

CS 2731 / ISSP 2230

# Introduction to Natural Language Processing

Session 1: Course introduction and NLP basics

---

Michael Miller Yoder

January 8, 2024

# Overview: Course introduction and NLP basics

- Introductions
- What is NLP?
- Course logistics

# About Michael Miller Yoder

- You can call me "Michael"
- Adjunct faculty, Pitt CS Department
- Postdoc, Software and Societal Systems Department at Carnegie Mellon University
- PhD, Language Technologies Institute at CMU (2021)
- **Research interests:**
  - NLP
  - computational social science
  - hate speech
  - language and identity



# Introductions

1. What is your name?
2. What is your program/year/research interests?
3. Is there anything that makes you interested in NLP or excited to take this class?

# What is natural language processing (NLP)?

---

# NLP is Everywhere

Did you ever wonder how web search engines work...



...or how Google can anticipate what you're searching for?

**That's NLP!**

# NLP is Everywhere

Did you ever wonder how digital assistants work?



**That's NLP!**

# NLP is Everywhere

Did you ever wonder how the government is spying on your every word?

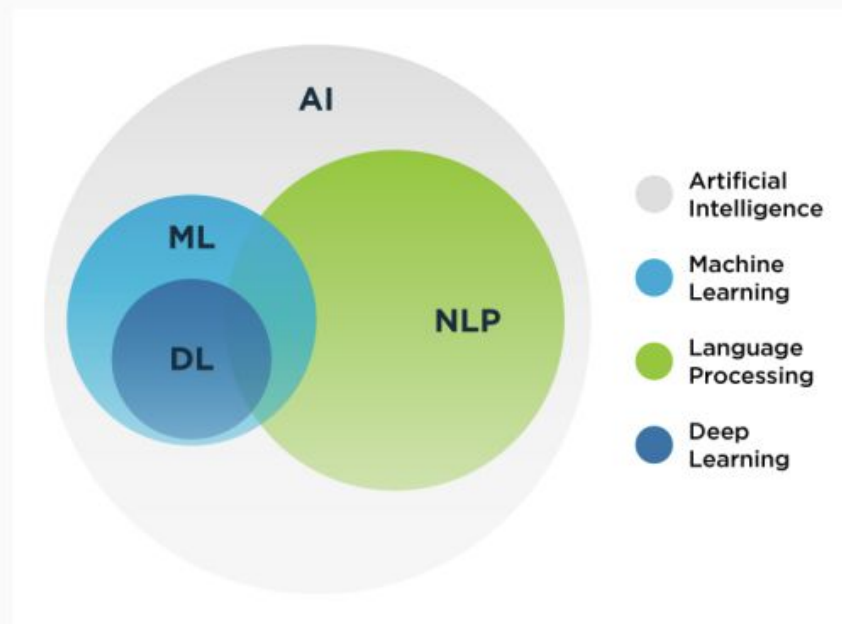


That's also NLP!



# NLP is the Computational Analysis and Synthesis of Speech and Language

- NLP is one of the most important AI fields today
- It is about processing language with computers
- Engineering focus—solving practical problems



# A brief history of NLP



Sentences in Russian are punched into standard cards for feeding into the electronic data processing machine for translation into English

*The Georgetown-IBM Experiment.  
Credit: John Hutchins*

- 1950s: **foundations**
  - Turing Test ("Computing Machinery and Intelligence" paper)
  - Georgetown-IBM Experiment translating Russian to English

# A brief history of NLP



Sentences in Russian are punched into standard cards for feeding into the electronic data processing machine for translation into English

*The Georgetown-IBM Experiment.  
Credit: John Hutchins*

- 1950s: **foundations**
  - Turing Test ("Computing Machinery and Intelligence" paper)
  - Georgetown-IBM Experiment translating Russian to English
- 1960s-1980s: **symbolic reasoning**
  - ELIZA, rule-based parsing, hand-built conceptual ontologies

# A brief history of NLP



Sentences in Russian are punched into standard cards for feeding into the electronic data processing machine for translation into English

*The Georgetown-IBM Experiment.  
Credit: John Hutchins*

- 1950s: **foundations**
  - Turing Test ("Computing Machinery and Intelligence" paper)
  - Georgetown-IBM Experiment translating Russian to English
- 1960s-1980s: **symbolic reasoning**
  - ELIZA, rule-based parsing, hand-built conceptual ontologies
- 1990s-2010s: **statistical NLP**
  - Learn patterns from large corpora (feature-based machine learning)

# A brief history of NLP



Sentences in Russian are punched into standard cards for feeding into the electronic data processing machine for translation into English

*The Georgetown-IBM Experiment.  
Credit: John Hutchins*

- 1950s: **foundations**
  - Turing Test ("Computing Machinery and Intelligence" paper)
  - Georgetown-IBM Experiment translating Russian to English
- 1960s-1980s: **symbolic reasoning**
  - ELIZA, rule-based parsing, hand-built conceptual ontologies
- 1990s-2010s: **statistical NLP**
  - Learn patterns from large corpora (feature-based machine learning)
- 2000s-today: **neural NLP**
  - SOTA on many tasks from "deep" layers of neural networks

# NLP and Computational Linguistics

- These terms are often used interchangeably
- If you want to make a distinction:
  - Computational linguistics is the scientific study of language using computers
  - Natural language processing is the development of computational tools to process human language (engineering-focused)
- "Natural language" = human languages (not programming languages)

# The other NLP 😂

Neuro-linguistic programming (pseudoscience)



The diagram illustrates the Neuro-Linguistic Programming (NLP) model. On the left, a human brain is shown with the left hemisphere in grayscale and the right hemisphere filled with a vibrant, multi-colored paint splatter. To the right, the title "NEURO LINGUISTIC PROGRAMMING" is displayed in large, bold, black capital letters. Below the title, a central brain icon is connected to three main components: "Neuro" (red line), "Linguistic" (yellow line), and "Programming" (green line). Each component has a sub-label: "Neuro" is "First Access Internal images Sounds and feelings"; "Linguistic" is "Linguistic Map Concrete word Description"; and "Programming" is "Behavioural response Neurological filtering Beliefs". A horizontal line with arrows at both ends passes through the brain icon, labeled "INPUT" on the left and "OUTPUT" on the right. In the top right corner, the "INNOVANS TECHNOLOGIES" logo is present, with "ISO 9001:2015 CERTIFIED" written below it. At the bottom, a yellow banner contains the text "NEURO-LINGUISTIC PROGRAMMING HELPS EMPLOYEE PERFORM BETTER" in bold, black capital letters.

NEURO LINGUISTIC PROGRAMMING

INPUT > > OUTPUT

**Neuro**  
First Access  
Internal images  
Sounds and feelings

**Linguistic**  
Linguistic Map  
Concrete word  
Description

**Programming**  
Behavioural response  
Neurological filtering  
Beliefs

**NEURO-LINGUISTIC PROGRAMMING HELPS EMPLOYEE PERFORM BETTER**

# Course objectives and overview

---



# Learning objectives

At the end of this course, a student will be able to structure an NLP system to achieve a desired outcome from language data.

# Learning objectives

When coming across a natural language problem, students will be able to:

- Recognize the class of tasks that a specific natural language task belongs to
- Explain the basics of language structure from linguistics (morphology, syntax, semantics, discourse) that are relevant to NLP
- Preprocess text into a machine-readable format
- Extract needed features from text for a variety of tasks
- Identify a suitable model to tackle the task
- Evaluate algorithms for that task
- Identify potential ethical pitfalls in an NLP system and how to potentially address them

# Core tasks and applications of NLP

## APPLICATIONS

machine  
translation

speech recognition  
& synthesis

chatbots

computational  
social science

summarization

information retrieval  
& extraction

question answering

# Core tasks and applications of NLP

## CORE TASKS

text  
classification

representation  
learning

language  
models

conditional  
language  
models

sequence  
labeling

syntactic  
parsing



## APPLICATIONS

machine  
translation

speech recognition  
& synthesis

chatbots

computational  
social science

summarization

information retrieval  
& extraction

question answering

# Core tasks and applications of NLP

## CORE TASKS

text  
classification

representation  
learning

language  
models

conditional  
language  
models

sequence  
labeling

syntactic  
parsing

MODULE 2

MODULE 3

MODULE 4

MODULE 5

## APPLICATIONS

machine  
translation

speech recognition  
& synthesis

chatbots

computational  
social science

MODULE 6

summarization

information retrieval  
& extraction

question answering

# Approaches covered in this course

For most NLP tasks, we will cover:

- Classic approaches: symbolic, statistical, feature-based approaches
- Contemporary approaches: neural network-based

# Resources

---

# Textbook (free)

- Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 3rd edition draft, 2023-01-07.
- **Available completely free online:**  
<https://web.stanford.edu/~jurafsky/slp3/>
- Why do the readings?
  - Learn better: get the information from readings and lectures
  - Spend class time more efficiently: come with questions
  - Reading quizzes due before class



# Lectures

- Cover the most important parts of the course content
- Students are expected to attend each lecture
- Slides will be provided in advance of each lecture for note-taking
- There are no current plans for recording lectures

# Infrastructure

- Website

- [https://michaelmilleryoder.github.io/cs2731\\_spring2024](https://michaelmilleryoder.github.io/cs2731_spring2024)
- <https://bit.ly/pittnlp>
- Up-to-date syllabus and schedule
- Lecture slides
- Homework assignment and project instructions

- Canvas

- Submit assignments  
(homeworks and project milestones)
- Post on discussion forums
- Receive course announcements
- Check your grade



# Programming languages and software

- Python will be the expected programming language used in assignments
- Python-based data science packages (numpy, pandas, jupyter, scikit-learn, pytorch) will be used and encouraged in both assignments and projects
- If you have zero familiarity with Python (no shame):
  - Check out the **Tutorials on Python and data science** section of the course website under **Learning resources**
- Let us know if you want to use other languages for assignments (it's probably fine)
- You can use whatever you want for the project

# Assessments

---

# Assessment overview

Assessment	Points	Percentage of grade
Homeworks (4 total)	224	44.8%
Project	222	44.4%
Reading quizzes	33	6.6%
Discussion posts	21	4.2%

No exams

# Homework assignments

- 4 total
- 11.2% of total course grade each
- Most will have a written component (working through an algorithm, e.g.) and a coding component
- Due ~15 to 18 days after they are released
- Descriptions will be on the course website
- Submitted through Canvas

# Project

NLP is inherently hands-on. The course project will demonstrate an ability to build a system that **makes a contribution** to NLP research or practice.

- Groups of 2-4, assigned by the instructor and TA based on interests, skills, and group preferences from students
  - There will be group member evaluation
- Self-selected topic, type of research contribution, and idea
  - Can fit with your research interests outside of this class
  - Come up with your own idea or choose one of the example project ideas
- Types of contributions: new dataset and/or annotations, new approach/application, new evaluation, new survey

# Project components

Component	Points	Percentage of course grade	Due
Interest survey response	5	1%	01-18
Project area and type of contribution	10	2%	02-08
Proposal and literature review	35	7%	02-22
Peer review	2	0.4%	02-29
Proposal presentation	<i>None</i>	<i>None</i>	03-04
Basic working system report	30	6%	04-04
Final presentation	<i>None</i>	<i>None</i>	04-24
Final report	140	28%	04-25



# Project survey

- [Google Form](#) that goes to the instructor
- Indicate interest in NLP research areas and type of contribution for the project
  - I don't expect you to be familiar with these! It's okay to just put your best guesses
- Indicate any prior experience in machine learning or NLP
- Indicate any group member preferences

**Due next Thursday, Jan 18 at 11:59pm ET**

# Reading quizzes

- On Canvas
- Quick checks for comprehension
- Designed to motivate you to do the reading
- Simple, auto-graded (generally multiple choice or short answer)
- "Muddiest point" free-text field (not graded)
  - Helps the instructor allocate class time in lecture
- Only 6.6% of your course grade total
- The lowest 2 quiz scores will be dropped
- **Due by 1pm on days with class sessions**
- Can't redo after they're due

**First will be due next Wednesday, Jan 17 at 1pm ET**

# Discussion posts

- There will be 3 required discussion posts from readings, 1 for extra credit
- Often on the social impact of NLP (bias, transparency, etc)
- Respond to a prompt
- Post on a discussion forum on Canvas
- Add your own ideas, respond to others
- Minimum 100 words with a substantive idea or response

# Policies

---

# Late work

- Contact the instructor and TA **before the deadline** if you need an extension due to unforeseen circumstances
- We are happy to extend deadlines for deaths and funerals, illnesses, mental health crises or episodes, weddings, important religious and national holidays, job interviews, and other circumstances
- No shame in asking. We care about your well-being more than we care about deadlines
- Unless you let us know beforehand (or an adverse event occurred very close to the deadline), **the late penalty is 2.5% per day**, including weekend days and holidays

# Academic integrity

- Students in this course will be expected to comply with the [University of Pittsburgh's Policy on Academic Integrity](#). Any student suspected of violating this obligation for any reason during the semester will be required to participate in the procedural process, initiated at the instructor level, as outlined in the University Guidelines on Academic Integrity
- Discussing tools, concepts, and formalisms is acceptable collaboration
- Sharing code is prohibited

# Generative AI policy

- You are welcome to use generative AI (ChatGPT, DALL-E, GitHub Copilot, etc)
  - Exposes you to the current capabilities and limitations of such systems
- Principles in how to use it
  - **Use as an aid, not for a finished product.** Generating ideas, study guides, bibliographies (watch for hallucinations, though) is ok. Drafting entire homework assignments or project reports, even if you revise the draft, is not ok.
  - **Cite its use.** Citing the generative AI's tool contribution to your work is required. See the [APA guidelines on how to cite ChatGPT](#).
  - **You are responsible for the work you turn in.** LLMs and other generative AI systems can and do generate biased, socially problematic language and assert unfounded claims.
- When in doubt, ask instructor if specific uses are ok. There will be no retaliation for asking.

# Disability rights

Many people have disabilities. **We view disabilities as deficits not in disabled people but in the institutions and societies that are structured to disadvantage disabled people.**

If you have a disability (visible or invisible), please let us know as soon as possible (you don't need to tell us the nature of the disability). You are encouraged to work with Disability Resources and Services (DRS), 140 William Pitt Union, (412) 648-7890, drsrecep@pitt.edu, (412) 228-5347 for P3 ASL users, as early as possible in the term. DRS will work with you to determine reasonable accommodations for this course. This might include lecture materials that are usable by people with visual disabilities, sign language interpretation, captioning, flexible due dates, etc.



# Maintaining scholarly discourse

In this course we will be discussing some complex issues. It is essential that we **approach this endeavor with our minds open** to evidence that may conflict with our presuppositions. Moreover, **it is vital that we treat each other's opinions and comments with courtesy even when they diverge and conflict with our own**. We must avoid personal attacks and the use of ad hominem arguments to invalidate each other's positions. Instead, we must develop a culture of civil argumentation, wherein all positions have the right to be defended and argued against in intellectually reasoned ways. It is this standard that everyone must accept in order to stay in this class; a standard that applies to all inquiry in the university, but whose observance is especially important in a course whose subject matter is so emotionally charged.

Questions?