro...

100% REAL

Teacher Strikes Idle Kids

China to Orbit Human on Oct. 15

Juvenile Court to Try Shooting Defendant
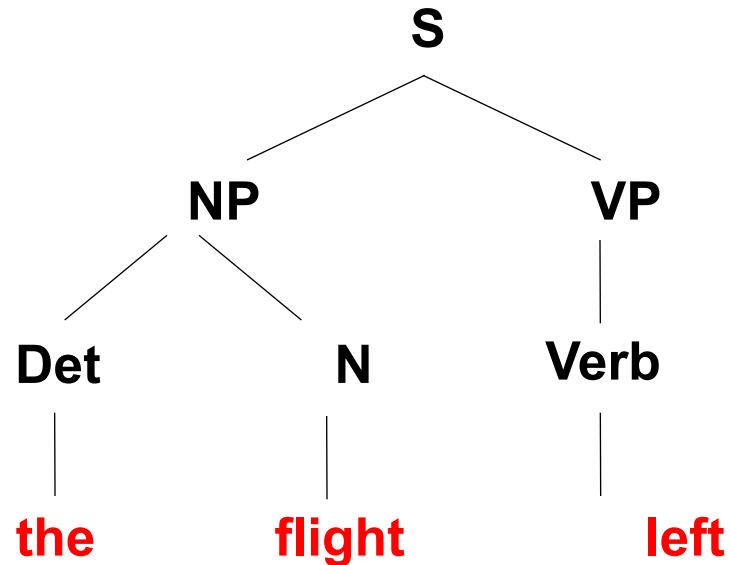
Clinton Wins on Budget, but More Lies Ahead

Local High School Dropouts Cut in Half

Police: Crack Found in Man's Buttocks

# Parsing

- Parsing is the process of taking a string and a grammar and returning a parse tree or trees for that string
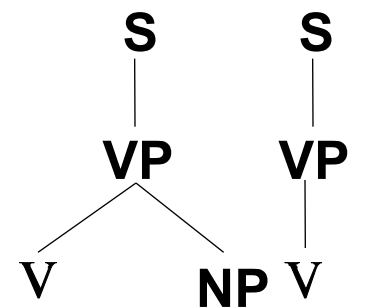
the flight left

```
            S
          /   \
        NP     VP
       /  \     |
     Det   N   Verb
      |    |    |
     the flight left
```
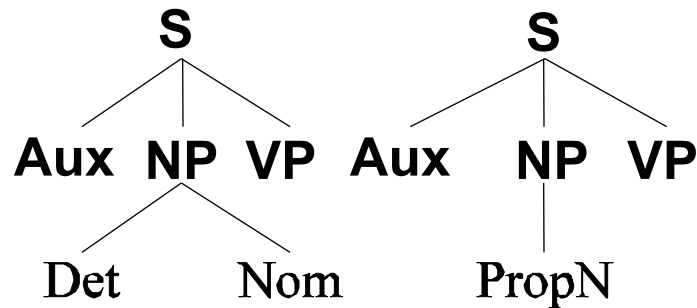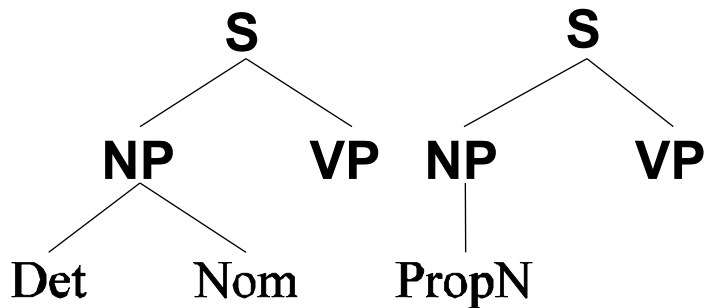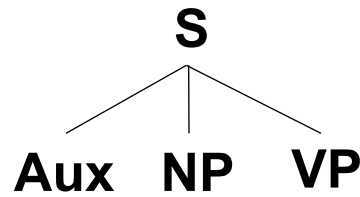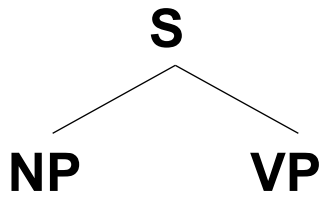
# Parsing involves search

- As with most everything of interest, parsing involves a search which involves the making of choices

- We'll look at some basic methods to give you an idea of the problem

# Top-Down Parsing

- Since we're trying to find trees rooted with an S (Sentence) start with the rules that give us an S.

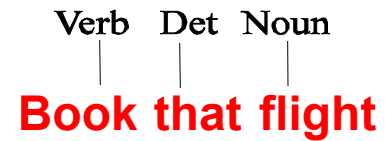- Then work your way down from there to the words.
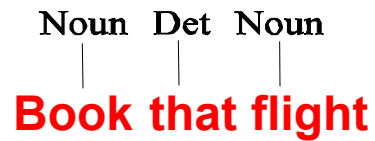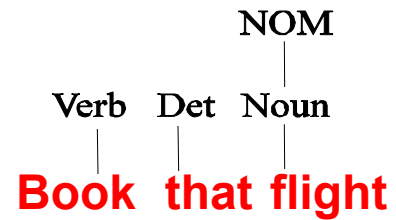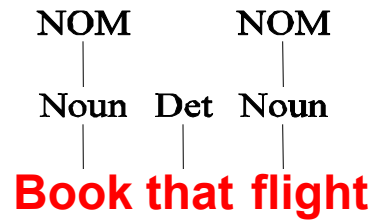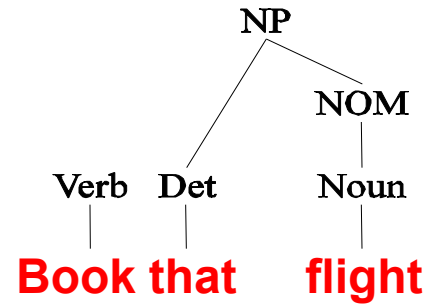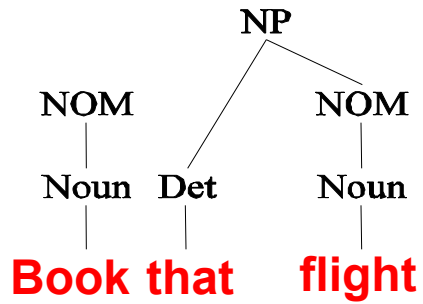
# Top-Down Space

# Bottom-Up Parsing

- Of course, we also want trees that cover the input words. So start with trees that link up with the words in the right way.

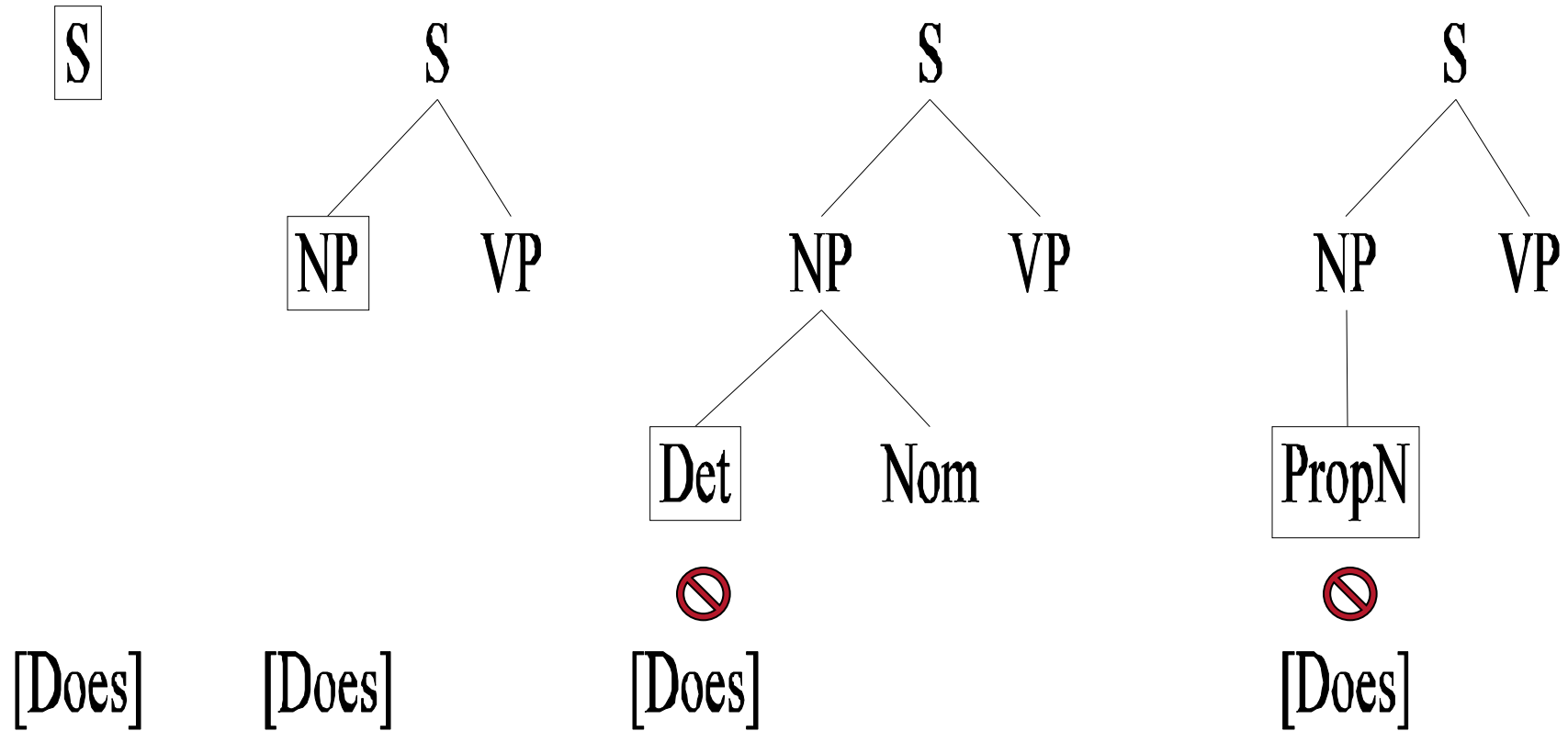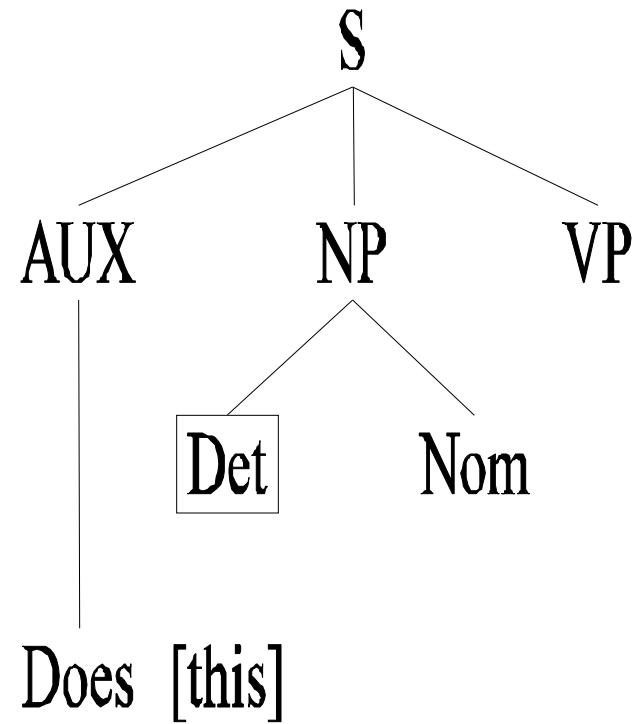- Then work your way up from there.
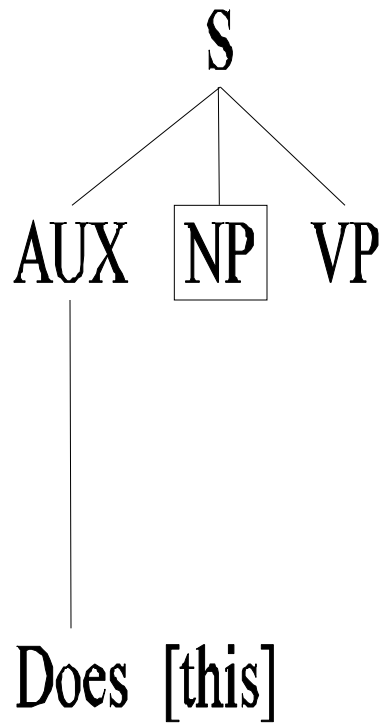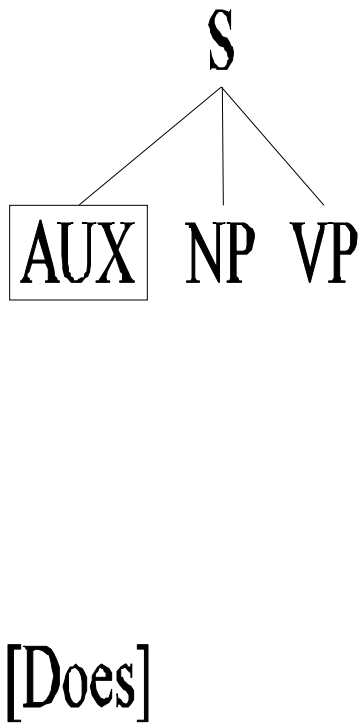
# Bottom-Up Space

# Control

- We need to keep track of the search space and have a strategy to make choices

  - We need to systematically explore everything to make sure we find the right parse for a sentence


  - Which node to try to expand next?
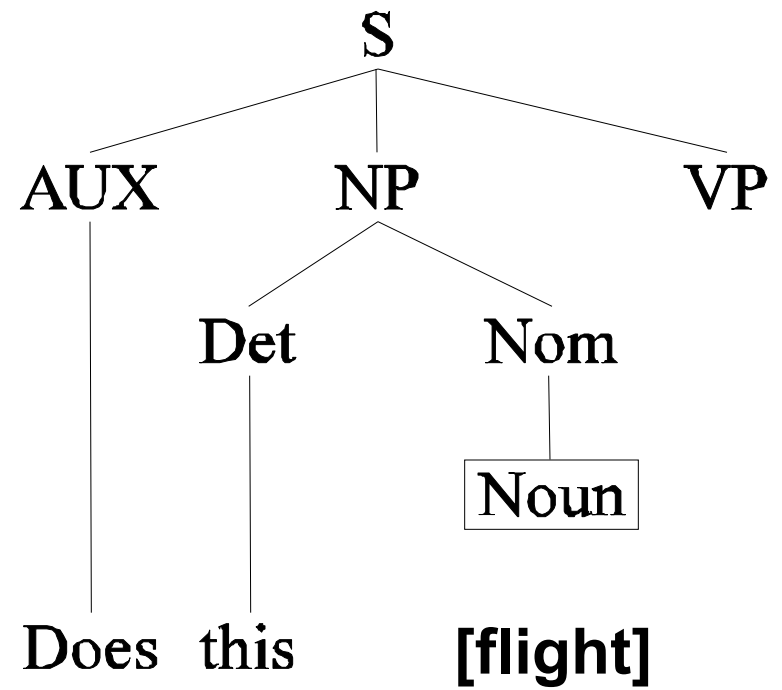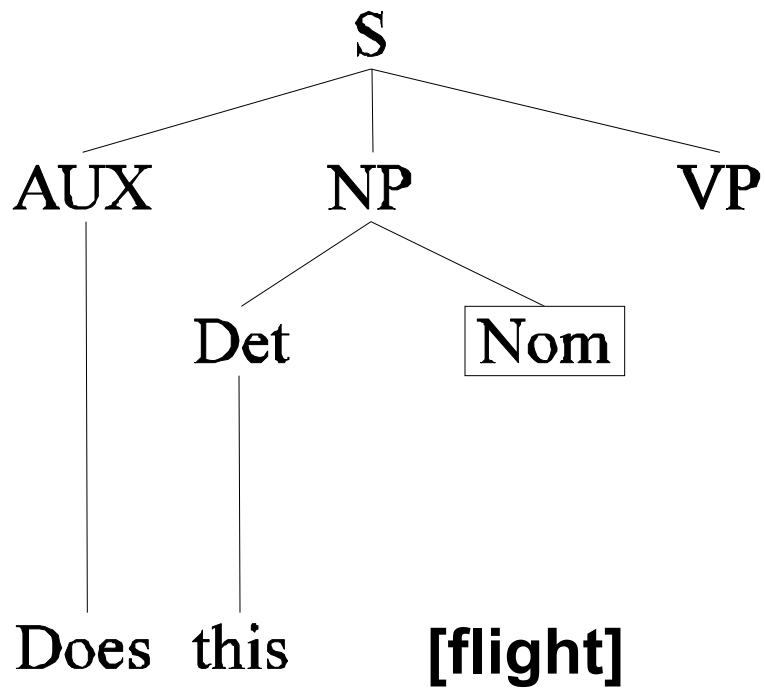  - Which grammar rule to use to expand a node?

# Top-Down, Depth-First, Left-to-Right Search

S

S
NP    VP

S
NP    VP
Det    Nom
🚫

S
NP    VP
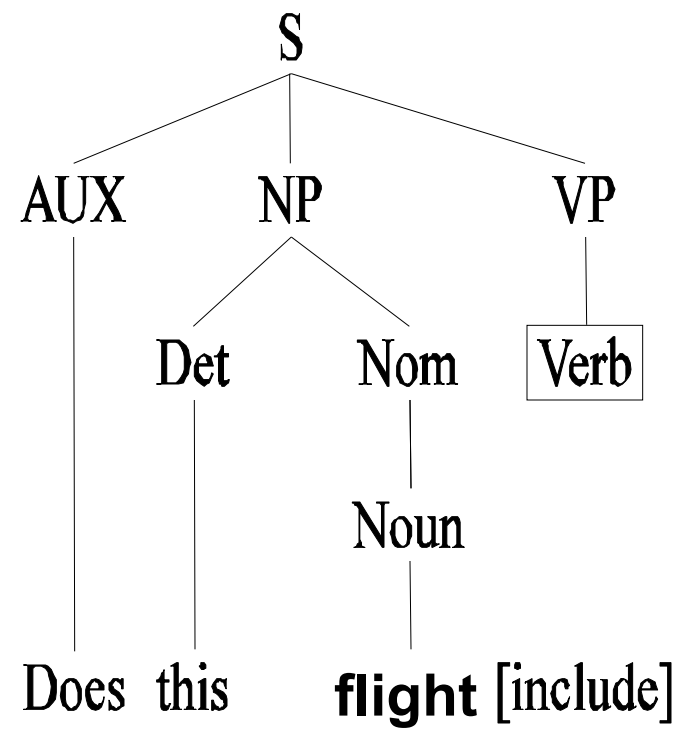PropN
🚫

[Does]    [Does]    [Does]    [Does]

# Example

# Example

# Example

# Efficient parsing

- That should give you the general idea of how a parser can work by exploring hypotheses systematically

- But really we need to do much more to make parsing efficient … this leads into dynamic programming, memoization, and other tricky stuff that I won't mention further here.

# How to choose between parses?

- Probabilistic methods!

- Augment the grammar rules with probabilities

- Modify the parser to keep only most probable parses

- At the end, return the most probable parse

# A statistical scientific revolution

- Computational Linguistics before 1990:
  - Hand-built parsers, hand-built dialogue systems
  - High precision, low coverage methods


- Computational Linguistics after 1995:
  - Automatically trained parsers, unsupervised clustering, statistical machine translation
  - High coverage, low precision methods
  - Build models exploiting **data**

# Demos!

# If you might like NLP / CompLing …

- learn Java or Python (and play with JavaNLP or NLTK)
- get some exposure to linguistics (LING1, …)
- and to logic, probability, statistics, linear algebra
- study AI and machine learning (CS121, CS221, CS229)
- read Jurafsky & Martin or Manning & Schütze
- Take
  - Ling 180/CS124: From Languages to Information
  - Ling 284/CS224N: Natural Language Processing
  - Ling 281/CS224S: Speech Recognition & Synthesis
  - Ling /CS224U: Natural Language Processing

# One more for the road