

Natural Language Processing

COMP-599

Sept 5, 2017

Preliminaries

Instructor:	Jackie Chi Kit Cheung
Time and Loc.:	TR 16:05-17:25 in MAASS 217
Office hours:	T 14:30-15:45 or by appointment in MC108N
TAs:	Ali Emami, Jad Kabbara, Kian Kenyon-Dean, Krtin Kumar
Evaluation:	4 assignments (40%) 1 midterm (20%) 1 group project (40%)

The Course Is Full

If you've registered for more courses than you plan to take, please decide soon! Many students are trying to get into this course.

Due to resource and classroom size limits, I cannot extend the class size anymore.

General Policies

Lateness policy for assignments:

- < 15 minutes: no penalty
- 15 minutes – 24 hours: 10% absolute penalty
- > 24 hours: not accepted

Plagiarism: just don't do it.

Language policy: In accordance with McGill policy, you have the right to write essays and examinations in English or in French.

Course website: <http://cs.mcgill.ca/~jcheung/teaching/fall-2017/comp550/index.html>

Important announcements given **in-class** or **on the course website**, not on MyCourses

Assignments

Four assignments (10% each)

Involve readings, problem sets and programming component.

Programming component – hand in online through myCourses

Programming to be done in Python 2.7.

Non-programming components – hand in on paper in class

Midterm

Worth 20% of your final grade

Currently scheduled for Thu, November 9, 2017

Will be conducted in-class (80 minutes long). More details as we approach the midterm date.

Final Project

Worth 40%.

Experiment on some language data set

Summarize and review relevant papers

Report on experiments

Must be done in teams of two

Coming up with a project idea:

- Extend a model we see in class
- Work on a relevant topic of interest
- Consult a list of suggested projects, to be posted

Project Steps

Paper or project proposal

Progress update

Final submission

Due dates to be announced

Computational Linguistics and Natural Language Processing

Language is Everywhere

CBC News Events Weather Programs Video Audio Contact Us

NEW | Hiker Julien Landry rescued days after fleeing up a tree to avoid bear

Hiker climbed a tree after a mother bear charged him - with incredible unexpected consequences

CBC News - Posted on 2014-03-03 10:00 PM PT | Last Updated on 03/21/2014 12:01 PM PT

Stay Connected with CBC News



Could not load Julien Landry, 25, who is in stable condition after he climbed a tree to escape a mother bear in Trout Creek, B.C. (Facebook)

4 shares

- Facebook
- Twitter
- Reddit
- Google+
- Print
- Email

A Quebec man is in a stable condition in a Kelowna hospital spending several days injured and alone in the forest after a mother bear attack.

After a day's work in the orchards around near 5000 S.C., Julien Landry, 25, of Trois-Rivières, Que., was in the Trout Creek canyon when a bear charged, forcing him to climb a tree.

It is not clear whether the bear and her cubs were looking for berries but as they circled the tree below, Landry climbed in the branches for hours, growing increasingly disoriented.

"Eventually he fell asleep because he'd been working all day in the orchards," said RCMP Const. Jacques Lefebvre. "When he fell asleep he fell down off the tree and landed on some rocks in the creek."

Lying unconscious in the creek, it was a day and a half before Landry awoke. He eventually managed to drag himself out of the water but was too weak to walk.

A search and rescue team including an RCMP helicopter and a plane could not find him.

It was three more days before another hiker spotted Landry, who was unable to move because he had buried himself in dirt to keep warm.

Landry suffered a concussion, bleeding in his head and broken vertebrae and was rushed to undergo emergency surgery. Doctors say he is in good recovery.

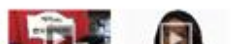
"I don't think he could have gotten himself out of there," said Lefebvre.



18.
 Shall I compare thee to a Summers day?
 Thou art more lovely and more temperate:
 Rough winds do shake the darling buds of Maie,
 And Sommers lease hath all too short a date:
 Sometime too hot the eye of heaven shines,
 And often is his gold complexion dimm'd,
 And every faire from faire some-time declines,
 By chance, or natures changing course vntim'd;
 But thy eternal Sommer shall not fade,
 Nor loose possession of that faire thou ow'st,
 Nor shall death brag thou wand'rst in his shade,
 When in eternal lines to time thou grow'st,
 So long as men can breathe or eyes can see,
 So long lives this, and this gives life to thee,

- Related Stories
- How to survive a bear encounter
 - Outrageous bear attack survivor was grabbed from forest
 - Furiously angry survivor bear attack by stomping on claws
 - B.C. man kills grizzly that attacked him
 - He's eating my brains. I can feel it" warns bear

2:33 Scientists have some surprising news about going in the ocean



0:40 Orphaned bear cub was rescued in June after he hibernated alone



Languages Are Diverse

6000+ languages in the world

language

langue

भाषा

語言

idioma

Sprache

lingua

→ The Great Language Game

<http://greatlanguagegame.com/> (My high score is 1300)

Computational Linguistics (CL)

Modelling natural language with computational models and techniques

Domains of natural language

Acoustic signals, phonemes, words, syntax, semantics, ...

Speech vs. text

Natural language understanding (or comprehension) vs. natural language generation (or production)

Computational Linguistics (CL)

Modelling natural language with computational models and techniques

Goals

Language technology applications

Scientific understanding of how language works

Computational Linguistics (CL)

Modelling natural language with computational models and techniques

Methodology and techniques

Gathering data: language resources

Evaluation

Statistical methods and machine learning

Rule-based methods

Natural Language Processing

Sometimes, **computational linguistics** and **natural language processing (NLP)** are used interchangeably.

Slight difference in emphasis:

NLP

Goal: practical
technologies

Engineering

CL

Goal: how language
actually works

Science

Understanding and Generation

Natural language understanding (NLU)

Language to form usable by machines or humans

Natural language generation (NLG)

Traditionally, semantic formalism to text

More recently, also text to text

Most work in NLP is in NLU

c.f. linguistics, where most theories deal primarily with production

Personal Assistant App

Understanding

Call a taxi to take me to the airport in 30 minutes.

What is the weather forecast for tomorrow?

Generation

Machine Translation

I like natural language processing.



Automatische Sprachverarbeitung gefällt mir.

Understanding

Generation

Recommendation System

A system chats with you to discover what you like, and recommends an event to check out this weekend.

Understanding

Generation

Computational Linguistics

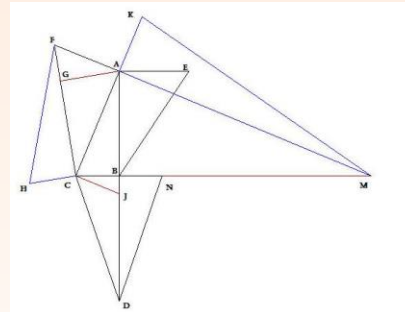
Besides new language technologies, there are other reasons to study CL and NLP as well.

The Nature of Language

First language acquisition

Chomsky proposed a **universal grammar**

Is language an “instinct”?



Do children have enough linguistic input to learn their mother tongue?

Train a model to find out!

The Nature of Language

Language processing

Some sentences are supposed to be grammatically correct, but are difficult to process.

Formal mathematical models to account for this.

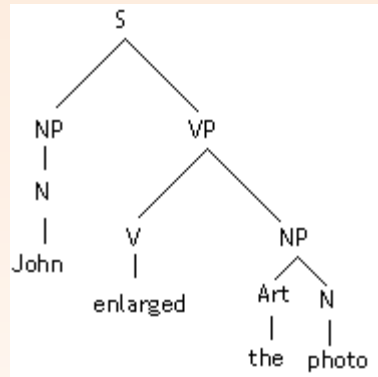
The rat escaped.

The rat the cat caught escaped.

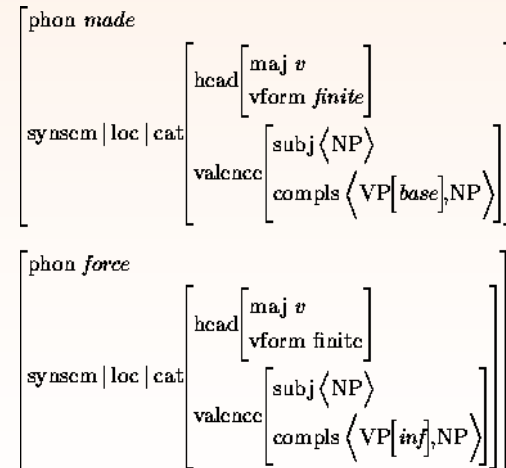
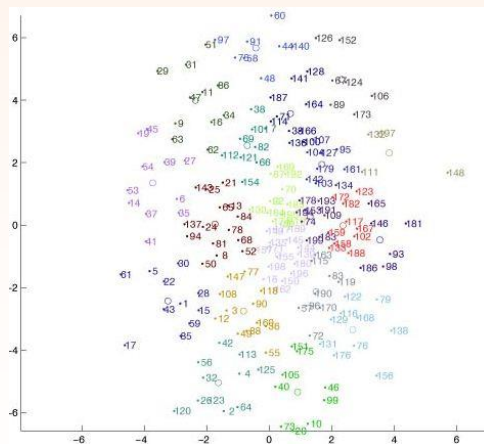
*?? The rat the cat **the dog chased** caught escaped.*

Mathematical Foundations of CL

We describe language with various formal systems.



cat + z > cats					
cat + z	*SS	Agree	Max	Dep	Ident
catiz				*!	
catis				*!	*
catz		*!			
cat			*!		
☞ cats					*



Mathematical Foundations of CL

Mathematical properties of formal systems and algorithms

Can they be efficiently learned from data?

Efficiently recovered from a sentence?

Complexity analysis

Implications for algorithm design

Types of Language

Text

Much of traditional NLP work has been on news text.

Clean, formal, standard English, but very limited!

More recent work on diversifying into multiple domains

Political texts, text messages, Twitter

Speech

Messier: disfluencies, non-standard language

Automatic speech recognition (ASR)

Text-to-speech generation

Domains of Language

The grammar of a language has traditionally been divided into multiple levels.

Phonetics

Phonology

Morphology

Syntax

Semantics

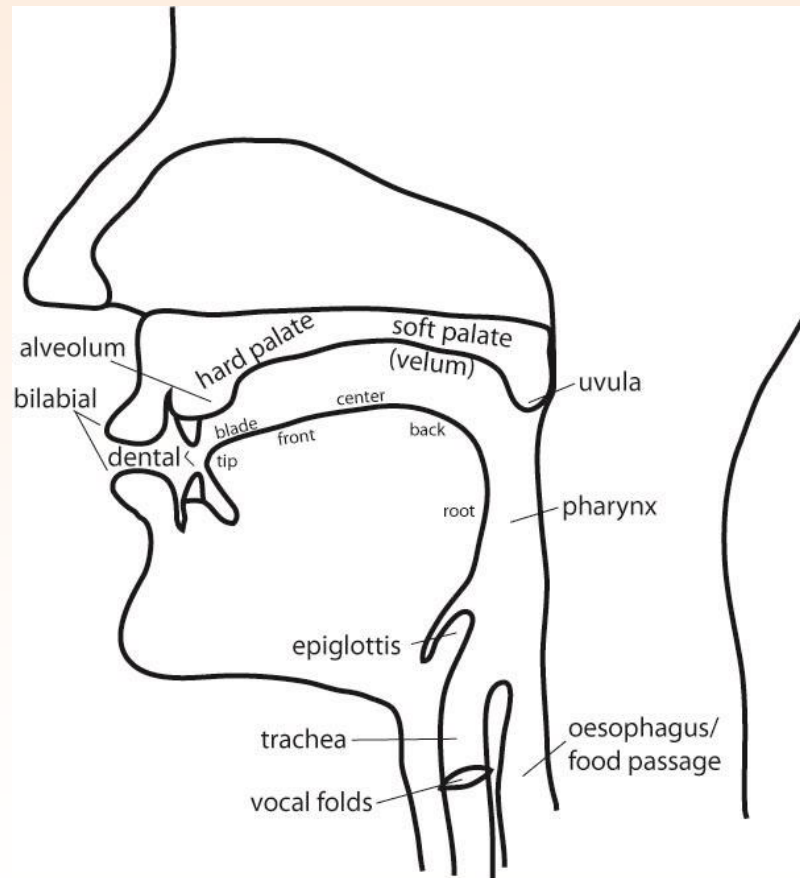
Pragmatics

Discourse

Phonetics

Study of the speech sounds that make up language

Articulation, transmission, perception



peach

[phi:tsh]

Involves closing of the lips, building up of pressure in the oral cavity, release with aspiration, ...

Vowel can be described by its formants, ...

Phonology

Study of the rules that govern sound patterns and how they are organized

peach [pi:tʃ]

speech [spi:tʃ]

beach [bi:tʃ]

The p in peach and speech are the same phoneme, but they actually are phonetically distinct!

Morphology

Word formation and meaning

antidisestablishmentarianism

anti- dis- establish -ment -arian -ism

establish

establish**ment**

establishment**arian**

establishmentarian**ism**

disestablishmentarianism

antidisestablishmentarianism

Syntax

Study of the structure of language

*I a woman saw park in the.

I saw a woman in the park.

There are two meanings for the sentence above! What are they? This is called **ambiguity**.

Semantics

Study of the meaning of language

bank

Ambiguity in the **sense** of the word



Semantics

Ross wants to marry a Swedish woman.

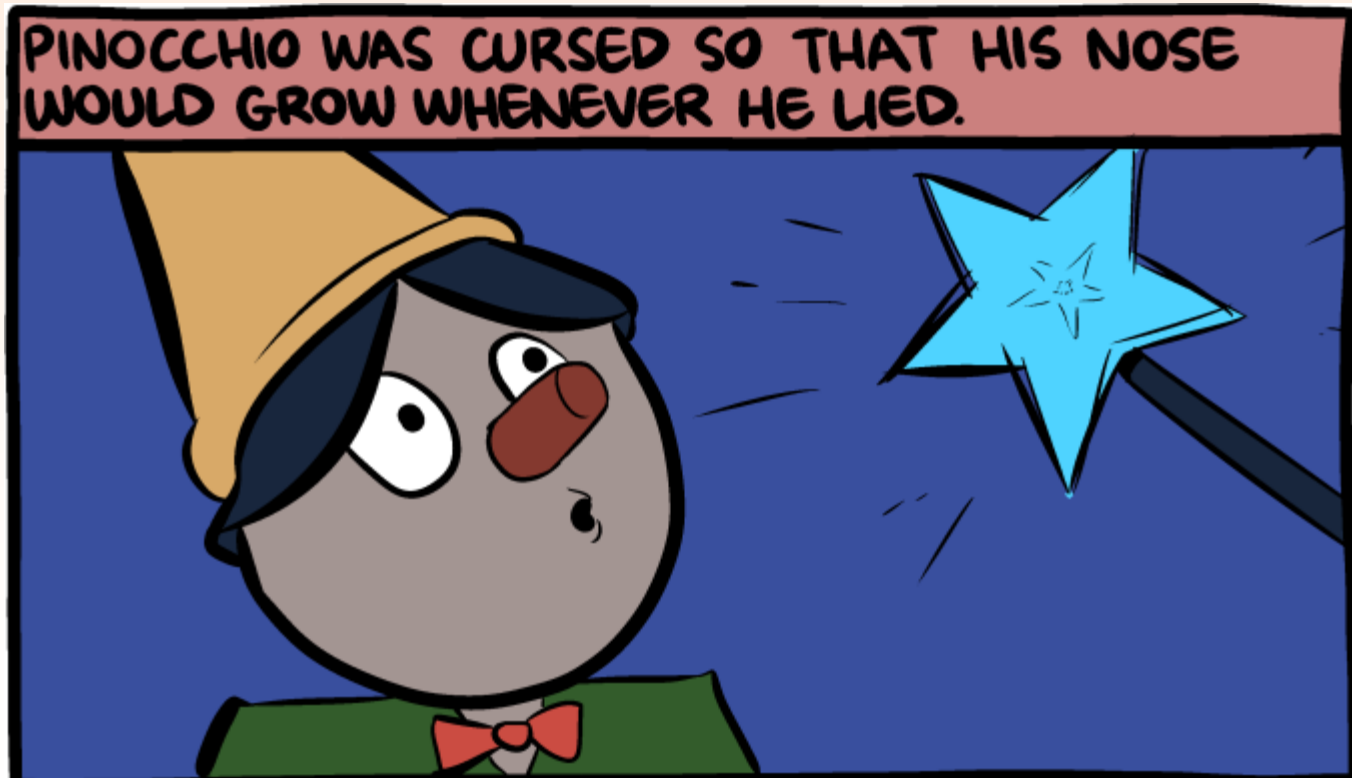


Pragmatics

Study of the meaning of language in context.

→ Literal meaning (semantics) vs. meaning in context:

<http://www.smbc-comics.com/index.php?id=3730>



Pragmatics



Pragmatics



Pragmatics



Pragmatics - Deixis

Interpretation of expressions can depend on **extralinguistic** context

e.g., pronouns

I think cilantro tastes great!

The entity referred to (the **antecedent**) by *I* depends on who is saying this sentence.

Discourse

Study of the structure of larger spans of language (i.e., beyond individual clauses or sentences)

I am angry at her.

She lost my cell phone.

I am angry at her.

The rabbit jumped and ate two carrots.

Questions

1. What is the difference between phonetics and phonology?
2. What are two possible readings of this phrase?
What level does the ambiguity act at? (i.e., lexical, syntactic, semantic, discourse)
 - *old men and women*

Topics in COMP-550

Progress through the subfields, roughly organized by the level of linguistic analysis

Morphology -> Syntax -> Semantics -> Discourse

NLP problems:

- Language modelling, part-of-speech tagging, parsing, word sense disambiguation, semantic parsing, coreference resolution, discourse coherence modelling

Focus on:

Basic linguistics needed to understand NLP issues

Algorithms and problem setups

Machine Learning in COMP-550

Interspersed throughout the course, and introduced as necessary

Machine learning topics we will cover:

- Feature extraction
- Sequence and structure prediction algorithms
- Probabilistic graphical models
- Linear discriminative models
- Neural networks and deep learning

Applications in COMP-550

Last three weeks of the course focus on language technology applications and advanced topics:

- Automatic summarization

- Machine translation

- Evaluation issues in NLP

Course Objectives

Understand the broad topics, applications and common terminology in the field

Prepare you for research or employment in CL/NLP

- Learn some basic linguistics

- Learn the basic algorithms

- Be able to read an NLP paper

Understand the challenges in CL/NLP

- Answer questions like “Is it easy or hard to...”

Plan for the Next Week

I will be away at a conference for the next week

Thursday's class:

- Lecture by TA Krtin Kumar on finite state machines for morphology

Tuesday's class:

- Python tutorial + a presentation of a NLP research project by TA Jad Kabbara

This means no office hours next Tuesday. E-mail me if you need to discuss anything.