

# CSC 485H/2501H: Computational linguistics, Fall 2024 - Assignment 2

**Due date:** 17:00 on Thursday, November 7, 2024.

## General Instructions

- Late assignments will not be accepted without a valid medical certificate or other documentation of an emergency.
- For CSC485 students, this assignment is worth 25% of your final grade. For CSC2501 students, it is worth 10%.
- Read the whole assignment carefully.
- Type the written parts of the assignment.
- Your work must be your own. Do not copy code or text from the problems. If you need assistance, contact the teaching assistant.
- Any clarifications to the assignment will be posted on the course page regularly.
- The starter code directory contains the starter code for the assignment. Refer to code files in that directory.
- When implementing code, follow the starter code and the assignment instructions, including the implementation details.
- Fill in your name, student ID, and course number at the top of each file you submit.



## 0. Warming up with WordNet and NLTK (4 marks)

### (a) Deepest function (1 mark)

- Implement the `deepest` function in `q0.py` to find the synset in WordNet with the largest maximum depth and report both the synset and its depth on each of its paths to a root hyperonym. Hint: use `wn.all_synsets` and `synset.max_depth` methods.

## (b) Superdefn function (2 marks)

- Implement the `superdefn` function in `q0.py` that takes a synset `s` and returns a list consisting of all of the tokens in the definitions of `s`, its hyperonyms, and its hyponyms. Use `word_tokenize` as shown in chapter 3 of the NLTK book.

## (c) Stop\_tokenize function (1 mark)

- Implement the `stop_tokenize` function in `q0.py` that takes a string, tokenizes it using `word_tokenize`, removes any tokens that occur in NLTK's list of English stop words and also removes any tokens that consist entirely of punctuation characters. Use Python's punctuation characters from the `string` module. Maintain the original case in the return value.

## 1. The Lesk algorithm (3 marks)

### (a) Mfs function (1 mark)

- Implement the `mfs` function in `q1.py` that takes a word and a sentence and returns the most frequent sense of the word in the sentence. Note that `wordnet.synset` returns a `Synset` object.

### (b) Lesk function (1 mark)

- In the `lesk` function in `q1.py`, implement the Lesk algorithm as specified in Algorithm 1, including the calculation of the overlap between the intersection of the bags of words in the definition and examples of two senses. The function returns the cardinality of the overlap. Use `stop_tokenize` to tokenize the input.

### (c) Lesk\_ext function (3 marks)

- In the `lesk_ext` function in `q1.py`, implement a version of Algorithm 1 where the `signature` also includes the words in the definition and examples of sense's hyponyms, holonyms, and meronyms. Use `stop_tokenize` as before.

### (d) Justification for lesk\_ext (2 marks)

- Explain why the extension in `lesk_ext` is helpful. Consider the likely sizes of the overlaps.



### (e) Lesk\_cos function (4 marks)

- In the `lesk_cos` function in `q1.py`, implement a variant of `lesk_ext` that uses `CosSim` instead of `Overlap`. Modify `signature` and `context` to be vector-valued and construct the vectors as described. Use `stop_tokenize` to get the tokens for the signature.

### (f) Lesk\_cos\_oneside function (2 marks)

- In the `lesk_cos_oneside` function in `q1.py`, implement a variant of `lesk_cos` that, when constructing the vectors for the signature and context, does not include words that occur only in the signature. Use `stop_tokenize` to get the tokens for the signature.

### (g) Comparison of `lesk_cos_oneside` and `lesk_cos` (3 marks)

- Compare how well `lesk_cos_oneside` and `lesk_cos` perform. Justify your answer with examples.

### (h) CosSim with k

- If we use `CosSim` for vector similarity, how is it related to the set intersection? (No implementation required)

### (i) Lesk\_w2v function

- In the `lesk_w2v` function, implement a variant of `lesk_ext` where the vectors for the signature and context are `Word2Vec` vectors for the words in the signature and set. Use `stop_tokenize` to get the tokens for the signature.

### (j) Lowercasing tokens (2 marks)

- Alter your code so that all tokens are lowercased before they are used for any of the comparisons, vector lookups, etc. Analyze how this alters the different methods' performance and explain why. Do not submit this lowercased version.



## 2. Word sense disambiguation with BERT (22 marks)

### (a) Context necessity (4 marks)

- Is context really necessary? Give an example of a sentence where word order-invariant methods such as those implemented for Q1 will never be able to completely disambiguate. Explain the more general pattern and why these methods cannot provide the correct sense for each ambiguous word.

### (b) Gather\_sense\_vectors function (10 marks)

- Implement `gather_sense_vectors` in `q2.py` to assign sense vectors as described.

### (c) Sorting corpus (4 marks)

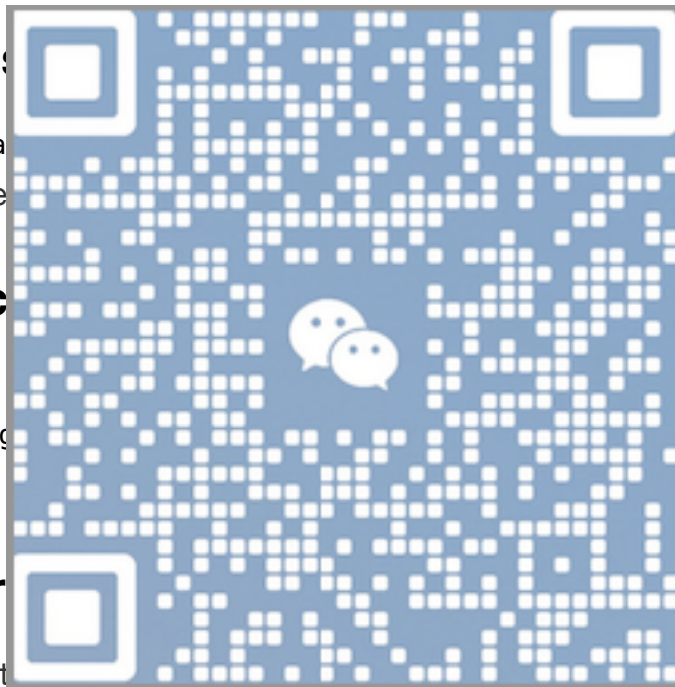
- In the docstring for `gather_sense_vectors`, explain why sorting the corpus by length before batching is much faster.

### (d) Bert\_1nn function (4 marks)

- Implement `bert_1nn` to return the sense vector for the word in the sentence given sense vectors produced by `gather_sense_vectors`. Refer to the docstring about loop usage.

### (e) Issues with ar (4 marks)

- Think of at least one other method to use the code for this assignment to disambiguate arbitrary sentences. Consider either the Lesk variants from Q1 or the BERT-based method here (or both).



## 3. Understanding transformers through causal tracing (16 marks)

### (a) Get\_forward\_hooks function (3 marks)

- Implement `get_forward_hooks`.

## (b) Causal\_trace\_analysis function (5 marks)

- Implement `causal_trace_analysis` to compute the impact of states, MLP and attention.

## (c) Causal tracing result report (1 mark)

- Report your generated causal tracing result plots for the prompt “The Eiffel Tower is located in the city of” with the output “Paris” in your report.

## (d) Model size impact on causal tracing (3 marks)

- Experiment with different sizes of GPT-2 models (e.g., small, medium, large, and XL) to examine how model size impacts causal tracing patterns. Address the following in your report:
  - At what model size does one observe that the causal tracing pattern no longer appears?
  - Discuss potential reasons for the change in tracing patterns as the model size increases.

## (e) Prompt types

- Using GPT-2 XL, experiment with different prompt types that result in a causal tracing pattern similar to the one observed in the report. Discuss your findings with examples and discuss what characterizes this similarity.

## (f) Absent or diminished causal tracing patterns (2 marks)

- For GPT-2 XL, explore prompt types where the causal tracing pattern is absent or significantly different. Discuss any trends or patterns you identified and reflect on the broader implications of how language models process, store and generate factual information obtained from pretraining.

## What to submit

- Submit electronically via MarkUs.
- Submit a total of five required files:
  - `a2written.pdf` : a PDF document containing answers to questions 0a, 1d, 1f, 1h, 2a, and 2d. Also include a typed copy of the Student Conduct declaration and sign it by typing your name.
  - `q0.py` : the entire file with your implementations filled in.