McMaster University is located on the traditional Territories of the Mississauga and Haudenosaunee Nations, and within the lands protected by the "Dish With One Spoon" wampum agreement.

# Session Recording and Privacy

*This session is being recorded with the intention of being shared publicly via the web for future audiences. In respect of your privacy, participant lists will not be shared outside of this session, nor will question or chat transcripts.*

*Questions asked via the chat box will be read by the facilitator without identifying you. Note that you may be identifiable when asking a question during the session in an audio or visual format.*

***Closed Captioning*** *is available for this session—please click the "CC" button at the bottom of the Zoom screen to turn these on. Please feel free to reach out to* [scds@mcmaster.ca](mailto:scds@mcmaster.ca) *with any other access requests.*

Lewis & Ruth
**Sherman Centre**
for Digital Scholarship
scds.ca

McMaster University | Library

# Code of Conduct

*The Sherman Centre and the McMaster University Library are committed to fostering a supportive and inclusive environment for its presenters and participants.*

*As a participant in this session, you agree to support and help cultivate an experience that is collaborative, respectful, and inclusive, as well as free of harassment, discrimination, and oppression. We reserve the right to remove participants who exhibit harassing, malicious, or persistently disruptive behaviour.*

*Please refer to our code of conduct webpage for more information:*
*https://scds.ca/events/code-of-conduct/*

# Certificate Program

*The Sherman Centre offers a Certificate of Completion that rewards synchronous participation in 7 workshops. We also offer concentrations in Data Analysis and Visualization, Digital Scholarship, and Research Data Management.*
*Learn more about the Certificate Program: https://scds.ca/certificate-program*

# Attendance Confirmation

*If you would like to be considered for a certificate, verify your participation in today's workshop by completing the form at: https://u.mcmaster.ca/verification*

*An organizer will enter the code into the session chat window.*

Lewis & Ruth
**Sherman Centre**
for Digital Scholarship
scds.ca

McMaster University | Library

# Outline & Schedule

| Segment | Time Allotted | Key Topics / Activities |
|---|---|---|
| **Introductory remarks** | 20 minutes | Introduction to text preparation and analysis<br>Overview of concepts and methods<br>Key considerations for different source materials and analyses |
| **Named Entity Recognition** | 35 minutes | Introduction to Google Colab & Jupyter Notebooks<br>Get the data<br>Introduction to NER and hands-on exercise |
| **--Break--** | **10 minutes** | **--Break--** |
| **Sentiment Analysis** | 30 minutes | Introduction and hands-on exercise<br>Constellate demonstration |
| **Topic Modeling** | 30 minutes | Introduction and hands-on exercise |
| **Q & A; Final thoughts (lecture + discussion)** | 20 minutes | Questions & final thoughts<br>Where to learn more |

Workshop landing page: scds.github.io/dmds23-24/textanalyses.html

# Learning Objectives

**By the end of this module, you will be able to:**

- List the common methodological approaches used in text preparation and analysis and identify when and how to use them based on source materials and analysis objectives.

- Apply prepared computational techniques to perform common text preparation steps and introductory analyses.

- Identify resources and tutorials for further learning and analyses.

Workshop landing page: scds.github.io/dmds23-24/textanalyses.html

# An Introduction

# Natural Language Processing
# (is a big family)

**Text and speech recognition / processing**
OCR, speech recognition, text-to-speech

**Lexical semantics**
Named entity recognition, Sentiment analysis, word sense disambiguation

**Morphological analysis**
Stemming, Lemmatization, Part of Speech Tagging

**Relational semantics**
Relationship extraction, Semantic parsing

**Syntactic analysis**
Parsing, Sentence breaking

**Discourse semantics**
Discourse analysis, Topic segmentation, Argument mining

# NLP Workflows

Collection

Initial analysis / review

Preprocessing

Exploration / Analyses

Results / Output

Dissemination

# NLP Workflows

Collection

Initial analysis / review

Preprocessing

Exploration / Analyses

Results / Output

Dissemination

**Common sources:**

Downloaded / scraped born-digital text
Digitized + OCRed analog text
Transcribed audio (human or machine)

# NLP Workflows

Collection

Initial analysis / review

Preprocessing

Exploration / Analyses

Results / Output

Dissemination

**Common activities:**

Identify spelling / OCR error patterns
Evaluate OCR accuracy
Identify unnecessary parts
Check metadata
Inspect text formatting / encoding

# NLP Workflows
## … *are iterative*

# NLP Workflows
## *... are iterative*

# NLP Workflows
## … *are iterative*

# Text preparation and analysis are task specific

Your approaches should be informed by:

1. **Your analysis objectives**

2. **Your source materials and their common traits, inconsistencies, errors**

3. **Your abilities, time, interests, and familiarity with tools**

# Considerations

**1. Your analysis objectives**

- Do you have a defined research question or are you experimenting?

- What analyses are required to meet your objectives and create desired outputs?

- Are your methods sensitive to particular types of errors and imprecision?

- For which applications were the methods developed? How were they trained/validated? Are they appropriate for your purposes?

# Considerations

## 2. Your source materials and their common traits, inconsistencies, errors

- Born-digital vs. digitized
- The quality of the source materials
- The methods used to digitize materials and create text
- The structure of the materials and the text within
- The nature of communication within the materials
- Which (if any) processing operations can be automated?

# Considerations

## 3. Your abilities, time, interests, and familiarity with tools

- With which tools are you familiar? Do feasible solutions exist within those?

- How much time and interest do you have to learn new approaches and tools?

- Do you have time to explore, test, and iterate?

- Can you apply your acquired knowledge & workflows to future projects?

# Getting started:
Google Colab & Jupyter notebooks
Get your data

Go to u.mcmaster.ca/dmds-text-2324 to download your data and view the workshops for today's workshop.
Follow along with Devon's instructions

# Orientation to Google Colab / Jupyter Notebooks

Named Entity Recognition

# Analyzing Texts with Named Entity Recognition

Four months `DATE` after she had gone to Paris `GPE` , Mary Wollstonecraft `PERSON` met at the house of a merchant, with whose wife she had become intimate, an American `NORP` named Gilbert Imlay `PERSON` . He won her affections. That was in April, 1793 `DATE` . He had no means, and she had home embarrassments, for which she was unwilling that he should become in any way responsible. A part of the new dream in some minds then was of a love too pure to need or bear the bondage of authority. The mere forced union of marriage ties implied, it was said, a distrust of fidelity. When Gilbert Imlay `PERSON` would have married Mary Wollstonecraft `PERSON` , she herself refused to bind him; she would keep him legally exempt from her responsibilities towards the father, sisters, brothers, whom she was supporting. She took his name and called herself his wife, when the French Convention `ORG` , indignant at the conduct of the British Government `ORG` , issue a decree from the effects of which she would escape as the wife of a citizen of the United States `GPE` . But she did not marry. She witnessed many of the horrors that came of the loosened

# Named Entity Recognition (NER) in Practice

# How Named Entity Recognition (NER) Works

**Training dataset**

"..and soon the white walls and flowery garden of Fort William, the Hudson Bay Company's trading post. The rockery in the centre of the garden would have gladdened the heart of an Ontario gardener. I believe that wealthy people there have had large fragments of Lake Superior rock brought down to adorn their lawns and gardens. We found friends at the fort in the factor and his family, with whom we spent a pleasant half-hour. Mr. McIntyre is well known, and many will owe him gratitude for kindness as long as Fort William or the Canada Pacific Railway remains in their memory."



"bell hooks"

# Try it out in Jupyter Notebooks... Make a copy!

+ Code   + Text

```
[ ]  # Import Counter to count named entities
     from collections import Counter

     # Import SpaCy library
     import spacy
     from spacy import displacy

     # Import matplotlib.pyplot to create bar graph
     import matplotlib.pyplot as plt
```

```
[ ]  # Assign the filename to a variable
     filename = 'wollstonecraft.txt'

     # Make the text of the file available to our script
     ner_text = open(filename).read()
```

```
[ ]  # Instantiate NLP pipeline — load transformer corpus
     nlp = spacy.load('en_core_web_trf')

     # For faster but less accurate results, you can use nlp = spacy.load('en_core_web_sm')

     # Create the Doc object by passing it through the text pipeline (nlp)
     doc = nlp(ner_text)
```

```
[ ]  for ent in doc.ents:
         print(ent.text, ent.start_char, ent.end_char, ent.label_, spacy.explain(ent.label_))
```

# Interpreting the Results

- **PERSON** - People (including fictional ones)
- **NORP** - Nationalities, or religious or political groups
- **GPE** - Geopolitical Entity, e.g. city, country, states
- **LOC** - Non GPE locations, mountain ranges, bodies of water
- **FAC** - Buildings, airports, highways, bridges, etc.
- **ORG** - Companies, agencies institutions
- **EVENT** - battles, wars, sports events, etc.

# Sentiment Analysis

# Sentiment Analysis

- Go to our shared materials for this workshop:
  u.mcmaster.ca/dmds-text-2324

- Open the file **sentiment-analysis-2324.ipynb** and save a copy to your Google Drive.

- Follow along with Jay

# Topic Modeling

# Discerning Corpus "Topics" with Topic Modeling

# Try it out in Jupyter Notebooks... Make a copy!

+ Code   + Text

```
[ ]  # Install pyLDAvis with pip for visualization
     !pip install pyLDAvis
```

```
[ ]  # Import internal libraries: glob for grabbing docs from directory
     import glob

     # Import external libraries: gensim for preprocessing and LDA
     import gensim
     import gensim.corpora as corpora
     from gensim.utils import simple_preprocess
     from gensim.models import CoherenceModel

     # Import external libraries: spaCy for lemmatization, NLTK for stopwords
     import spacy
     import nltk
     nltk.download('stopwords')

     # Import external libraries: pyLDA for vis
     import pyLDAvis
     import pyLDAvis.gensim_models as gensimvis
```

```
[ ]  # Read files from directory and create list from contents
     file_list = glob.glob('./russelltexts' + '/*.txt') # directory containing text (.txt) files

     texts = []

     for filename in file_list:
         with open(filename, mode = 'r', encoding = 'mac-roman') as f: # specify encoding as appropriate
             texts.append(f.read())
```

# Questions &
# Final thoughts

—

# Some final thoughts

- Begin with your goals in mind

- Experiment and iterate

- Understand your methods

- Start small and scale up

- Document your sources, methods, rationale, and outcomes **as you develop them**