



DMDS

February 16 | 1:30-4:30pm
Virtual Workshop

Computational Approaches to Text Prep & Analysis

u.mcmaster.ca/scds-events

Devon Mordell
Jay Brodeur

SCDS
■■■■

Library

McMaster
University 



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.

Laslovarga, “Webster Falls in Winter, Waterdown, Hamilton, Ontario, Canada - Spencer Gorge / Webster's Falls Conservation Area,” 23 January 2011, Wikimedia Commons - https://commons.wikimedia.org/wiki/File:Waterdawn_Webster_Falls_in_Winter8.jpg

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 with any other access requests.

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.

Outline & Schedule

Segment	Time Allotted	Key Topics / Activities
Introductory remarks	20 minutes	Introduction to text preparation and analysis Overview of concepts and methods Key considerations for different source materials and analyses
Text prep with OpenRefine	40 minutes	Introduction to OpenRefine Manual cleanup (e.g. find and replace) Faceting
Getting Programmatic with Python	20 minutes	Overview of programmatic approaches The 'what' and 'when' to program Using Python for text preparation
--Break--	10 minutes	--Break--
Hands-on sampling of text analysis methods	75 minutes	Named entity recognition Topic modeling Sentiment analysis
Q & A; Final thoughts (lecture + discussion)	10 minutes	Questions & final thoughts Where to learn more

Workshop landing page: scds.github.io/dmds-22-23/ComputationalText.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.
- Explain the benefits and challenges of applying a scripted or semi-scripted approach to text preparation and analysis; identify situations where scripting your work will be beneficial.
- Apply prepared computational techniques to perform common text preparation steps and introductory analyses.

Workshop landing page: scds.github.io/dmds-22-23/ComputationalText.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

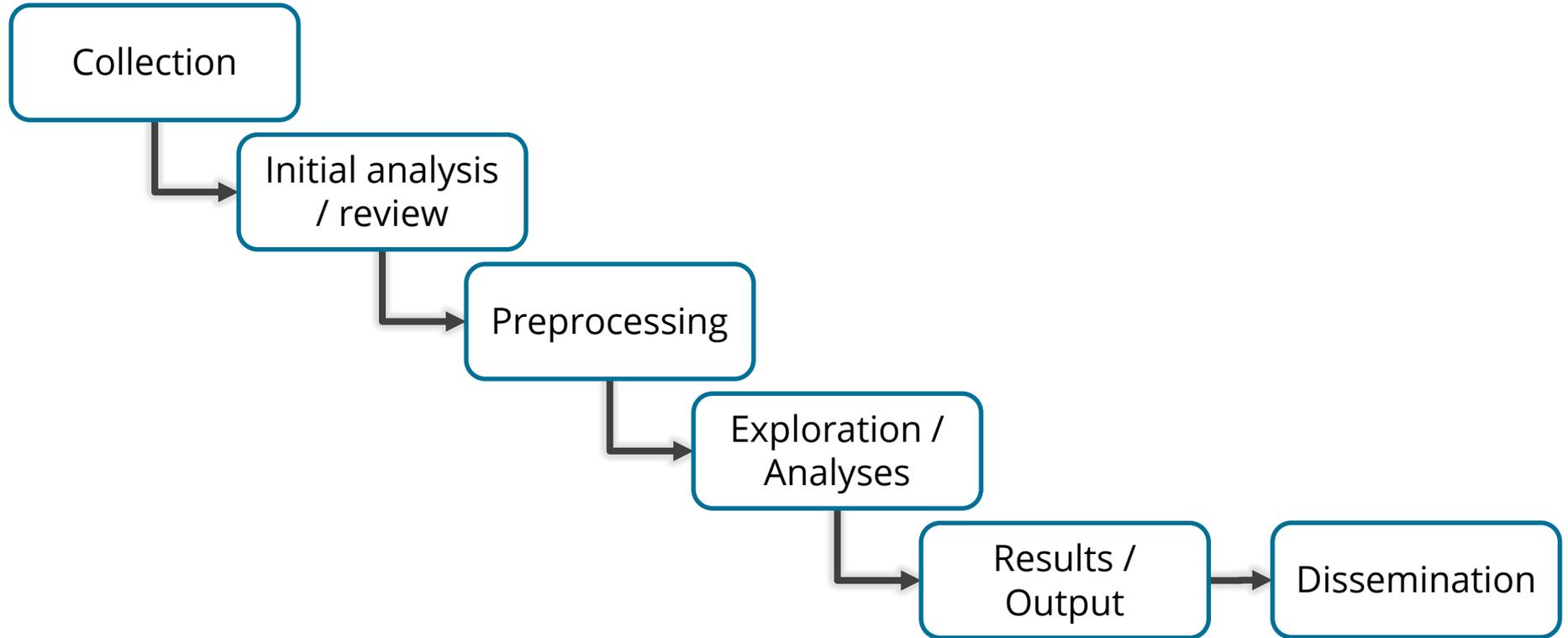
Syntactic analysis

Parsing, Sentence breaking

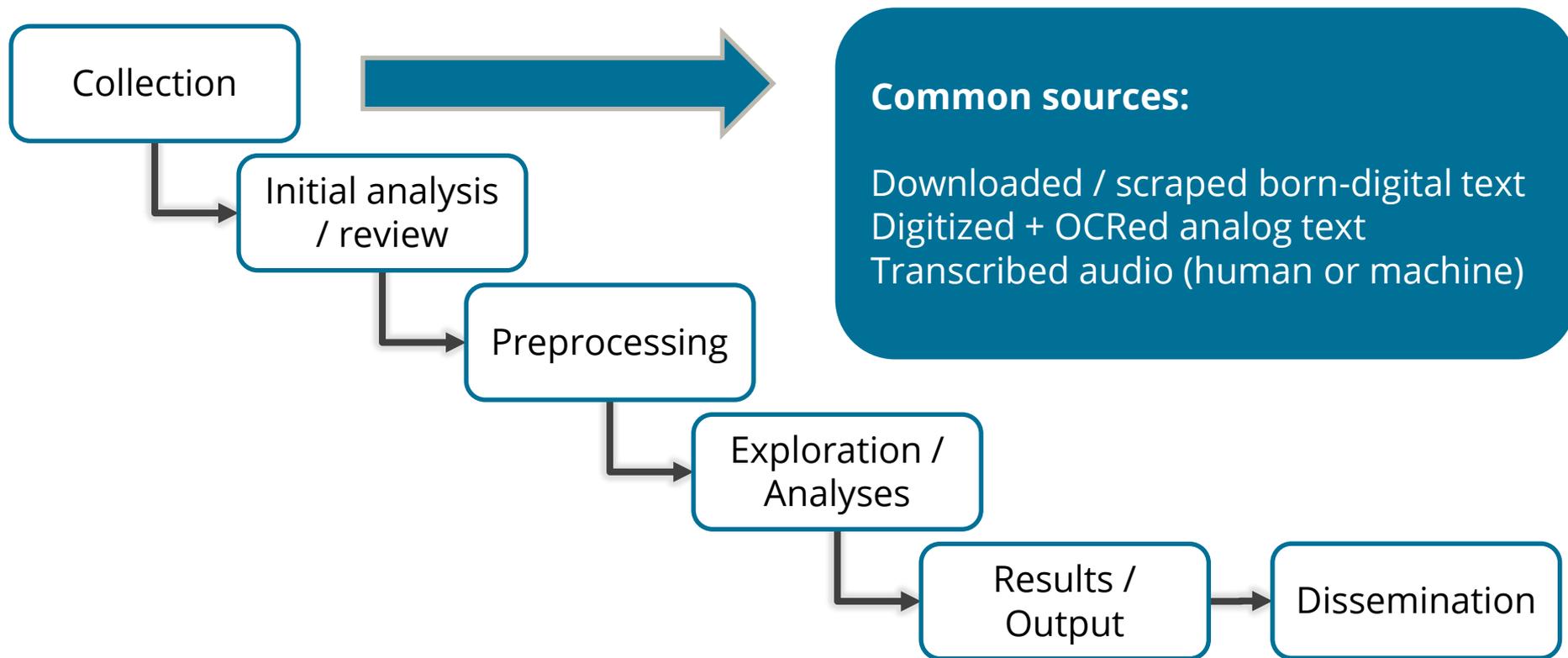
Discourse semantics

Discourse analysis, Topic segmentation,
Argument mining

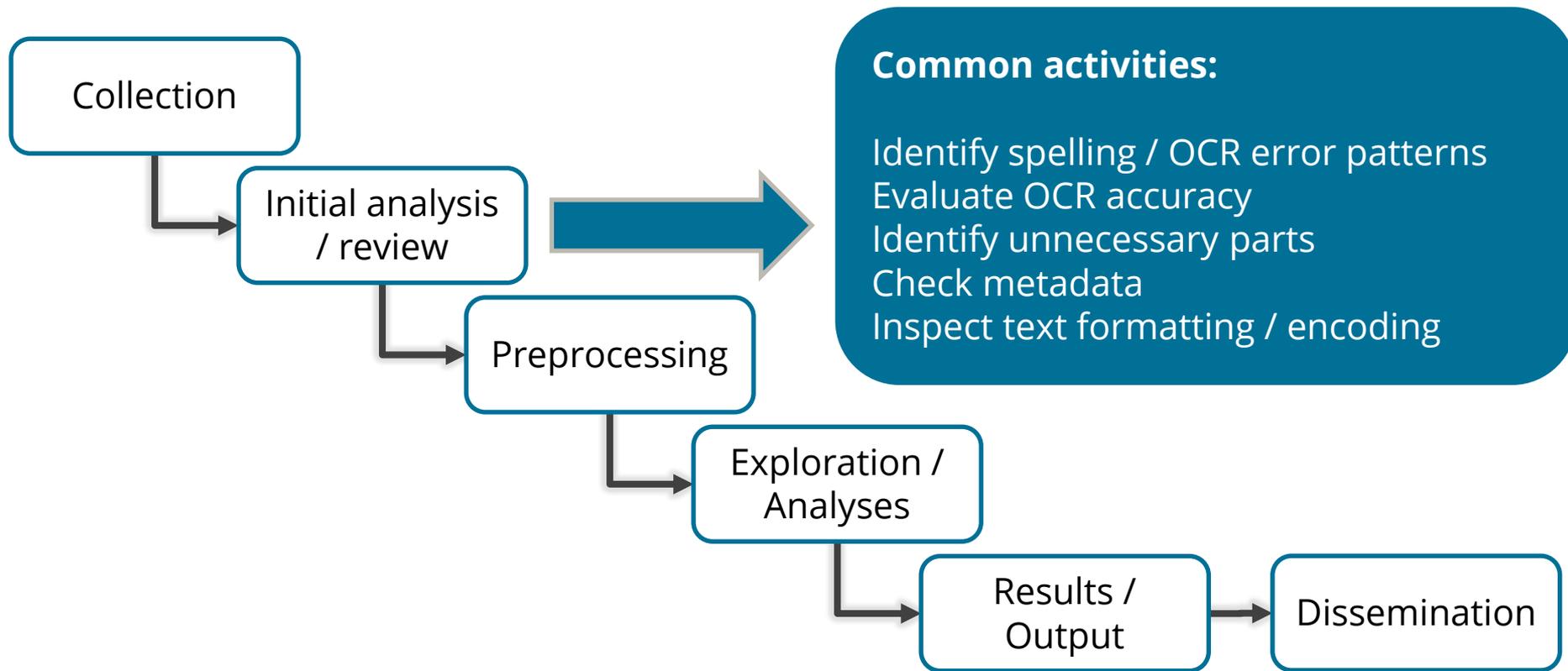
NLP Workflows



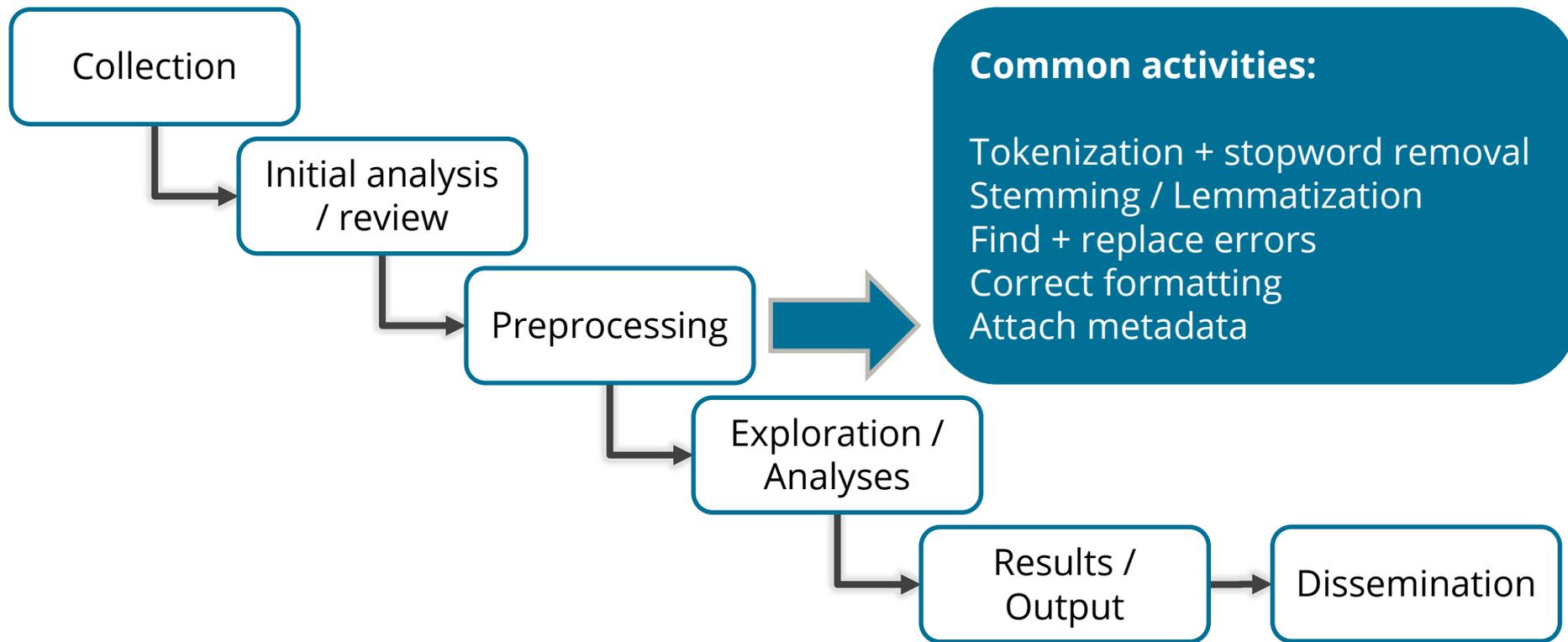
NLP Workflows



NLP Workflows

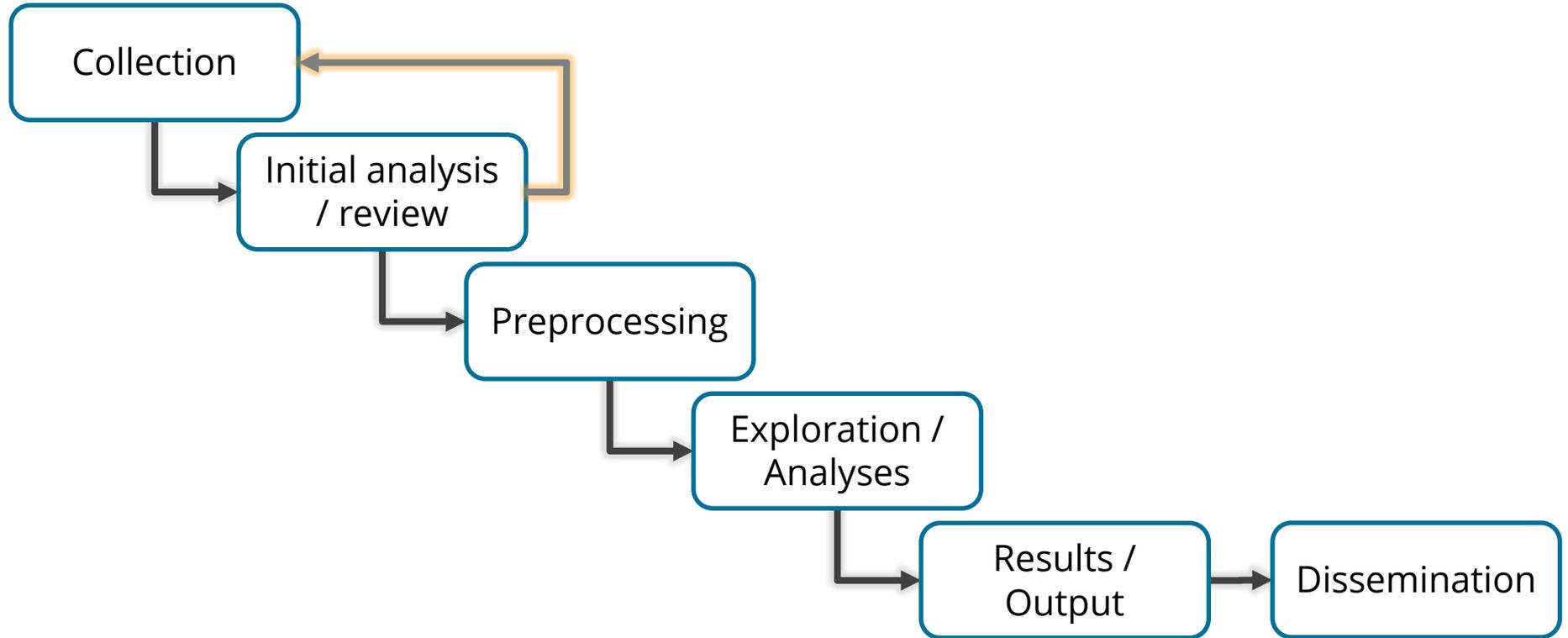


NLP Workflows



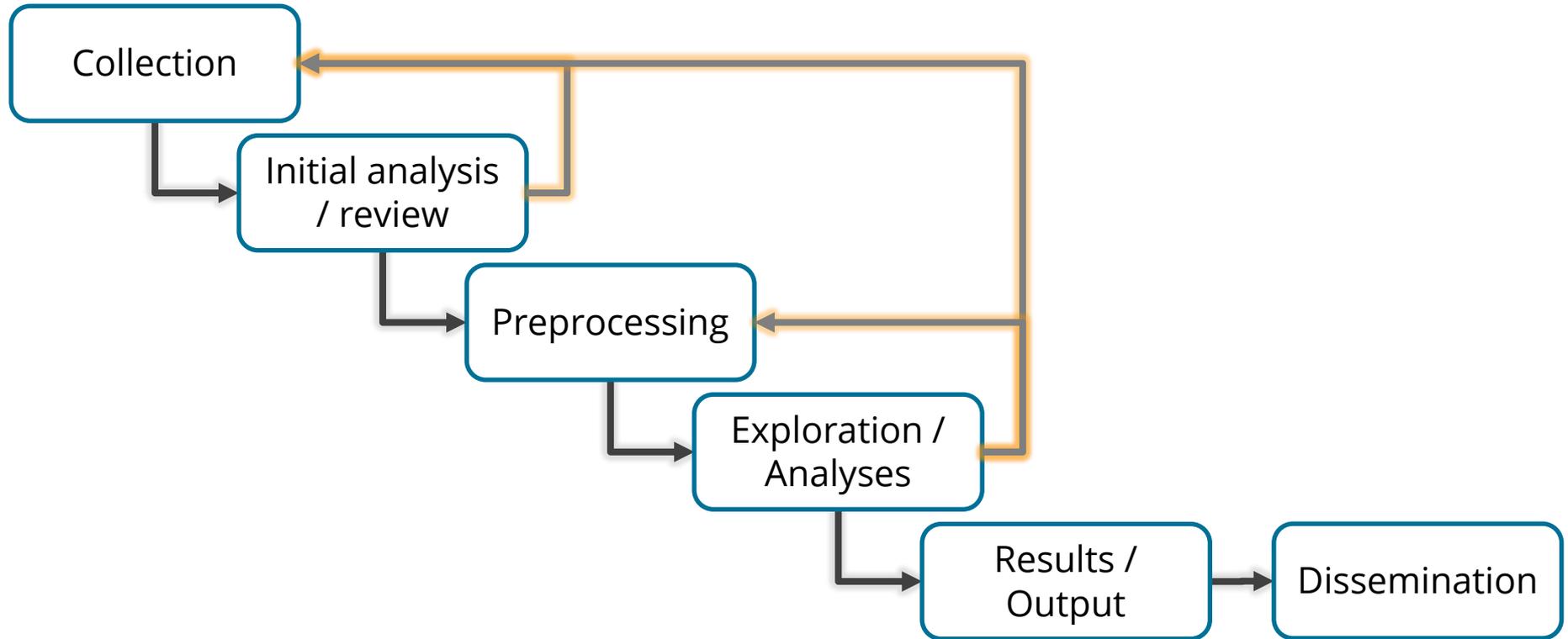
NLP Workflows

... are iterative



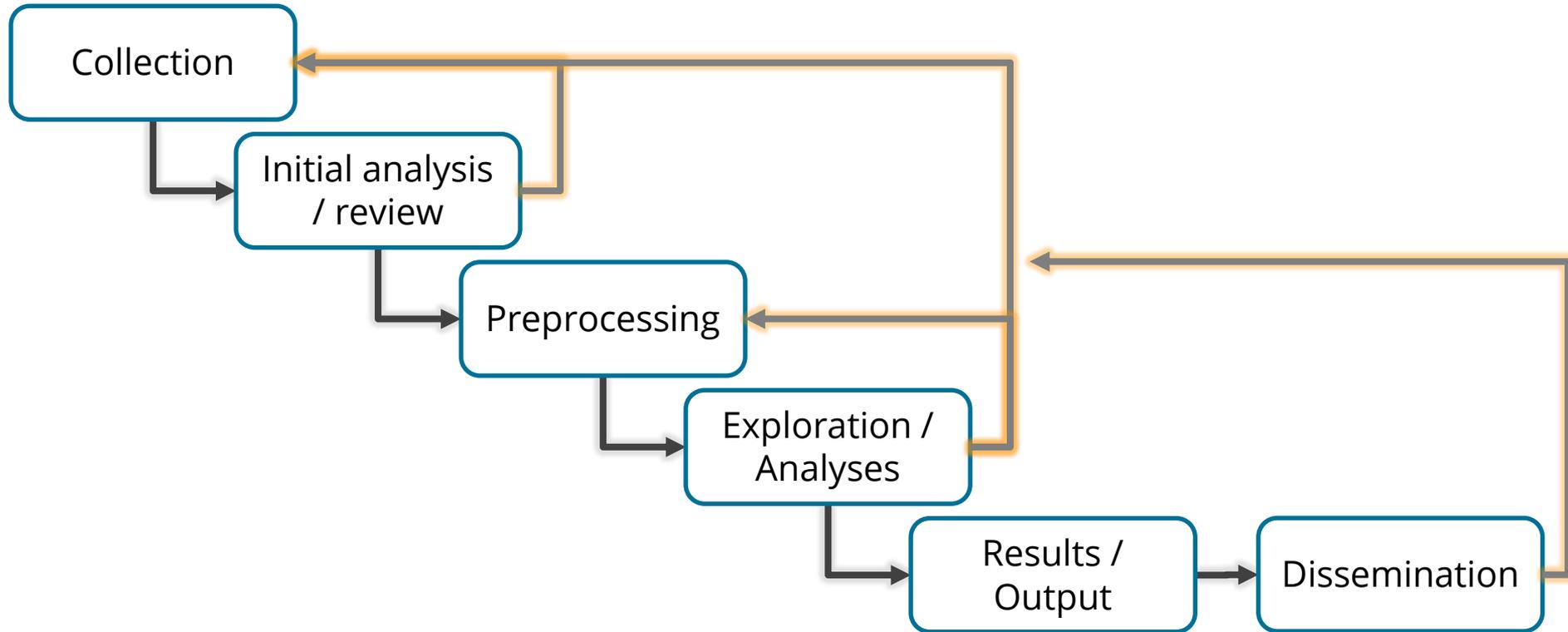
NLP Workflows

... are iterative

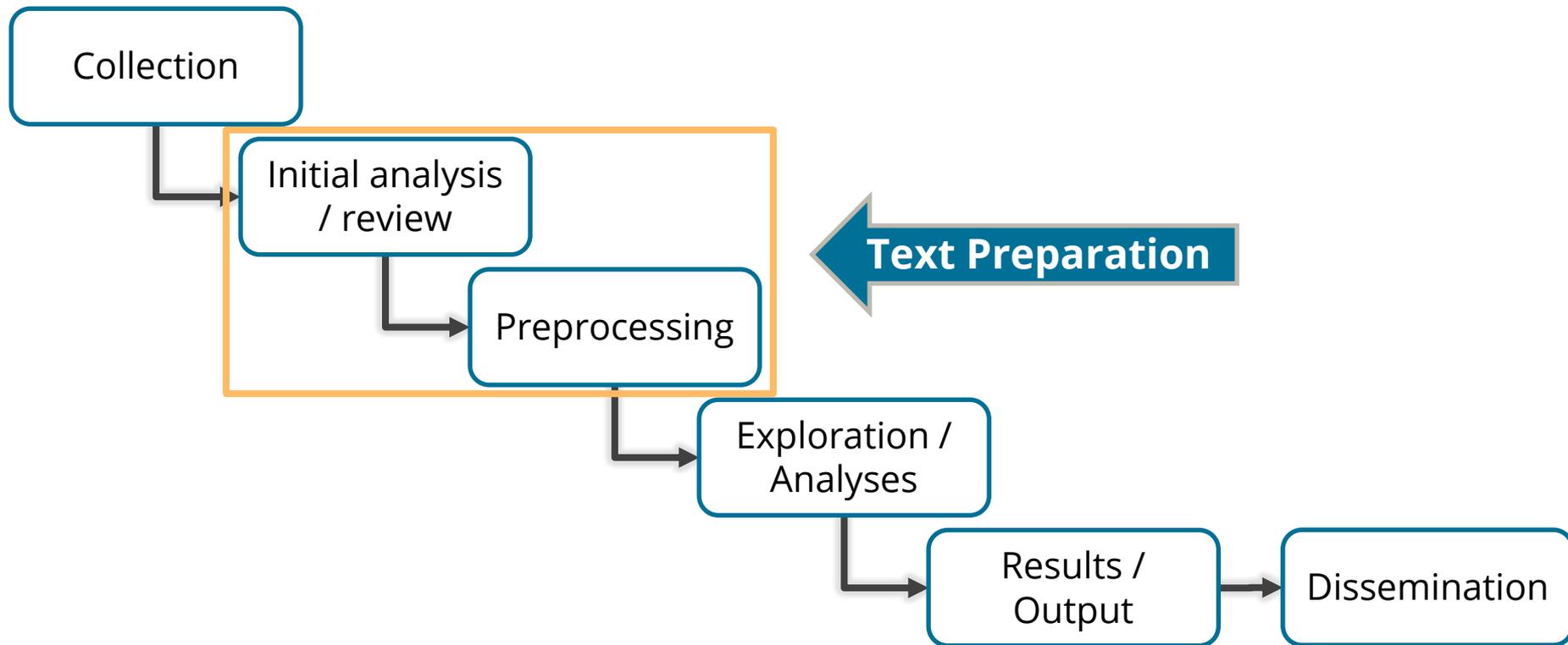


NLP Workflows

... are iterative



NLP Workflows



So, why 'prep'
your text?

Common OCR Issues



SOME COMMENTS ON CORREGGIO IN CONNECTION WITH HIS PICTURES IN DRESDEN.

A few years ago, it would have been hard to tell whether Correggio's *Night* or Raphael's *Madonna Di San Sisto* was the favourite picture of the Dresden Gallery.

The little sanctuary where the Virgin with Saint Sixtus floats above the

pseudo-altar was then crowded with worshippers as it is now, and Correggio's picture had quite as large and devout a following. But some change in popular taste has evidently taken place, for few people now linger before the *Night*.

What inference is to be drawn? Was the enthusiasm for Correggio merely a fashion which has had its season? He is certainly no longer admired as he was in the first few decades of this century, in the day when no gentleman could afford to be without his theory of the "Correggiosity of Correggio." The explanation is not far to seek.

The enthusiasm for Correggio dates from the time when, all the possible variations having been played upon the themes introduced by Raphael and Michelangelo, the Caracci betook themselves to a comparatively unlaboured field, and founded upon Correggio their school of painting, and thus succeeded in lending a new life to Italian art. Most people, however, appreciate only what is of their own day, and Correggio's interpreters proved far more interesting to their contemporaries than the master himself. The Caracci, Domenichino, Guercino, Guido Reni, and Lanfranco used up all the aesthetic capacity of their admirers, who believed in Correggio as the Catholic peasant doubtless believes in God, although he makes his offerings to the Saints. Furthermore, it was by no means easy to know the master himself. Correggio lived to be scarcely forty. Of his pictures then known the earliest dates from his twenty-first year, and in a career of barely twenty years no painter could have painted enough to fill the various collections of Europe. But in the third decade of this century, the few whose word was law in matters of taste suddenly turned away from Guido, Lanfranco, and their like, and gave themselves up to an unbridled enthusiasm for the Caracci and for their master, Correggio. Later, even the Caracci dropped out of sight, and Correggio stood alone.

The Madonna with St. Francis, No. 150.
The Madonna with St. Sebastian, No. 151.

The Nativity, called the "Night," No. 152.
The Madonna with St. George, No. 153.



i IPS^S^ffcls OME comments on correg ncRftJH GI0 IN CONNECTION with 1^58^^ HIS PICTURES IN DRESDEN. Spl^^T^ES A few years ago, it would have been hard u|^J^^fev^S to tell whether Correggio's *Night* or E^g^M Raphael's *Madonna Di San Sisto* was the ^L^^^M favourite picture of the Dresden Gallery. mmSSSmS^mSSm The little sanctuary where the Virgin with Saint *Sixtus* floats above the pseudo-altar was then crowded with worshippers as it is now, and Correggio's picture had quite as large and devout a fol lowing. But some change in popular taste has evidently taken place, for few people now linger before the *Night*. What inference is to be drawn? Was the enthusiasm for Correggio merely a fashion which has had its season? He is certainly no longer admired as he was in the first few decades of this century, in the day when no gentleman could afford to be without his theory of the "Correggiosity of Correggio." The explanation is not far to seek. The enthusiasm for Correggio dates from the time when, all the possible variations having been played upon the themes introduced by Raphael and Michelangelo, the Caracci betook themselves to a comparatively unlaboured field, and founded upon Correggio their school of painting, and thus succeeded in lending a new life to Italian art. Most people, however, appreciate only what is of their own day, ana Correggio's in terpreters proved far more interesting to their contemporaries than the master himself. The Caracci, Domenichino, Guercino, Guido Reni, and Lanfranco used up all the aesthetic capacity of their admirers, who believed in Correggio as the Catholic peasant doubtless believes in God, although he makes his offerings to the Saints. Furthermore, it was by no means easy to know the master himself. Correggio lived to be scarcely forty. Of his

Berenson, B. (1892). Some Comments on Correggio in Connection with His Pictures in Dresden. *The Knight Errant*, 1(3), 73-85. doi:10.2307/25515893

Common Transcription Issues

10

00:06:56.910 --> 00:07:07.200

aaa: We work and study on the traditional territory shared between the holden has shown a confederacy and the addition of a nations, which is acknowledged in the dish with one spoon off of.

11

00:07:08.220 --> 00:07:19.440

aaa: The wampanoag uses the symbolism of a dish to represent the territory and one spoon to represent that the people are to share the resources of this land and take only what they need.

We work and study on the traditional territory shared between the Haudenosaunee confederacy and the Anishinabe nations, which was acknowledged in the Dish with One Spoon Wampum belt. The wampum uses the symbolism of a dish to represent the territory, and one spoon to represent that the people are to share the resources of the land and only take what they need

Srsly, this stuff can [#lackconsistency](#)



Born-digital text (especially from SM) may be well-structured, but can also:

- contain a lot of spelling errors (sometimes intentionally) and non-words
- use non-traditional representations and abbreviations
- contain non-textual data like markup and embedded scripts

Some text prep considerations

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?

Hands-on text prep with OpenRefine

OpenRefine – for text preparation???

- Graphical interface (GUI)
- Non-destructive editing
- Self-documenting
- Reproducibility of steps



Using OpenRefine for text prep is best suited to....

- Scanned print documents:
 - good contrast
 - clearly defined boundaries
 - no or few tables, images or equations
- e.g. typed correspondence, minutes, manuscripts, reports, etc.

conversations and correspondence with municipal employees and regional and local MOE staff, a subsurface soils drilling program, a ground water monitoring program, a surface water monitoring program, and a landfill gas monitoring program. The work carried out under each of these parts of the work program is outlined in Sections 2.1 through 2.5.

2.1 DESK TOP INVENTORY

A series of air photos from 1953 to 1987 were used to identify the extent and process of filling over the period the landfill was operational. Historical water quality data for the Bay of Quinte and the Moira River provided by the MOE were reviewed. Recent study results on the Bay of Quinte were provided by the Bay of Quinte Remedial Action Plan. Geological and hydrogeological information was provided by the Ministry of Transportation through their work on the construction of Highway 62 on Zwick's Island. All of the background data reviewed was used in the design and interpretation of the drilling and monitoring programs.

2.2 SITE VISIT

A site visit was conducted on April 2, 1990 in order to relate the data collected during the desk top study to actual field conditions, and to collect additional data on the physical setting of the site. In addition to meeting with the MOE, GLL staff also met with a Municipal employee who worked at the landfill in the 1960's. This meeting provided first-hand information with respect to the nature of the refuse and the filling locations. The information obtained through the site visit assisted GLL in finalizing the drilling program as well as providing input to the health and safety protocols which would be followed during the course of field work.

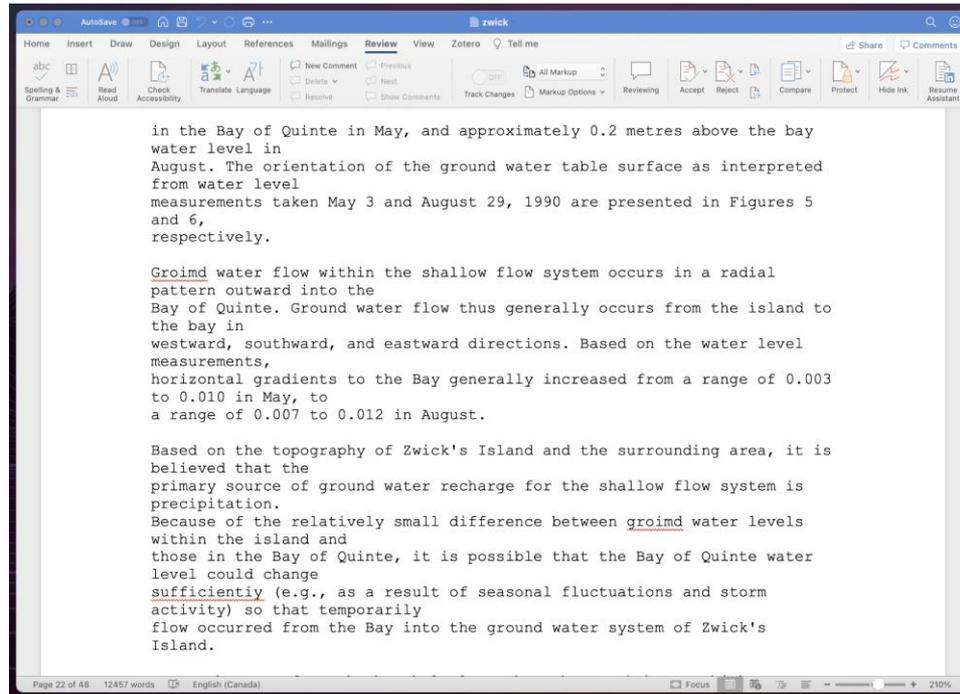
During the site visit surface water drainage and pathways were observed and noted. Conductivity measurements of surface water were taken at several locations around the site along with observations of iron staining and vegetation loss in roadside and drainage ditches in the northeast corner of the Island. Observations were also made as to the occurrence of ground settlement, locations of exposed refuse, and evidence of leachate seeps.

During the course of the field visit, GLL staff visited the City of Belleville Town Offices and obtained historical maps of Zwick's Island. This information assisted in establishing the original shoreline of the island and the landfilling locations.

The Dataset

- “Zwick's Island landfill environmental investigations” (1991)
 - Copied and pasted from full text on Internet Archive
- Transformations:
 - removed preamble
 - removed tabular data

Initial Data Analysis: in MS Word



A closer look at our errors...

- Pay attention to surrounding letters
(i.e. note context, not just the errors)
- Try to observe and record patterns

undenake

conceptmal

smdy

smdy

smdy

backgroimd

acmal

coUect

Mimicipal

fmalizing

Dtiring

aroidm

stammg

landfiUed

moimted

Figiu

Stratigraphie

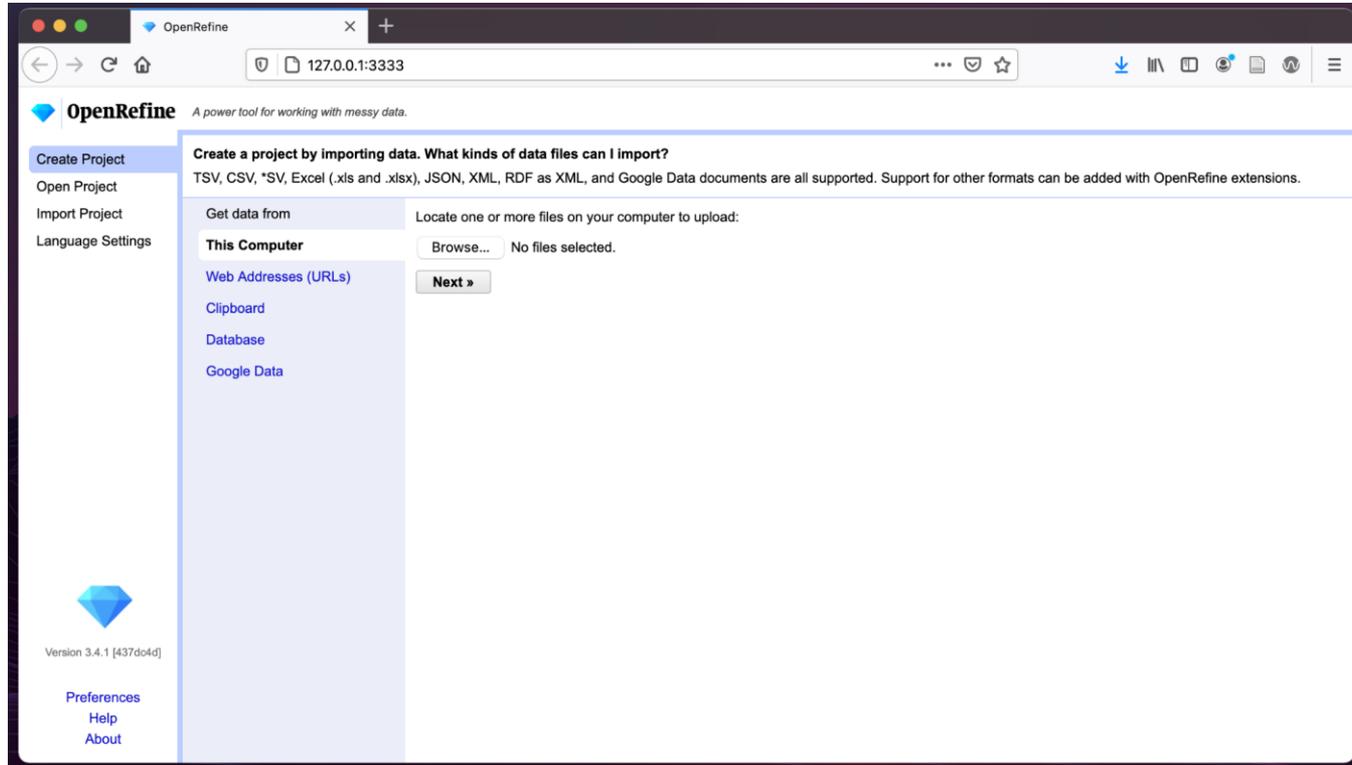
concurrentLly

Ruoride

Qjloride

Organo

Open OpenRefine



Initial Data Analysis: in OpenRefine

The screenshot displays the OpenRefine web interface in a browser window. The browser's address bar shows the URL `127.0.0.1:3333/project?project=1633045624437`. The page title is "OpenRefine zwick-demo". The main content area shows a list of 1151 rows of data, with the first 36 rows visible. The data is organized into columns, with the first column containing a list of numbered items. The interface includes a navigation bar at the top with "Open...", "Export", and "Help" buttons. A sidebar on the left provides instructions on using facets and filters. The data list is as follows:

All	Column 1
	1. ZWICK'S ISLAND LANDFILL
	2. ENVIRONMENTAL INVESTIGATIONS
	3. FINAL REPORT
	4. OCTOBER 1991
	5. TABLE OF CONTENTS
	6. Letter of Transmittal
	7. PAGE
	8. LQ INTRODUCTION 1
	9. LI BACKGROUND 1
	10. L2 OBJECTIVES 2
	11. 2.0 APPROACH 2
	12. 2.1 DESK TOP INVENTORY 3
	13. 2.2 SITE VISIT- 3
	14. 2.3 SUBSURFACE INVESTIGATIONS 4
	15. 2.4 SURFACE WATER MONITORING PROGRAM 7
	16. 2.5 LANDFILL GAS INVESTIGATIONS 8
	17. 3.0 PHYSICAL SETTING 9
	18. 3.1 GEOGRAPHIC SETTING 9
	19. 3.2 SITE HISTORY AND PRESENT USE 9
	20. 3.3 SITE HYDROLOGY 10
	21. 3.4 SUBSURFACE CONDITIONS 10
	22. 4.0 LANDFILL GAS 12
	23. 4.1 POTENTIAL FOR LANDFILL GAS 12
	24. 4.2 OCCURRENCE OF LANDFILL GAS 13
	25. 5.0 GROUND WATER AND LEACHATE 14
	26. 5.1 PHYSICAL HYDROGEOLOGY 14
	27. 5.2 GROUND WATER QUALITY 17
	28. 6.0 SURFACE WATER 21
	29. 6.1 SURFACE WATER QUALITY PARAMETERS OF INTEREST 21
	30. 6.2 WATER QUALITY CRITERIA 22
	31. 6.3 SURFACE WATER QUALITY 23
	32. 7.0 IMPACTS 27
	33. 7.1 IDENTIFICATION OF RECEPTORS, CONTAMINANT
	34. PATHWAYS AND CONTAMINANT LOADINGS 27
	35. 7.1.1 Receptors 27
	36. 7.1.2 Pathways 28

Prepare Dataset for Text Analysis

Tokenize, trim and remove blank rows

1151 rows

Show as: **rows** records Show: 5 10 25 50 rows

All	Column 1
1. Facet	FILL
2. Text filter	STIGATIONS
3. Edit cells	Transform...
4. Edit column	Common transforms
5. Transpose	Fill down
6. Sort...	Blank down
7. View	Split multi-valued cells...
8. Reconcile	Join multi-valued cells...
9. 2.0 APPROACH 2	Cluster and edit...
10. 2.1 DESK TOP INVENTO	Replace
11. 2.2 SITE VISIT - 3	
12. 2.3 SUBSURFACE INVESTIGATIONS 4	

13607 rows

Show as: **rows** records Show: 5 10 25 50 rows

All	Column 1
1. Facet	
2. Text filter	
3. Edit cells	Transform...
4. Edit column	Common transforms
5. Transpose	Fill down
6. Sort...	Blank down
7. View	Split multi-valued cells...
8. Reconcile	Join multi-valued cells...
9. OCTOBER	Cluster and edit...
10. 1991	Replace
11. TABLE	
12. OF	
13. CONTENTS	
14. Letter	
15. of	
16. Transmittal	

Trim leading and trailing whitespace

Collapse consecutive whitespace

Unescape HTML entities

Replace Smart quotes with ascii

To titlecase

To uppercase

To lowercase

To number

To date

To text

To null

To empty string

Filter and Facet to Find Errors

A screenshot of a software interface showing a context menu for 'Column 1'. The menu is open, displaying various options for filtering and faceting data. The 'View' option is selected, and a sub-menu is visible showing 'Customized facets' and 'Word facet'. The 'Word facet' option is highlighted. Below the menu, a list of words is visible, including 'monitoring', 'contamination', 'remedial', 'mainly', 'municipal', 'former', 'some', 'commercial', 'system', 'managed', 'comer,', and 'Ramada'.

- Column 1
 - Facet
 - Text facet
 - Numeric facet
 - Timeline facet
 - Scatterplot facet
 - Custom text facet...
 - Custom Numeric Facet...
 - Customized facets
 - Word facet
 - Duplicates facet
 - Numeric log facet
 - 1-bounded numeric log facet
 - Text length facet
 - Log of text length facet
 - Unicode char-code facet
 - Facet by error
 - Facet by null
 - Facet by empty string
 - Facet by blank (null or empty string)
 - Text filter
 - Edit cells
 - Edit column
 - Transpose
 - Sort...
 - View
 - Reconcile

monitoring
contamination
remedial
mainly
municipal
former
some
commercial
system
managed
comer,
Ramada

A screenshot of a software interface showing a 'Facet / Filter' panel. The panel is titled 'Facet / Filter' and includes 'Undo / Redo 4 / 4' and buttons for 'Refresh', 'Reset All', and 'Remove All'. Below the buttons, there are two facet panels for 'Column 1'. The first panel shows a search input with the text 'm', a checked 'case sensitive' checkbox, and an unchecked 'regular expression' checkbox. The second panel shows a list of 384 choices, sorted by 'name count'. The list includes words like 'decomposition', 'define', 'demonstrate', 'demonstrates', 'determination', 'determine', 'determined', 'determining', 'dumping', 'dumppmg', and 'emissions', each with a count next to it.

Facet / Filter Undo / Redo 4 / 4

Refresh Reset All Remove All

Column 1 invert reset

m

case sensitive regular expression

Column 1 change

384 choices Sort by: name count

decomposition 3
define 1
demonstrate 1
demonstrates 1
determination 1
determine 13
determined 1
determining 1
dumping 1
dumppmg 1
emissions 1

Find and Replace with Text Filter

The screenshot shows a data table with a 'Facet / Filter' panel on the left. The filter is set to 'Column 1' with the text 'smdy'. The table header indicates '8 matching rows (12457 total)'. A context menu is open over the table, listing various actions. The 'Replace' option at the bottom of the menu is highlighted.

All	Column 1
1082.	
1089.	
1209.	
4539.	
6835.	
10145.	
10849.	
12155.	

- Facet
- Text filter
- Edit cells
- Edit column
- Transpose
- Sort...
- View
- Reconcile
- Replace

The 'Replace' dialog box is shown with the following settings:

- Find:** smdy
- case insensitive
- whole word
- regular expression
- Replace with:** study
- use \n for new lines, \t for tabulation, \\n for \n, \\t for \t.

Buttons: OK, Cancel

Find and Replace with Text Filter

Filter by “no” as case sensitive, then word facet & include

Facet / Filter Undo / Redo 5 / 6

Refresh Reset All Remove All

Column 1 invert reset

no

case sensitive regular expression

Column 1 change

34 choices Sort by: name count

not 27

noted 5

noted. 1

noticeably 1

Ontano 1 include

Organo 1

organo 1

organochlorine 2

Organochlorine 2

Phenol 2

phenol 3

Replace

Find: no

case insensitive whole word regular expression

Leave blank to add the replacement string after each character.
Check "regular expression" to find special characters (new lines, tabulations...) or complex patterns.

Replace with: rio

use \n for new lines, \t for tabulation, \\n for \n, \\t for \t.

If "regular expression" option is checked and finding pattern contains groups delimited with parentheses, \$0 will return the complete string matching the pattern, and \$1, \$2... the 1st, 2nd... group.

OK Cancel

Find and Replace with GREL

Filter by "oim" (oun) → `value.replace('im', 'un')`

A screenshot of a data table interface. The table has a column header 'Column 1' and several rows of data. A context menu is open over the first row, showing options like 'Facet', 'Text filter', 'Edit cells', 'Edit column', 'Transpose', 'Sort...', 'View', and 'Reconcile'. The 'Edit cells' option is expanded, showing a sub-menu with 'Transform...', 'Common transforms', 'Fill down', 'Blank down', 'Split multi-valued cells...', 'Join multi-valued cells...', 'Cluster and edit...', and 'Replace'. The 'Replace' option is highlighted.

Column 1
backgroundd
groimd
groimd
aroimd

A screenshot of the 'Custom text transform on column Column 1' dialog box. The 'Expression' field contains the GREL code `value.replace('im', 'un')`. The 'Language' is set to 'General Refine Expression Language (GREL)'. A 'Preview' tab is active, showing a table with the transformed data. The 'On error' section has 'keep original' selected. The 'Re-transform up to 10 times until no change' checkbox is unchecked. 'OK' and 'Cancel' buttons are at the bottom.

Custom text transform on column Column 1

Expression: `value.replace('im', 'un')` Language: General Refine Expression Language (GREL) No syntax error.

Preview History Starred Help

row	value	value.replace('im','un')
1252.	backgroimd	background
1406.	aroimd	around
1679.	truck-moimted	truck-mounted
3853.	encoimtered	encountered
4827.	Groimd	Ground
4931.	groimd	ground

On error: keep original set to blank store error Re-transform up to 10 times until no change

OK Cancel

Find and Replace with GREL

Try it out: filter by “`tiy`” (`tly`) → `value.replace('tiy', 'tly')`

Custom text transform on column Column 1

Expression `value.replace('tiy','tly')` Language General Refine Expression Language (GREL) ⓘ No syntax error.

Preview History Starred Help

row	value	value.replace('tiy','tly')
4749.	consistently	consistently
4956.	sufficiently	sufficiently
7589.	significantly	significantly
8454.	slightly	slightly
8516.	slightly	slightly
9546.	significantly	significantly
...

On error keep original set to blank store error Re-transform up to 10 times until no change

OK Cancel

Find and Replace with Regular Expression (Regex)

Filter by "mg" → try mg\$ with "regular expression" checked ("ing" as "mg" error)

The screenshot shows a data visualization interface with a facet filter. The filter is named 'Column 1' and contains the regular expression 'mg\$'. The 'regular expression' checkbox is checked, while 'case sensitive' is unchecked. The filter results in 5 matching rows out of a total of 12457. The table below shows the first five rows of the filtered data.

Facet / Filter		Undo / Redo 9 / 9	
Refresh		Reset All Remove All	
Column 1		invert reset	
mg\$			
<input type="checkbox"/>	case sensitive	<input checked="" type="checkbox"/>	regular expression
5 matching rows (12457 total)			
Show as: rows records		Show: 5 1	
All	Column 1		
☆	🗨	1414.	stammg
☆	🗨	2988.	dumpmg
☆	🗨	4151.	ventmg
☆	🗨	4386.	bemg
☆	🗨	12248.	reconstructmg

Quick Guide to Regex

^ start of expression

\$ end of expression

E.g. **^T\$** will only return cells with “T”

^mn will only return cells that **start** with “mn”

ent\$ will only return cells that **end** with “ent”

Quick Guide to Regex

[string] - contains any of the letters

[^string] - does not contain the letters

E.g. **[iou]m** will return words that contain "im," "om" and "um"

ti[^o] will exclude "tion"

Find and Replace with Regex ctd.

Try it out:

^mt → `value.replace('m' , 'in')`

[a-z]U → `value.replace('U' , 'll')` [with case sensitive checked]

Others...?

Reconstitute your Document

Custom Tabular Exporter

Content **Download** Upload Option Code

Line-based text formats

- Tab-separated values (TSV)
- Comma-separated values (CSV)
- Custom separator _____

Line separator

Character encoding

Always quote text

Other formats

- Excel (.xls)
- Excel in XML (.xlsx)
- HTML table

Preview **Download**

Cancel

Export your "Recipe" of Tasks

Extract Operation History

Extract and save parts of your operation history as JSON that you can apply to this or other projects in the future.

- Split multi-valued cells in column Column 1
- Text transform on cells in column Column 1 using expression value.trim()
Star row 4
- Remove rows
- Text transform on cells in column Column 1 using expression value.replace("smdy","study")
- Text transform on cells in column Column 1 using expression value.replace("no","rio")
- Text transform on cells in column Column 1 using expression value.replace("no","nic")
- Text transform on cells in column Column 1 using expression grel.value.replace('im','un')
- Text transform on cells in column Column 1 using expression grel.value.replace('iy','ly')
- Text transform on cells in column Column 1 using expression grel.value.replace('mg','ing')
- Mass edit cells in column Column 1
Star row 11473

Select All Unselect All

Close

```

{
  "facets": [
    {
      "type": "text",
      "name": "Column 1",
      "columnName": "Column 1",
      "query": "staining",
      "mode": "regex",
      "caseSensitive": false,
      "invert": false
    }
  ],
  "mode": "row-based"
},
{
  "columnName": "Column 1",
  "expression": "value",
  "edits": [
    {
      "from": [
        "staining"
      ],
      "fromBlank": false,
      "fromError": false,
      "to": "staining"
    }
  ]
},
{
  "description": "Mass edit cells in column C
}

```

Break time!

Let's take 10

Programmatic approaches with Python

—

Text prep and analysis as a continuum of mediation

Completely
manual



Completely
automated

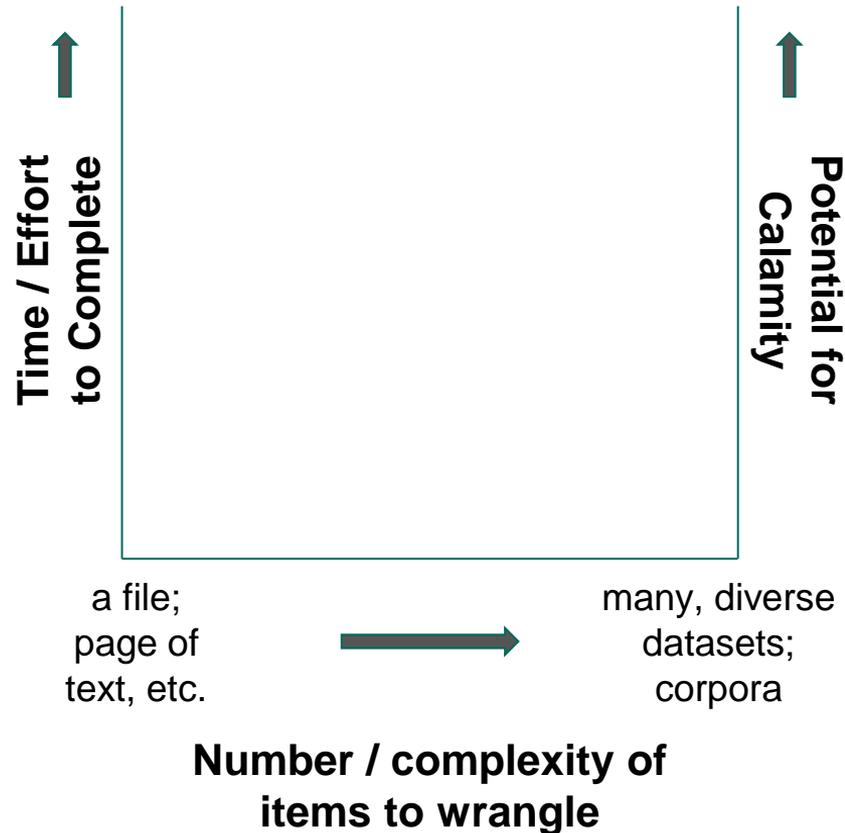
Text prep and analysis as a continuum of mediation

Completely
manual

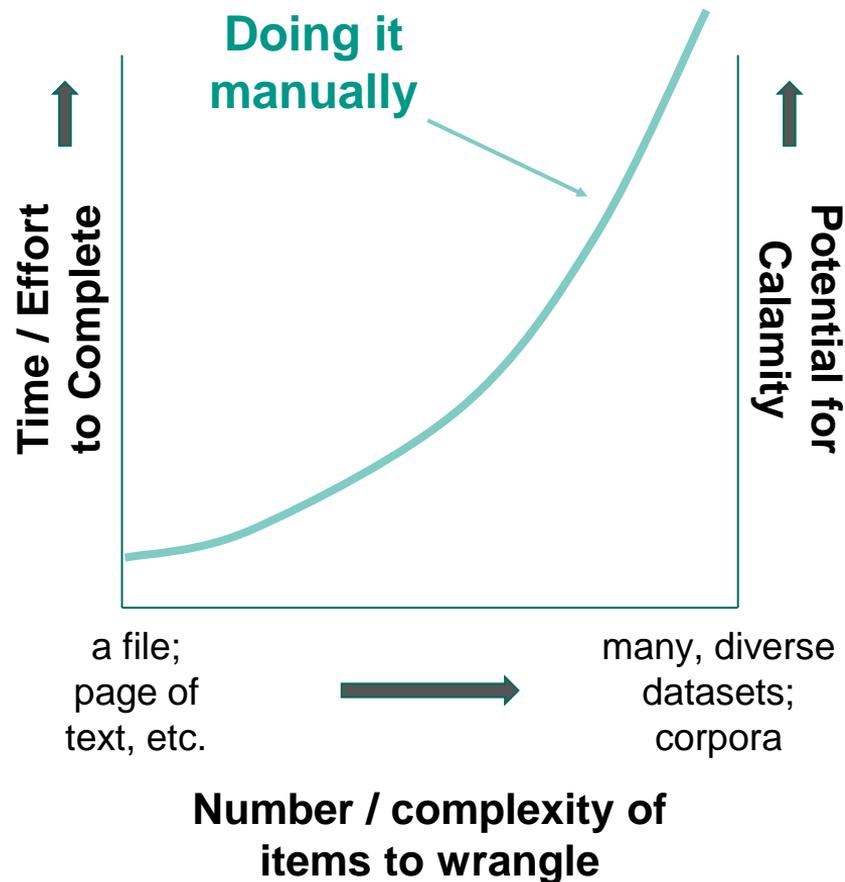


Completely
automated

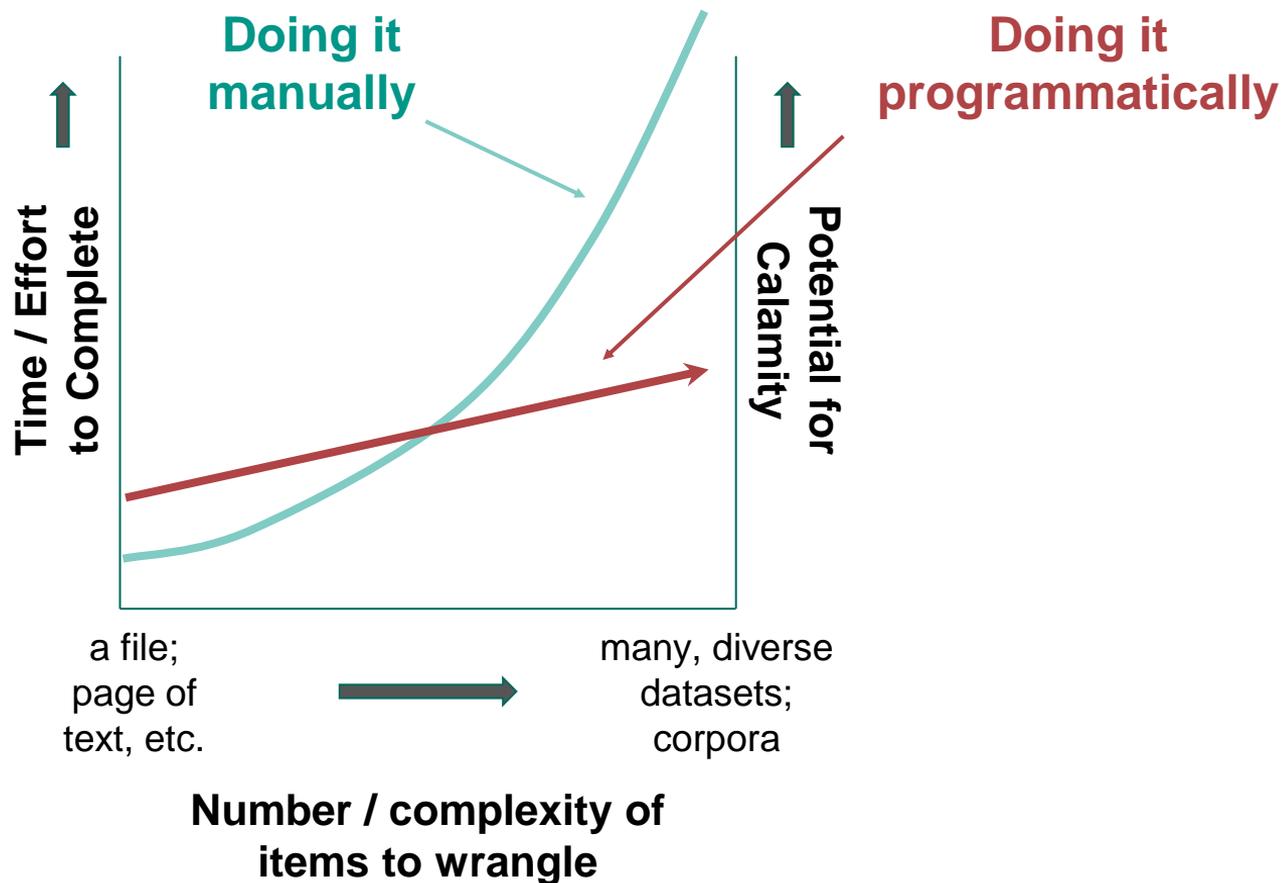
So, when to let the computer take over?



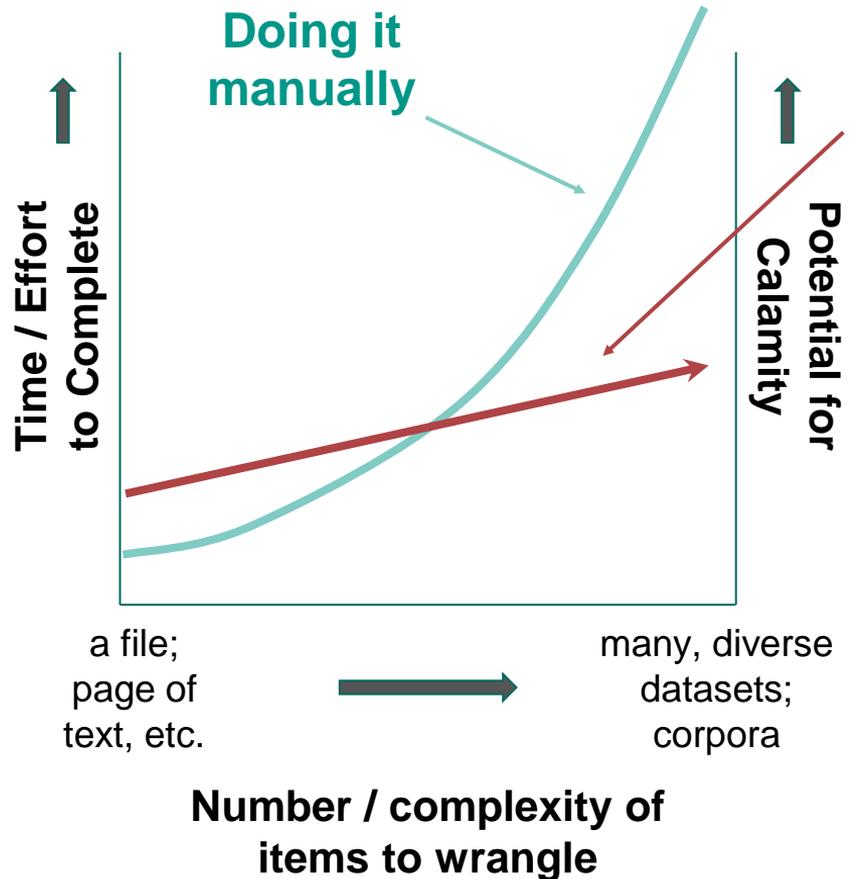
Spreadsheets: The frenemy of research



Spreadsheets: The frenemy of research



Spreadsheets: The frenemy of research



Doing it programmatically

HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE? (ACROSS FIVE YEARS)

	HOW OFTEN YOU DO THE TASK					
	50/DAY	5/DAY	DAILY	WEEKLY	MONTHLY	YEARLY
1 SECOND	1 DAY	2 HOURS	30 MINUTES	4 MINUTES	1 MINUTE	5 SECONDS
5 SECONDS	5 DAYS	12 HOURS	2 HOURS	21 MINUTES	5 MINUTES	25 SECONDS
30 SECONDS	4 WEEKS	3 DAYS	12 HOURS	2 HOURS	30 MINUTES	2 MINUTES
1 MINUTE	8 WEEKS	6 DAYS	1 DAY	4 HOURS	1 HOUR	5 MINUTES
5 MINUTES	9 MONTHS	4 WEEKS	6 DAYS	21 HOURS	5 HOURS	25 MINUTES
30 MINUTES		6 MONTHS	5 WEEKS	5 DAYS	1 DAY	2 HOURS
1 HOUR		10 MONTHS	2 MONTHS	10 DAYS	2 DAYS	5 HOURS
6 HOURS				2 MONTHS	2 WEEKS	1 DAY
1 DAY					8 WEEKS	5 DAYS

HOW MUCH TIME YOU SHAVE OFF

Relevant xkcd:

<https://xkcd.com/1205/>

Reasons to code some/all of your approach

- To save you time
- To scale your approaches
- To reduce analytical toil
- Because (some people think) it is fun
- To build your own 'toolkit' of analytical scripts, functions, modules
- To enhance tractability, transparency, reproducibility, and reuse

To our Jupyter Notebook

Go to u.mcmaster.ca/dmnds-text-prep and save a copy to your Google Drive.

Follow along with Jay's instructions

A hands-on sampler

—

How Named Entity Recognition (NER) Works



— Text to annotate —

Three years after the passage of the Fugitive Slave Act of 1850, A. D. Shadd moved his family to the United Canadas (Canada West), settling in North Buxton, Ontario. In 1858, he became one of the first black men to be elected to political office in Canada, when he was elected to the position of Counsellor of Raleigh Township, Ontario.

— Annotations —

named entities X

— Language —
English

Submit

Named Entity Recognition:

1	Mary Ann Shadd was born in Wilmington, Delaware, on October 9, 1823, the eldest of 13 children to Abraham Doras Shadd (1801 - 1882) and Harriet Burton Parnell, who were free African - Americans.	PERSON	CITY	STATE_OR_PROVINCE	DATE	NUMBER	PERSON	DATE	DATE
2	Abraham D. Shadd was a grandson of Hans Schad, alias John Shadd, a native of Hesse - Cassel who had entered the United States serving as a Hessian soldier with the British Army during the French and Indian War.	PERSON	PERSON	PERSON	COUNTRY	MISC	TITLE		
3	Hans Schad was wounded and left in the care of two African - American women, mother and daughter, both named Elizabeth Jackson.	PERSON	NUMBER	NATIONALITY	PERSON				
4	The Hessian soldier and the daughter were married in January 1756 and their first son was born six months later.	MISC	TITLE	DATE	ORDINAL	TIME			
5	[5] A. D. Shadd was a son of Jeremiah Shadd, John's younger son, who was a Wilmington butcher.	NUMBER	PERSON	PERSON	PERSON	CITY	TITLE		
6	Abraham Shadd was trained as a shoemaker [6] and had a shop in Wilmington and later in the nearby town of West Chester, Pennsylvania.	PERSON	TITLE	NUMBER	CITY	CITY	STATE_OR_PROVINCE		
7	In both places he was active as a conductor on the Underground Railroad and in other civil rights activities, being an active member of the American Anti-Slavery Society, and, in	TITLE	LOCATION	ORGANIZATION					

Try it out in Jupyter Notebooks...

+ 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_))
```

Discerning Corpus "Topics" with Topic Modeling

Run 50 iterations Iterations: 150 Train with 25 topics

Topic Documents Topic Correlations Time Series Vocabulary Downloads

[0] pop culture fiction representations industrial film artificial simultaneously far natural

[1] neuromantic cowboy sexism concept terms same does feminist suggests prosthesis

[2] identity realm representation cyberspace one's possibilities others over opposes perceive

[3] human itself began still felt gestures limited garment became matrix

[4] lather rinse repeat specialized highly brain task becomes lateral ability

[5] body physical virtual form despite bodies information role instance does

Documents are sorted by their proportion of the currently selected topic, biased to prefer longer documents.

[2/8.5%] Much like a science fiction film, my work is situated in a hybrid space-time, simultaneously part of the present and the future. The figure of the cyborg and the realm of cyberspace are central to the work—both of which are similarly here and yet, not-here. That is, their mundane existence in dai...

[1/8.0%] I intend to address issues surrounding technological interactions with the body, and examine how these interactions are amplified in film and literature. Binary relationships are imposed and exaggerated in these often-oversimplified representations of technology: male/female, mind/body, transcend...

[36/7.2%] Grenville is nothing if not thorough; pop culture, industrial antiques and artworks alike constitute the collection. The juxtaposition of (arti)fact with fiction outlines the parallel developments in the three realms, and alludes to the nebulous boundaries between them. The overall impression of ...

Use a different collection:
Documents d-corp.txt
Stoplist No file selected.

Try it out in Jupyter Notebooks...

+ 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())
```

Sentiment Analysis – Jay

Go to u.mcmaster.ca/dmds-sentimental and save a copy to your Google Drive.

Follow along with Jay's instructions

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**