

# CS 6120/4120: Graduate/Undergraduate Natural Language Processing, Fall 2023

Instructor: Prof. Felix Muzny<sup>1</sup> (pronunciation: "Muse-knee"; pronouns: they/them and he/him<sup>2</sup>)

Contact: [f.muzny@northeastern.edu](mailto:f.muzny@northeastern.edu)

Office: Meserve 307A

TAs:

Ankit Ramakrishnan : he/him; UNK-ith

Harshitha Somala: she/her

Nidhi Bodar: she/her

Contacting course staff: we prefer that you pose all general questions on the [course piazza](#)

Piazza sign-up: <https://piazza.com/northeastern/fall2023/cs41206120merged202410>

Contacting Prof. Felix: reach out via email or before/after class. Please do not message them via teams.

Credit: 4 credits, Mondays & Thursdays, 11:45am - 1:25pm, Dodge Hall 050

Zoom link for synchronous remote lectures: go to "zoom meetings" on the left-hand side of our Canvas course

Office Hours: see course website

Lecture Format: Participating in synchronous course content (lectures, office hours) will be important to your overall learning and is **expected**. If you need to attend lecture remotely on any given day, fill out the remote attendance form (available on the course website) and attend via Zoom. If you are attending remotely, do expect to be actively participating in all course content and collaborating with your peers both during lecture and activities. Not all lectures will have a remote option—special events and days will be in person only, as marked on the course calendar.

I will be masking for a minimum of the first two weeks of the semester and would appreciate you to do the same.

As an on-ground class, it is expected by the university that the primary mode of instruction is in-person, on-ground. Do what you need to do to be fed/hydrated/prepared to be present in lecture.

---

*The world has changed dramatically in certain ways in our "post-COVID" world. There are a few lessons from that time that I'd like to emphasize:*

- *Life happens, **communicate** as proactively as you can and we will work together to make a plan.*
  - *Practicing proactive communication is a valuable skill that will help you in your post-university life.*
- *If you are sick in any way, do not come to class in person, even if you have received a negative covid test. If you are well enough to attend class remotely, this is the best alternative option.*
  - *I will do the same, as will our TAs.*

---

<sup>1</sup> Call them "Felix", "Professor Muzny", or "Professor Felix"

<sup>2</sup> This means that they are happy with either they/them or he/him or a mixture of both sets of pronouns

---

This is a shared syllabus for both CS 4120 and CS 6120. Specifics for 4120 will be noted in **green times new roman**, specifics for 6120 will be noted in **blue consolas**.

## Course Overview

NLP is about getting computers to perform useful and interesting tasks involving spoken and written human language. NLP is sometimes referred to as Computational Linguistics to emphasize the fact that it involves the combination of CS methods with research insights from Linguistics (the study of human language). Practical applications of NLP include question answering, machine translation, information extraction, and interactive dialog systems (both written and spoken). Modern NLP systems rely heavily on methods involving probability, linear algebra, and calculus --- often in combination with machine learning methods.

We'll be exploring both applications and the computational methods behind them. You should be prepared to get your hands dirty in terms of the math, programming, and data that comprise the behind the scenes components of NLP systems.

Expect a heavy emphasis on data, how systems are deployed in the real world, and engaging in critical thought about provided reading materials.

NLP is not *just* generative AI, and you can expect a significant fraction of this course to inspect other aspects of NLP (even though we'll cover generative AI models like ChatGPT, GPT-3, GPT-4, etc).

### Course Goals

1. Develop an understanding of the general problems that people who work on NLP study and the strategies they use to solve them.
2. Understand the role of data, machine learning, and neural networks in NLP systems.
3. Understand the ethical considerations and potentials for bias in NLP systems.
4. Be able to implement models to solve some "standard" NLP problems.
5. Be able to formulate potential starting points given a new problem with NLP elements.
6. Understand some of the motivating linguistic phenomena that make NLP problems hard and why these can be hard phenomena for computers to approach.

### Topics

- Language models, large and small
- Probabilistic (non-neural) methods & models
- Issues of ethics & bias in NLP
- Data sets and characteristics
- Words, word counting, lexicons
- Text classification with language models
- Text classification with single layer neural networks
- Vector semantics & word embeddings
- Part-of-speech tagging
- Viterbi algorithm & dynamic programming
- Machine translation
- Speech-to-text and text-to-speech\*
- Information extraction\*

- Question answering systems\*
- Dependency parsing\*

\* if time allows

## Textbook & Course Configuration

1. We'll be using draft chapters from the 3rd Edition of *Speech and Language Processing* by Dan Jurafsky and James H. Martin. You don't need to buy the current edition, draft pdfs of the new chapters are available from [the textbook website](#). You can also download (and print, if you desire) the entire book from the website. We will also link a pdf of this text from the course website.
  - a. We will be using the draft version from [January 7th, 2023](#).
2. We will supplement this text occasionally with readings from:
  - a. [Eisenstein, Jacob. \*Introduction to Natural Language Processing\*. MIT Press, 2019.](#)
  - b. [Kohn, Philipp. "Neural Machine Translation." arXiv preprint arXiv:1709.07809, 2017.](#)
  - c. (and others, to be linked from the course website)

## Websites & Technology

Make sure that you have access to all of the following websites and software:

- **Canvas:** We'll be using Canvas for some quizzes and links to homework submissions. Homework submissions will be done through Gradescope.
- **Gradescope:** Gradescope is where you will submit homework and some quizzes. You will also see your grades, feedback, and submit regrade requests via gradescope. You can find the link to Gradescope on Canvas.
- **Piazza:** This is our course discussion forum. This is where we will discuss relevant topics and answer your homework and content questions that come up outside of class. If you send us a content question via email, we'll likely ask you to post it to Piazza instead!
- **Python 3:** We'll be writing our homework coding assignments using python 3. Many assignments will ask that you are familiar with object-oriented programming in python and how to run unit tests on standalone .py files.
  - For help with python environment management: post your problem (screenshots are always encouraged) on piazza or come to office hours
  - If you have less python experience: we're holding a special python review session on Friday, September 8th, and make sure to come to office hours in the first week
- **Jupyter Notebooks:** Many of the coding activities that we complete in class will be distributed as Jupyter Notebooks. You can install jupyter notebooks either by installing [Anaconda](#) or via the [command line](#).
- **IDEs:** You can develop your code using whatever your preferred IDE is. For some coding assignments, **you might submit all or part of your solutions as .py files** (make sure to run any converted-from-jupyter-notebooks-.py-files again before you turn them in)! If you installed Anaconda, it comes with Spyder, which is an IDE that can be used to write and run .py files.
  - If you have less experience working with python and **both** Jupyter Notebooks and .py files, we **\*highly\*** encourage you to make time to come to office hours in the first few weeks of the course.

## Classroom Environment & Expectations

- **Preparation:** When there are readings assigned, it is the expectation that you do them before the first class meeting in the following week. This course will be a great opportunity for those of you who are interested in NLP & research to start flexing those muscles, and the best way for us to go down those paths is for you to develop a solid foundation.
- **Attendance:** You are expected to attend lecture *synchronously* whenever possible. **You are expected to attend *in person* whenever this is a reasonable choice.** We will be doing interactive activities during lecture as well as covering the material necessary for you to complete your homework and quizzes. Past students have reported the in class materials to be very valuable to their success in this course. See the grading rubric for notes on credit for in class activities.
- **Classroom environment:** It is unusually common in Computer science classes for some students to ask questions that are not really questions so much as opportunities to demonstrate knowledge of vocabulary or facts beyond the topic at hand. This can have a discouraging effect on other students who are not familiar with those terms, causing them to worry that they are less prepared to do well in the class (this is rarely the case—knowing terms outside the scope of the course is not a good predictor of success). If you find yourself wanting to make such a question or comment, please come talk to me about the topic after class or during office hours—I'm always happy to discuss tangentially related topics at those times! Please be particularly cognizant of using acronyms that we haven't yet talked about as a group.
- **Accommodation letters:** If you have an accommodation letter, please email it to me at your earliest convenience so that I can make sure this class is meeting your needs.
- **Name and pronouns:** If your name and pronouns are not in alignment with those listed on our class roster, please let me know either in person or via email so that I can ensure you are correctly addressed in this class.
  - If you wish to add, change, or update your pronouns in Canvas, go to "Account" > "Profile" > "Edit Profile", then add, change, or update your pronouns and display name.
  - If you wish to change or update your name here at Northeastern as a whole, find [instructions with the registrar here](#).
- **Class expenses:** If obtaining any material for use in our class presents a financial hardship for you, please let me know and I will work with you to locate the resources that you need to succeed in this class. You are **not expected to pay for compute time** for any components of this course. Please talk to me if you find yourself in a situation where you are thinking about doing this.
- **Feedback:** Please don't hesitate to reach out to me if any aspect of this course or class community could be improved.
- **Illness:** If you are ill in any way, you are expected to attend class remotely or make up the material asynchronously on your own and by attending office hours.

## Late Policy

All homework should be turned in on time whenever possible. All homework may be turned in up to 2 days (48 hours) late for a 10% penalty. For example, if homework is due on Wednesday at 9pm, it may be turned in as late as Friday at 9pm.

***Once a semester, you may turn in your homework up to 48 hours late without penalty or explanation. This will be automatically applied the first time you turn in your homework late. This policy only applies to homeworks and does not apply to the final project or quizzes, which cannot be turned in late.***

You will not receive credit/extra credit for not using your late pass.

## Extensions

Extensions for any work beyond the regular late policy will be given based on proactive communication with Prof. Felix. **Whenever possible, this should occur at least 24 hours before the posted deadline.** The sooner that you reach out, the easier this will be.

Email Felix ([f.muzny@northeastern.edu](mailto:f.muzny@northeastern.edu)) with the following information:

- 1) Which assignment are you requesting an extension on & why you are requesting an extension.
- 2) When are you requesting the extension until.
- 3) What is your plan for how this extension will impact the due dates for the other assignments in this course.

You don't need to write an essay, just be sure to include the above information. This extension policy is based on our mutual understanding that living during and recovering from a pandemic can be difficult, we're all doing our best, and the easiest way for you to succeed in this course is proactive communication. If a situation arises that makes it impossible to reach out 24 hours before the deadline, don't panic—send Prof. Felix an email when you can and we'll discuss your options together.

## AI Collaboration Policy

This is the AI collaboration policy produced as a result of our work together on the first day of term.

### **1. How should LLMs be allowed to be used in this course? Justify your answers.**

LLMs should be used to help understand concepts that need to be re-clarified. They should be used to supplement information from lecture, the instructional staff, fellow students, and the textbook. LLMs should not be viewed as a truthful source of information, by default (information given by them should be verified). It is the responsibility of the student to verify the output of the LLMs that they use.

Any code or text that was originally generated by an LLM should be cited as such if turned in as part of an assignment. (see example below)

Github co-pilot may be used to aid inline code completion and need not be cited.

### **2. How should LLMs not be allowed to be used in this course? Justify your answers.**

LLMs cannot be used to generate end-to-end code or text answers to homework assignments or tasks. The output of an LLM should not be passed off as your own work. LLMs should not be used to generate the majority of any assignment or subtask in an assignment, unless specifically directed to do so in the instructions.

3. On a scale from 1- least permissive (no large language models) to 5 - most permissive (large language models are always okay), what level of permissiveness does your group want for this course?

3.8

4. Translate your answer from question #3 into a single sentence.

**3.0 - 3.9** - Medium-permissive policy. Since this is a class on NLP and we will be learning about GPTs, it makes sense for there to be a certain level of permissiveness. We should not expect to get assignment answers from LLMs. We should be allowed to use gpt-style tools for most but not all areas of the course with creativity, adapting their responses and citing their usage. LLMs should be used to help understand content rather than solve homework problems.

**Example:**

**Question given:** Explain the effects of tokenization in NLP applications and give an example tokenize function in python.

**Answer submitted:**

**LLM(s) used:** ChatGPT

**prompt(s):** "what is tokenization", gave an overview of tokenization, "show an example of splitting a string on periods in python", "show an example of splitting a string on periods and spaces in python"

Tokenization is the process of breaking an input text into the units (generally words or smaller) that will serve as input to our applications. Different languages have different corner cases that must be considered. How we tokenize has direct impact on what NLP applications are able to do—this affects what vocabulary the application is considering and therefore both what it is able to evaluate and often what it is able to produce.

```
import re

def tokenize(text: str) -> list:
    # begin ChatGPT
    # Split the string on periods and spaces using a regular expression
    split_string = re.split(r'[\.\s]', input_string)
    # end ChatGPT

    # we likely want to do something about contractions too
    # the rest of the code
    return my_fully_split_string
```

## Collaboration Policy

The work that you turn in should be your own. We encourage you to collaborate with your classmates, but remember that collaboration looks very different than working on a pair or group project.

Here are three big-picture points to remember when collaborating with your classmates:

- **Strategies:** You may talk with your classmates about *general strategies* but you may not talk about *specific solutions*.

- **Explaining concepts:** You may talk with your classmates about how certain techniques work *in general* but not how to write any part (or sub-part) of the solution needed for the homework.
- **A good rule of thumb:** don't show your assignments to other people; don't look at other people's assignments (this makes it very hard to come up with your own solution afterwards); don't write code together unless the assignment explicitly states that you may work in groups. This includes working through solutions on whiteboards as well as talking through the solutions verbally.

You are expected to use the internet as a place for online resources, such as documentation, not as a place to get solutions to your assignments.

The finer-grained details:

- **Do not search for a solution online:** You may not actively search for a solution to the problem from the internet. This includes posting to sources like StackExchange, Reddit, Chegg, etc.
  - **StackExchange Clarification:** Searching for basic techniques in python is fine. If you want to post and ask "How do convert a float to an integer" that's fine. What you **cannot** do is post things like "Here's the function my prof gave me to write. I need to convert this temperature in celcius to fahrenheit".
- **Plagiarism:** assignments and code should be your own. You should not need to consult sources beyond the class notes, posted lecture notes, examples, and resources, and python and its associated libraries' documentation. Make sure that you understand how this relates to the AI collaboration policy for this course.
  - We will practice properly documenting other resources that you consult throughout the course of this class.
- **Tutors:** you should always consult the course instructional staff if you need extra help. They are here specifically to help you! You should never have anyone else write code for you. This includes tutors, friends, strangers, friends of friends, or anyone who is not you (even friends who have taken this course previously!).
- **When in doubt, ask:** If you have doubts about this policy or would like to discuss specific cases, please ask Prof. Felix.

A single Collaboration Policy violation will result in a 0 on the assignment in question that may not be dropped. Multiple violations will result in an F in the course.

Violations will result in being reported to the appropriate university-level committee. For graduate students, consequences of violating collaboration policies vary but often result in being barred from being eligible to TA for the duration of your time at Northeastern in addition to other consequences.

The university's academic integrity policy discusses actions regarded as violations and consequences for undergraduate students in particular. <http://www.northeastern.edu/osccr/academic-integrity>

## Grading & Assignments

If you enrolled in this course after deadlines have passed **and** you contact Prof. Felix ASAP, we will work with you to adjust deadlines as needed.

Category	Due Dates & Points	Grade Percentage
Homework	Due on Wednesdays at 9pm.	50%

	<p>No homework grades will be dropped. Homeworks are graded on correctness and depth of interaction with the material. Students in CS 6120 will be completing a small number of additional problems for each homework assignment. Students in CS 4120 will not be required to do these problems, nor will they earn extra credit for completing them (though you are always welcome to attempt them and won't lose points for any incorrect attempts.)</p> <p>Homeworks cannot be updated and re-submitted after receiving your grade.</p>	
Labs	<p>Labs are due on Fridays at 9pm.</p> <p>Labs will consist of problems from the in-class lab work plus a few additional problems. The additional problems may address reading material, lecture material, or a combination of both.</p> <p>Labs <b>may not</b> be turned in late.</p> <p>All labs may be updated and re-submitted within two weeks of receiving each lab grade.</p>	20%
Research presentations	<p>Students in CS 6120 will be reading and presenting one research paper either as individuals or in groups of two.</p> <p>This will count as one, additional, homework grade incorporated into the homework bucket.</p> <p>Papers and presentation schedule will be posted on the course website.</p> <p>Students in CS 4120 who wish to complete this activity must seek special permission from Prof. Felix by September 12th.</p>	*
Final Project	<p>Final projects may be completed individually or in groups of up to three people.</p> <p>Instructions for Final Projects will be posted on September 28th.</p>	30%



Final grades will be based on the following scale. Decimals will be rounded to the nearest integer.

<b>Letter Range</b>	
A	95 - 100
A-	90 - 94
B+	87 - 89
B	83 - 86
B-	80 - 82
C+	77 - 79
C	73 - 76
C-	70 - 72
D	60 - 69*
F	< 60

\* graduate students receive an F if their grade is 69 or under

## Calendar

The course calendar will be (subject to change at the instructor's discretion): [cs4\\_6120\\_calendar\\_f23](#)

## Classroom Recording

This course, or parts of this course, will be recorded for educational purposes. These recordings will be made available only to students enrolled in the course, instructor of record, and any teaching assistants assigned to the course.

If you have objections or would like to opt-out of recordings, please contact the instructor.

Only students who have arranged an accommodation with the Disability Resource Center may use mechanical or electronic transcribing, recording, or communication devices in the classroom. Students with disabilities who believe they may need such an accommodation may contact the Disabilities Resource Center.

## Accommodations

It is my job to create a classroom environment that is most conducive to you learning well. If you have accommodations from the [Disability Resource Center](#), please provide your letter to me early in the semester so that I can arrange for these accommodations. If you wish to receive accommodations and do not have a letter, please visit the DRC at 20 Dodge Hall or call (617) 373-2675.

## Student Names and Pronouns

We recognize that your legal information doesn't always align with how you identify. Students may update their first and middle names as well as gender marker [with the registrar](#), even if they are not your legal names or gender marker. Those names and gender marker are what would appear publicly in most university systems.

In the absence of such updates, what we see on most university systems by default are your legal name and gender marker.

## Classroom Environment

To create and preserve a classroom atmosphere that optimizes teaching and learning, all participants share a responsibility in creating a civil and non-disruptive forum for the discussion of ideas. Students are expected to conduct themselves at all times in a manner that does not disrupt teaching or learning. Your comments to others should be constructive and free from harassing statements. You are encouraged to disagree with other students and the instructor, but such disagreements need to be respectful and be based upon facts and documentation (rather than prejudices and personalities). The instructor reserves the right to interrupt conversations that deviate from these expectations. Repeated unprofessional or disrespectful conduct may result in a lower grade or more severe consequences.

Part of the learning process in this course is respectful engagement of ideas with others.

The [Code of Student Conduct can be found on the OSCCR website](#).

## Title IX

Title IX of the Education Amendments of 1972 protects individuals from sex or gender-based discrimination, including discrimination based on gender-identity, in educational programs and activities that receive federal financial assistance.

Northeastern's Title IX Policy prohibits Prohibited Offenses, which are defined as sexual harassment, sexual assault, relationship or domestic violence, and stalking.

The Title IX Policy applies to the entire community, including male, female, non-binary, and transgender students, faculty and staff.

If you or someone you know has been a survivor of a Prohibited Offense, confidential support and guidance can be found through [University Health and Counseling Services](#) staff and the [Center for Spiritual Dialogue and Service](#) clergy members.

By law, those employees are not required to report allegations of sex or gender-based discrimination to the University.

Reports can be made non-confidentially to the Title IX Coordinator within the Office for Gender Equity and Compliance at: [titleix@northeastern.edu](mailto:titleix@northeastern.edu) and/or through NUPD (Emergency 617.373.3333; Non-Emergency 617.373.2121).

Reporting Prohibited Offenses to NUPD does NOT commit the victim/affected party to future legal action.

Faculty members are considered "responsible employees" at Northeastern University, meaning they are required to report all allegations of sex or gender-based discrimination to the Title IX Coordinator.

In case of an emergency, please call 911.

Please visit <http://www.northeastern.edu/titleix> for a complete list of reporting options and resources both on- and off-campus.

## Religious Holidays

The course staff will make every effort to deal reasonably and fairly with all students who, because of religious obligations, have conflicts with scheduled exams, assignments or required attendance. In this class, contact the course staff at least 7 days in advance of the conflicting date to reschedule a homework or quiz due date.