# CS-585
# Natural Language Processing

Prof. Derrick Higgins

dhiggins1@iit.edu

# Today

1. About the course
2. About me
3. About you
4. About language and linguistics
5. Math

# THIS COURSE

# Goals

- **Breadth of coverage**: Familiarity with a wide range of task types and methods in natural language processing

- **Depth in key areas**: Mastery of critical concepts and algorithms for NLP

- **Preparedness for further study**: introduction to deep learning frameworks as applied to NLP, such that current research papers will be accessible

# Prerequisites

- Math
  - Linear algebra
  - Probability
- Programming
  - Python 3
  - Basic algorithms and data structures
  - Access to a Unix system

# Methods

- Exams
  - Open-book, multiple choice (mostly)
  - No electronics allowed
  - Midterm will cover material through October 5
  - Final will cover material from the entire course
- Class Project
  - We will create a new dataset for text categorization
  - Three parts
    1. Data labeling
    2. Annotation analysis
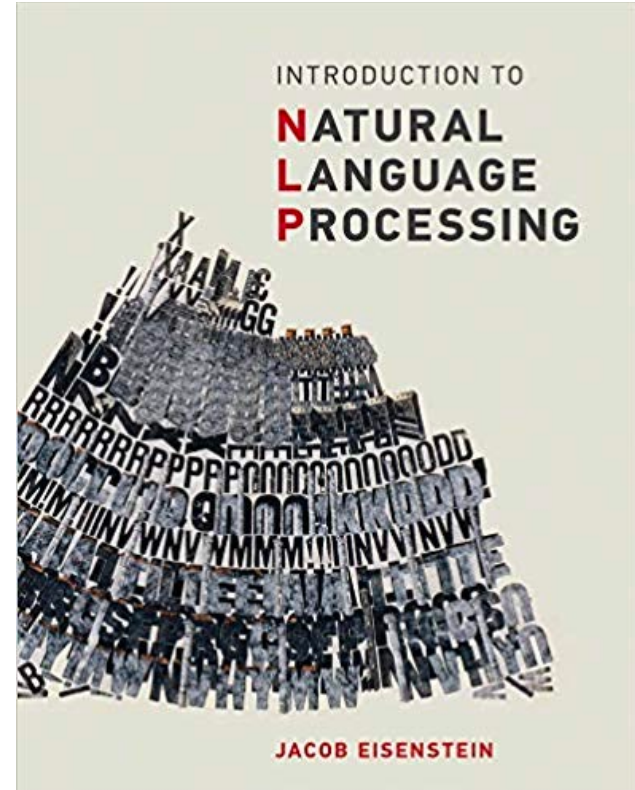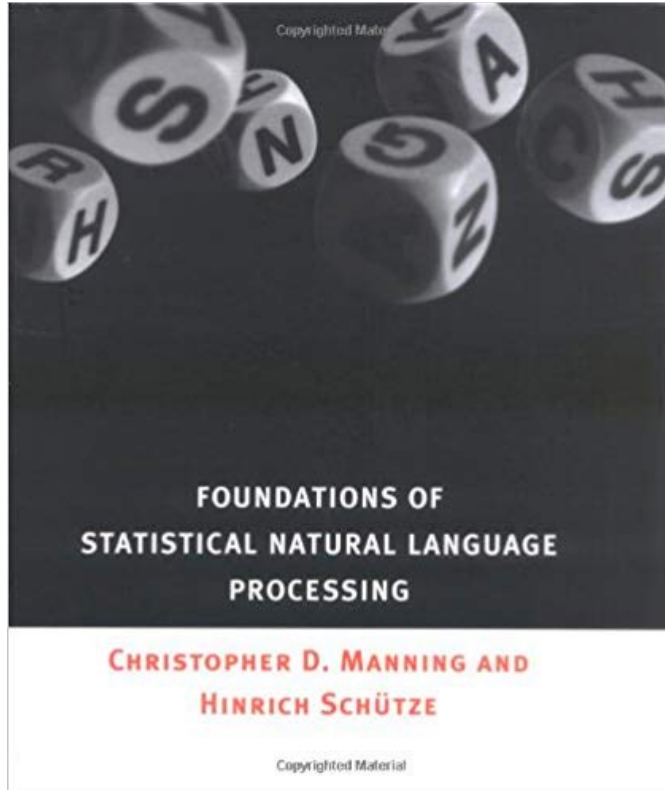    3. NLP Modeling

# Grading: Available Points

| Assignment | Points |
|---|---|
| Project part 1: Annotation | 175 |
| Project part 2: Analysis | 150 |
| Project part 3: Modeling | 175 |
| Midterm | 200 |
| Final | 300 |
| Total: | 1000 |

# Grading: Letter Grades

| Points | Grade |
|--------|-------|
| 900-1000 | A |
| 750-899 | B |
| 550-749 | C |
| 0-550 | E |

# Readings

# Lecture plan

1. Foundational concepts

2. Consideration of progressively higher levels of linguistic structure

3. Connection to neural networks at end of each unit

# Academic Honesty

- If you violate the academic honesty policy (such as unauthorized/undocumented collaboration, cheating, etc.), I **will** report it to the university

- Depending on the severity of the violation, it can result in
  - zero points on the respective assignment,
  - E in the course,
  - suspension from the university,
  - expulsion from the university

- Full guidelines: https://web.iit.edu/student-affairs/handbook/fine-print/code-academic-honesty

# Americans With Disabilities Act

Reasonable accommodations will be made for students with documented disabilities. In order to receive accommodations, students must obtain a letter of accommodation from the Center for Disability Resources.

The Center for Disability Resources (CDR) is located in 3424 S. State St., room 1C3-2 (on the first floor), telephone: 312.567.5744 or disabilities@iit.edu

ILLINOIS INSTITUTE OF TECHNOLOGY

*Transforming Lives. Inventing the Future.* **www.iit.edu**

# INTRODUCTION

# Me



Civis Analytics

# Open-ended Written Responses



- Cf. Shermis, Mark, Jill Burstein, Derrick Higgins & Klaus Zechner. (2009). Automated essay scoring: Writing assessment and instruction. International Encyclopedia of Education, Third Edition. United Kingdom, Elsevier.

# Data Science for Data Science

**02** /04

**AMERICAN FAMILY INSURANCE** ®

# Claims / Underwriting

American Family utilizes Arturo data throughout multiple parts of their business to improve underwriting performance, identify risk within their book by identifying changes, and predicting the right resources necessary to respond to claim events.

**ARTURO**

ILLINOIS INSTITUTE OF TECHNOLOGY

Transforming Lives. Inventing the Future. **www.iit.edu**

# GETTING TO KNOW YOU

ILLINOIS INSTITUTE
OF TECHNOLOGY
*Transforming Lives.Inventing the Future.* **www.iit.edu**

# Questions for you

How many of you

- …know python?

- …have worked with Unix shell?

- …have taken a machine learning course?

# LANGUAGE, LINGUISTICS AND NLP

# Some terminology

- **Text processing**: Engineering practices for transforming, normalizing, compressing or accessing textual data

- **Natural language processing**: The study of methods for exploiting or generating language represented as text, for practical tasks

- **Computational linguistics**: The use of computational tools to understand or learn the structure of human languages

- **Speech processing**: The study of methods for exploiting or generating language represented as audible waveforms, for practical tasks

# Adjacent fields



| Linguistics | Computer Science | Engineering |
|---|---|---|
| Pragmatics | Machine Learning | |
| Semantics | Natural Language Processing | Text Processing |
| Syntax | | |
| Morphology | | |
| Phonology | Speech Processing | Digital Signal Processing |
| Phonetics | | |

# Phonetics

The study of speech sounds

- Articulatory phonetics deals with the physiological speech process
- Acoustic phonetics deals with the sound waves produced

Applications:

- Speech recognition
- Speech synthesis
- Clinical speech pathology



https://commons.wikimedia.org/w/index.php?curid=14508443

ILLINOIS INSTITUTE
OF TECHNOLOGY
Transforming Lives. Inventing the Future. www.iit.edu

# Phonology

The structure and patterning of sounds within a language

- Segmental phonology deals with phonemes (minimal contrastive units)

- Suprasegmental phonology deals with tones, prosody and stress accent

- Subsegmental phonology deals with features of phonemes



https://commons.wikimedia.org/w/index.php?curid=18555461

Applications:

- Speech recognition

- Speech synthesis

# Morphology

The internal structure of words

- Morphemes include stems, prefixes, suffixes and infixes

Applications

- Stemming / lemmatization
- Compound breaking
- Inflection generation (NLG)

mis　treat　ing

pre　judge　s

bil　m　iyor　um

know　not　[progressive]　I

# Syntax

The structure of words and phrases within a sentence

- Different formalisms, coming from the American (phrase structure) and European (dependency grammar) structuralist traditions

Applications

- Part-of-speech tagging
- Entity extraction
- Syntactic parsing (CFG)
- Syntactic parsing (dependencies)



*All birds can fly*

# Semantics

The representation of meaning in language

- At different levels: lexical, sentential, textual
- Logical formalisms: reference and truth conditions

Applications:

- Word embedding/encoding
- Lexical resources
- Semantic role labeling

$\forall x(\text{bird}(x) \rightarrow \text{fly}(x))$

$\text{kill}(x, y) :=$

$\text{Cause}(x, \text{Become}(\neg\text{Alive}(y)))$

# Pragmatics

How language is used to achieve specific intentions

- Conversational implicatures: how I interpret what you say because of what I assume you're trying to do
- Speech acts

Applications:

- Speech act labeling
- Discourse structure parsing
- Dialogue systems

"I ate <u>most</u> of your cookies"

⊨

I did not eat <u>all</u> of your cookies

---

"Where does your brother live?"

⊨

I do not know where your brother lives

# Sociolinguistics

Language use patterns associated with particular groups, or language used to communicate status relative to a group

With friends

In class

With family

With strangers

With professional colleagues

Applications:

- Stylometrics / authorship attribution
- Forensic linguistics
- Natural language generation

# Historical Linguistics

Language change over time

- Lexical innovation
- Phonological change
- Language contact

Applications:

- Linguistic typology
- Digital humanities

ILLINOIS INSTITUTE
OF TECHNOLOGY
*Transforming Lives. Inventing the Future.* **www.iit.edu**

# Psycholinguistics

Language as a cognitive function

- Role of brain areas in language production and processing
- Language learning


http://arikaokrent.com/bio.html

Applications:

- Language pathology
- Assistive technology

ILLINOIS INSTITUTE
OF TECHNOLOGY
Transforming Lives. Inventing the Future. www.iit.edu

# And of course…

Not all NLP tasks relate to a single linguistic domain.

E.g., machine translation involves morphology, syntax, semantics and pragmatics (at least)

# Why is NLP hard?

- The "hidden structure" of language is ambiguous at all levels!

- Consider the simple proverb:

Time   flies   like   an   arrow

# Word sense ambiguity

*Time:* "abstract time", "a specific point in time", "to measure time"

*flies:* "moves through the air", "little pesky insects"

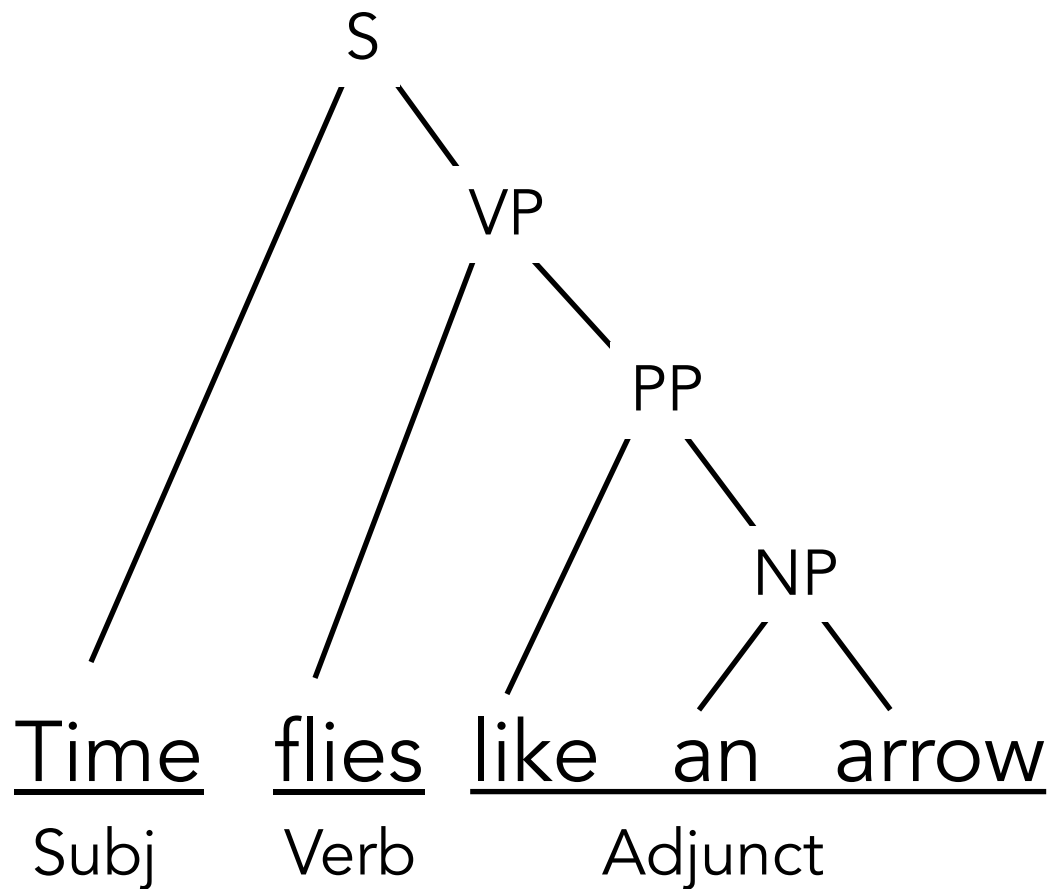*like:* "similar to", "have affection for"

*arrow:* "pointy stick shot from a bow", "to move straight towards a target"

Time   flies   like   an   arrow

# Part of speech ambiguity

|  |  | JJ |  |  |
|  |  | VB |  |  |
| VB | NNS | NN |  | VB |
| NN | VBZ | IN | DT | NN |
| Time | flies | like | an | arrow |

# Syntactic ambiguity



```
                    S
                   / \
                  /   VP
                 /   / \
                /   /   PP
               /   /   / \
              /   /   /   NP
             /   /   /   / \
          Time flies like an arrow
          Subj  Verb      Adjunct
```

# Syntactic ambiguity



S
VP
PP
NP

Time    flies    like    an    arrow

Verb    Obj        Adjunct
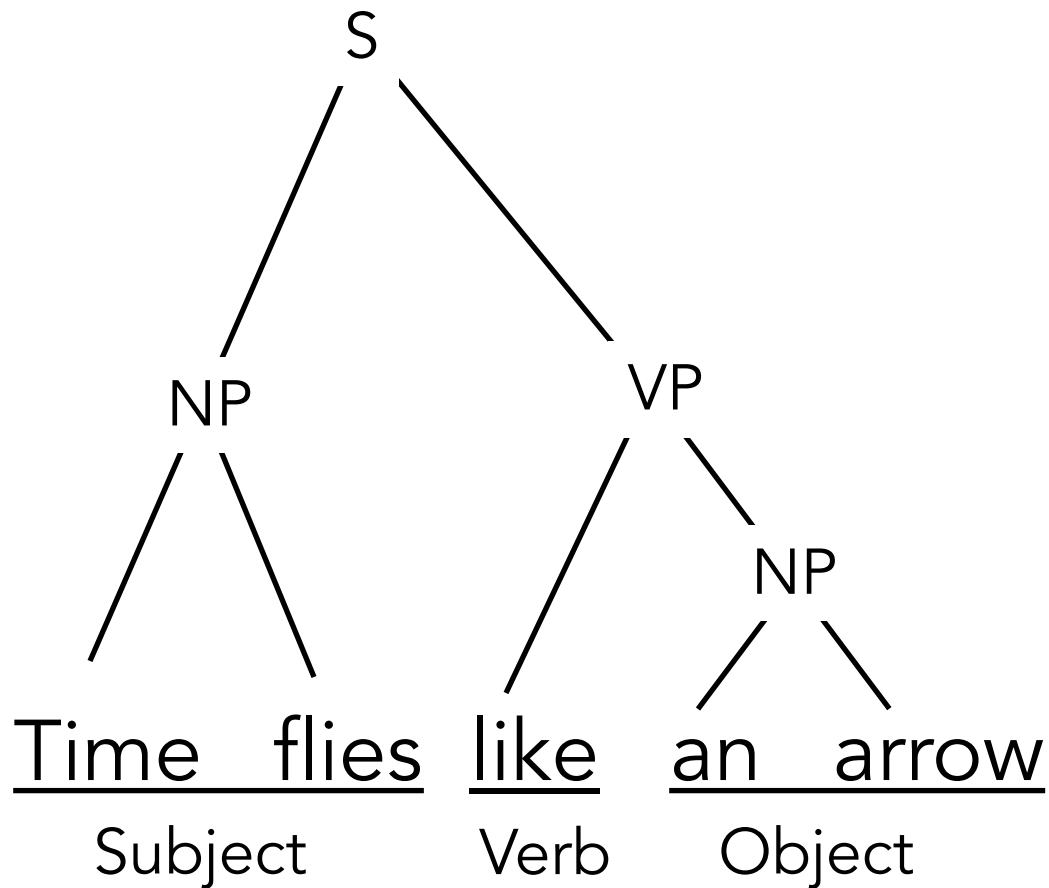
*...instead of timing them like a snail!*

# Syntactic ambiguity



...but fruit flies like a banana!
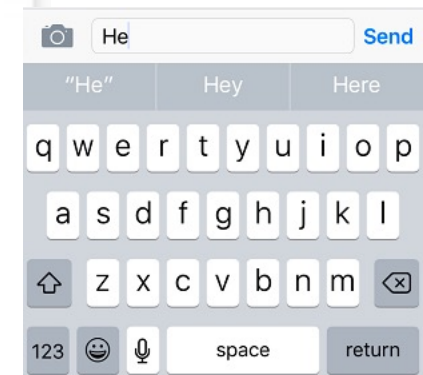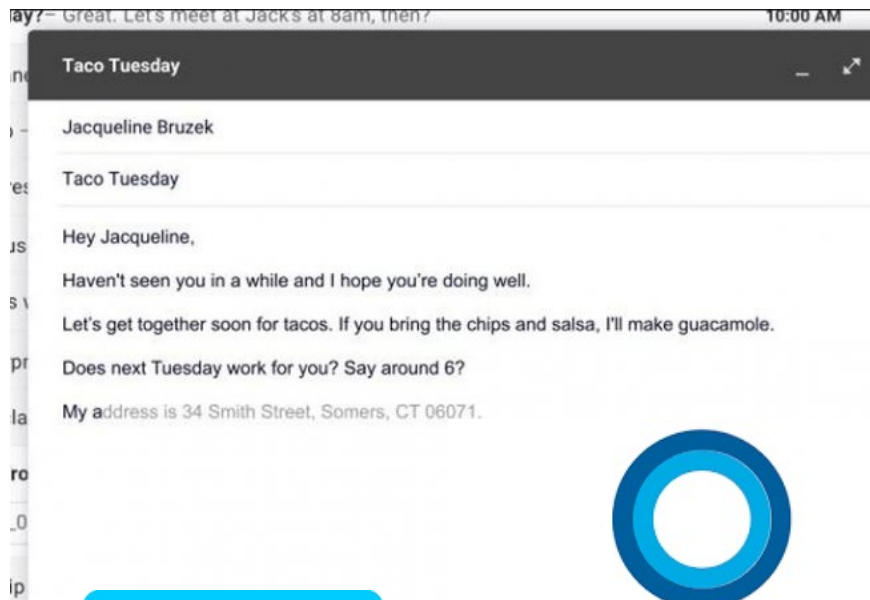
# Newspaper Headlines

- Ban on Nude Dancing on Governor's Desk
  *from a discussion of current legislation*
- Juvenile Court to Try Shooting Defendant
- Stolen Painting Found by Tree
- British Left Waffles on Falkland Islands
- Red Tape Holds Up New Bridges
- Kids Make Nutritious Snacks
- Hospitals Sued by Seven Foot Doctors

# A Changing Target

- Neologisms (= new words/phrases):
  - *cosmocrat, technocrat, davos man*
  - *megacryometeor*
  - *flash mob, carjack*
  - *googling, spam, blogger, wi-fi*
  - *kleptocracy, identity theft*
  - *just-in-time learning, egoboo*
- Also sentence structure, though it's subtler…

# Such a great time to get into NLP!

- There is so much we can do now!

# Such a great time to get into NLP!

- ## There is so much we <u>still can't</u> do!
  - Handle real-world knowledge and logical inferences
  - Deal with limited-data contexts and low resource languages
  - Transfer learning across tasks and domains (although we're getting better)
  - Integrate information across modalities: text, imagery, action sequences
  - Infer linguistic structure without manual labeling based on human judgements

# DEMO: SLACK

# Slack Channel Usage

- `#general` channel
  - Ask clarification questions publicly so that everyone can benefit from the answers
  - Email OK for personal concerns