

## Course Overview

### Lecture #1

## Introduction to Natural Language Processing

CMPSCI 585, Fall 2004

University of Massachusetts Amherst



Andrew McCallum

Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein

## Today's Main Points

- Why is NLP interesting and difficult, complex and ambiguous.
  - Why? How to humans resolve this ambiguity?
- The six “layers” of NLP.
- NLP history, an overview, current successes.
- Course mechanics; what you can expect

Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein

## 1967

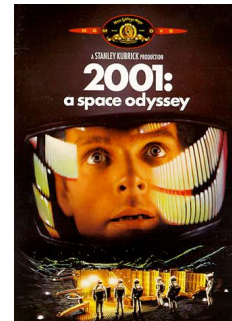


Stanley Kubrick,  
filmmaker  
1928 - 1999



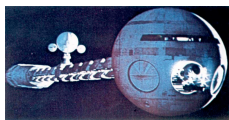
Arthur C. Clarke,  
author, futurist,  
1917 -

Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein



Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein

## HAL



Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein


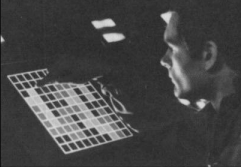
## HAL's Capabilities

- Display graphics
- Play chess
- *Natural language production and understanding*
  
- Vision
- Planning
- Learning
- ...


Andrew McCallum, ©1999, Amherst, including material from Chris Manning and Jacob Eisenstein

## Graphics

**HAL**





**Now**

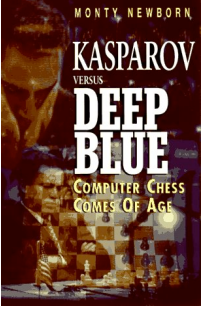


## Chess

**HAL**



**Now**



## Natural Language Understanding


**HAL**

*David Bowman:*  
Open the pod bay doors, Hal.

*HAL:*  
I'm sorry, Dave, I'm afraid I can't do that.

*David Bowman:*  
What are you talking about, Hal?

*...HAL:*  
I know that you and Frank were planning to disconnect me, and I'm afraid that's something I cannot allow to happen.




**Now**

# ?

Many useful tools, but none that come even close to HAL's ability to communicate in natural language.

## 1950



*Alan Turing*  
1912 - 1954

**Turing Test**  
"Computing Machinery and Intelligence"  
*Mind*, Vol. 59, No. 236, pp. 433-460, 1950

I propose to consider the question "Can machines think?" ... We can only see a short distance ahead, but we can see plenty there that needs to be done.

## Layers of Natural Language Processing

1. Phonetics & Phonology
2. Morphology
3. Syntax
4. Semantics
5. Pragmatics
6. Discourse

## 1. Phonetics & Phonology

The study of: language sounds, how they are physically formed; systems of discrete sounds, e.g. languages' syllable structure.

**dis-k&-'nekt**

**disconnect**

"It is easy to recognize speech."

"It is easy to wreck a nice beach."

JeetJet?

## 2. Morphology

The study of the sub-word units of meaning.

**disconnect**

“not” “to attach”

Even more necessary in some other languages, e.g. Turkish:

*uygarlastiramadiklarimizdanmissinizcasina*

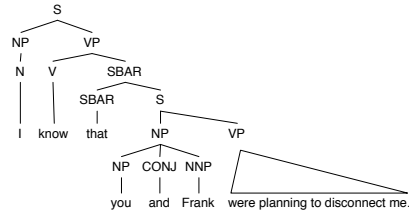
*uygar las tir ama dik lar imiz dan mis siniz casina*  
(behaving) as if you are among those whom we could not civilize

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## 3. Syntax

The study of the structural relationships between words.

I know that you and Frank were planning to disconnect me.



Not same structure:

You know me--Frank and I were planning to disconnect that.

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## 4. Semantics

The study of the literal meaning.

I know that you and Frank were planning to disconnect me.

ACTION = disconnect  
ACTOR = you and Frank  
OBJECT = me

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## 5. Pragmatics

The study of how language is used to accomplish goals.

What should you conclude from the fact I said something?  
How should you react?

I'm sorry Dave, I'm afraid I can't do that.

Includes notions of polite and indirect styles.

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## 6. Discourse

The study of linguistic units larger than a single utterance.

The structure of conversations:  
turn taking, thread of meaning.

David Bowman:  
Open the pod bay doors, Hal.  
HAL:  
I'm sorry, Dave, I'm afraid I can't do that.  
David Bowman:  
What are you talking about, Hal?  
...HAL:  
I know that you and Frank were planning to disconnect me, and I'm afraid that's something I cannot allow to happen.

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## Linguistic Rules

E.g. Morphology

To make a word plural, add “s”

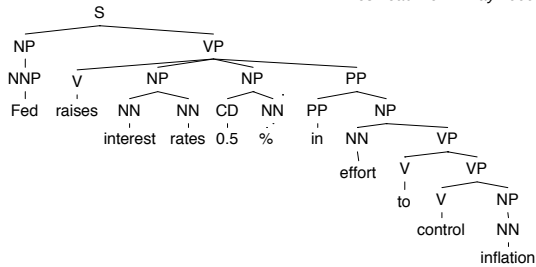
- dog → dogs
- baby → babies
- dish → dishes
- goose → geese
- child → children
- fish → fish (!)

Andrew McCulloch, UMass Amherst, including material from Chris Manning and Jason Eisner

## Inherent Ambiguity in Syntax

Fed raises interest rates 0.5%  
in effort to control inflation

NY Times headline 17 May 2000



Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## Where are the ambiguities?

Part-of-speech ambiguities

Syntactic attachment ambiguities

NNP VBZ VBZ VBZ CD NN  
NNS NNS NNS NNS

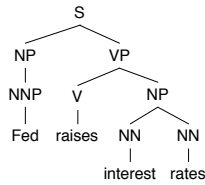
Fed raises interest rates 0.5 % in effort to control inflation

Word sense ambiguities: Fed → "federal agent"  
interest → a feeling of wanting to know or learn more

Semantic interpretation ambiguities above the word level.

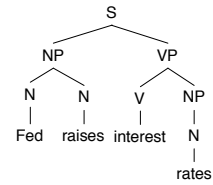
Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## Effects of V/N Ambiguity (1)



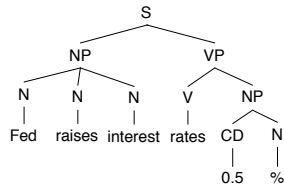
Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## Effects of V/N Ambiguity (2)



Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## Effects of V/N Ambiguity (3)



Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## Ambiguous Headlines

- Iraqi Head Seeks Arms
- Juvenile Court to Try Shooting Defendant
- Teacher Strikes Idle Kids
- Stolen Painting Found by Tree
- Kids Make Nutritious Snacks
- British Left Waffles on Falkland Islands
- Red Tape Holds Up New Bridges
- Clinton Wins on Budget, but More Lies Ahead
- Ban on Nude Dancing on Governor's Desk

Andrew McCallum, CMU, Andrew Senior, including material from Chris Manning and Jason Eisner

## What is grammatical and what isn't?

- John I believe Sally said Bill believed Sue saw.
- What did Sally whisper that she had secretly read?
- John wants very much for himself to win.
- Who did Jo think said John saw him?
- The boys read Mary's stories about each other.
- Mary, while John had had had had had had had had had was the correct answer.

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## What is grammatical and what isn't?

- John I believe Sally said Bill believed Sue saw.
- What did Sally whisper that she had secretly read?
- John wants very much for himself to win.
- Who did Jo think said John saw him?
- The boys read Mary's stories about each other.
- Mary, while John had had "had" had had "had had," "had had" was the correct answer.

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## Language Evolves

- Morphology
  - We learn new words all the time:  
bioterrorism, cyberstalker, infotainment,  
thumb candy, energy bar
- Part-of-speech
  - Historically: "kind" and "sort" were always *nouns*:  
"I knowe that sorte of men ryght well." [1560]
  - Now also used as *degree modifiers*:  
"I'm sort of hungry." [Present]  
"It sort o' stirs one up to hear about old times." [1833]

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## Natural Language Computing is hard because

- Natural language is:
  - highly ambiguous at all levels
  - complex and subtle
  - fuzzy, probabilistic
  - involves reasoning about the world
  - embedded a social system of people interacting
    - persuading, insulting and amusing them
    - changing over time

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## Probabilistic Models of Language

To handle this ambiguity and to integrate evidence from multiple levels we turn to:

- Bayesian Classifiers (not rules)
- Hidden Markov Models (not DFAs)
- Probabilistic Context Free Grammars
- Maximum Entropy models
- ...other tools of Machine Learning, AI, Statistics

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## Natural Language Processing

- Natural Language Processing (NLP) is the study of the computational treatment of natural languages:
  - Most commonly Natural Language Understanding
  - The complementary task is Natural Language Generation
- NLP draws on research in Linguistics, Theoretical Computer Science, Artificial Intelligence, Mathematics and Statistics, Psychology, etc.

Andrew McCallum, CMU, Andrew  
Inhalating material from Chris Manning and Jason Eisner

## What & Where is NLP

- Goals can be very far-reaching
  - True text understanding
  - Reasoning and decision-making from text
  - Real-time spoken dialog
- Or very down-to-earth
  - Searching the Web
  - Context-sensitive spelling correction
  - Analyzing reading-level or authorship statistically
  - Extracting company names and locations from news articles.
- These days, the later predominate (as NLP becomes increasingly practical, focused on performing measurably useful tasks *now*).
- Although language is complex, and ambiguity is pervasive, NLP can also be surprisingly easy sometimes:
  - rough text features often do half the job

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 1950s

- Early NLP on machines less powerful than pocket calculators.
- Foundational work on automata, formal languages, probabilities and information theory.
- First speech systems (Davis et al, Bell Labs).
- MT heavily funded by military, but basically just word substitution programs.
- Little understanding of natural language syntax, semantics, pragmatics.

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 1960s

- Alvey report (1966) ends funding for MT in America - the lack of real results realized
- ELIZA (MIT): Fraudulent NLP in a simple pattern matcher psychotherapist
  - It's true, I am unhappy.
  - *Do you think coming here will make you not to be unhappy?*
  - I need some help; that much is certain.
  - *What would it mean to you if you got some help?*
  - Perhaps I could learn to get along with my mother.
  - *Tell me more about your family.*
- Early corpora: Brown Corpus (Kudera and Francis)

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 1970s

- Winograd's SHRDLU (1971): existence proof of NLP (in tangled LISP code).
- Could interpret questions, statements commands.
  - Which cube is sitting on the table?
  - *The large green one which supports the red pyramid.*
  - Is there a large block behind the pyramid?
  - *Yes, three of them. A large red one, a large green cube, and the blue one.*
  - Put a small one onto the green cuube with supports a pyramid.
  - OK.

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 1980s

- Procedural --> Declarative (including logic programming)
- Separation of processing (parser) from description of linguistic knowledge.
- Representations of meaning: procedural semantics (SHRDLU), semantic nets (Schank), logic (perceived as answer; finally applicable to real languages (Montague)
- Perceived need for KR (Lenat and Cyc)
- Working MT in limited domains (METEO)

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 1990s

- Resurgence of finite-state methods for NLP: in practice they are incredibly effective.
- Speech recognition becomes widely usable.
- Large amounts of digital text become widely available and reorient the field. The Web.
- Resurgence of probabilistic/statistical methods, led by a few centers, especially IBM (speech, parsing, Candide MT system), often replacing logic for reasoning.
- Recognition of *ambiguity* as key problem.
- Emphasis on machine learning methods.

Andrew McCallum, CMU, Andrew  
McClosky, neural networks, Meaning and Sense Error

## Some brief history: 2000s

- A bit early to tell! But maybe:
  - Emphasis on meaning and knowledge representation.
  - Emphasis on discourse and dialog.
  - Strong integration of techniques, and levels: bringing together statistical NLP and sophisticated linguistic representations.
  - Increased emphasis on unsupervised learning.
  - More integration of NLP components into larger systems.

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP

A screenshot of a Google search for "natural language processing". The search results include several links to academic and research resources, such as "Natural Language Processing" from Microsoft, "IS's Natural Language Group" at USC, and "Foundations of Statistical Natural Language Processing" from MIT Press. There are also sponsored links for "Natural Language Search" and "NLP News".

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP: MSWord spelling correction, grammar checking

If you use Microsoft Word you have no doubt noticed red any misspelled words (or, to be exact, all words that did you know that you can correct these errors simply Microsoft Word will give you a list of the words that if a word you want appears in the list) you simply pick it f

A screenshot of the Microsoft Word spelling correction dropdown menu. The menu is open, showing a list of suggested words: "appears", "appease", "apparels", "appeals", and "appear". Below the list are buttons for "Ignore All", "Add", and "AutoCorrect". The word "appears" is highlighted, indicating it is the suggested correction.

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP: News categorization and summarization

A screenshot of Google News. The page displays a grid of news articles, each with a category label (e.g., "U.S.", "Business", "Sports") and a brief summary. The articles are organized into sections like "Top Stories" and "News Alerts". The layout is clean and easy to navigate, with clear categorization and summarization of the news content.

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP: Information Extraction: Find experts, employees

A screenshot of a resume for Dr. Andrew McCallum. The resume is structured with various sections, including "Other Titles Held", "Additional Current Employment", "Board Memberships and Affiliations", "Past Employment History", and "Education". The text is extracted from a document and presented in a clear, organized format, demonstrating the application of NLP for information extraction.

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP: Information Extraction: Job Openings

A screenshot of a job opening for "Ice Cream Guru" at foodscience.com. The job details are extracted and highlighted with red boxes and arrows. The extracted information includes: "Job Title: Ice Cream Guru", "Employer: foodscience.com", "Job Category: Travel/Hospitality", "Job Function: Food Services", "Job Location: Upper Midwest", "Contact Phone: 800-488-2611", "Date Extracted: January 8, 2001", "Source: www.foodscience.com/jobs\_midwest.htm", and "Other Company Jobs: foodscience.com-Job I".

Andrew McCallum, Editor, Andrew, including several Ben-Chen Manning and Jason Eisner

## Example Applications of NLP: Information Extraction: Job Openings

## Example Applications of NLP: Automatically Solving Crossword Puzzles

## Example Applications of NLP: Question Answering

## Example Applications of NLP: Machine Translation

## Example Applications of NLP: Automatically generate Harlequin Romance novels?

## Goals of the Course

- Introduce you to NLP problems and solutions.
- Relation to linguistics & statistics.
- Give you some hands-on practice with data and a handful of methods.
- At the end you should
  - Agree that language is subtle and interesting.
  - Feel some ownership over the formal & statistical models.
  - Be able to build some useful NLP system of your choosing.

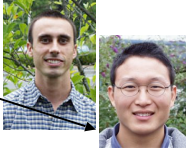


## This Class

- Assumes you come with some skills...
  - Some basic statistics, decent programming skills (in a language of your choice—although solutions will be in Java)
  - Some ability to learn missing knowledge
- Teaches key theory and methods for language modeling, tagging, parsing, etc.
- But it's something like an “AI Systems” class:
  - Hands on with data
  - Often practical issues dominate over theoretical niceties

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein

## Course Logistics

- Professor: Andrew McCallum
  - TA: Aron Culotta  
Gary Huang
- 
- Time: Tue/Thu 2:30-3:45pm
  - Mailing list: [cs585@cs.umass.edu](mailto:cs585@cs.umass.edu)
  - More information on Web site:  
<http://www.cs.umass.edu/~mccallum/courses/inlp2004>

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein

## Grading

- 5 short written homeworks
  - should take less than 30 minutes each
  - some hands-on experience
  - help you set expectations for the mid-term and final
- 3 programming assignments
  - no way to really internalize without doing it
  - should be fun!
- Final project: with a partner
  - chance to explore a special interest at end of term
- Midterm & Final, and classroom participation

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein

## Syllabus Outline

- Grammars and parsing
- Foundations (probability & info theory)
- Language models, Spam filtering.
- Collocations, word clustering, disambiguation.
- Finite state machines, Markov models, Part-of-speech tagging.
- Modern parsing techniques.
- Information extraction, Semantics, Question answering, Dialog systems.

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein

## Recommended Reading

- Manning & Schütze
  - Chapter 11, section 1  
Context Free Grammars, topic of next class
- Manning & Schütze
  - Chapter 3, for background on linguistics.

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein

Thank you!

Andrew McCallum, CMU, Stanford, including material from Chris Manning and Jacob Eisenstein