

1 FINANCE 751 Technical Note

This technical note informs software installation, transcript extraction, implementation and comparison methodologies to ascertain measures of corporate culture in NZX50 companies, and six Australian commercial banks, using 258 earnings call transcripts. This note describes the steps taken to implement Option-2, replicating the corporate culture results. Additionally, this technical note is for a MacOS operating system and assumes basic proficiency in package management and Python programming.

1.1 Installation

This section informs the installation of required software to facilitate analysis.

1. Install several software packages to run the StanfordCoreNLP to measure corporate culture from text files and develop transcript processing code. Anaconda is a distribution of the Python programming language, simplifying package management to develop code in the Python language. Microsoft Visual Studio Code is an integrated development environment, suitable for application building. The combination of both Anaconda and Microsoft Visual Studio Code enable programming environments to process the transcripts.
2. Secondly, clone the remote repository implementing the method described by Li et al., (2021) to your local directory. Read the instructions carefully for correct installation. In particular, change the `os.environ["CORENLP_HOME"]` variable in `global_options.py` file to the installation location of the `stanford-corenlp-full-2018-10-05` directory on your local device. Install the required python packages excluded from Anaconda using the pip package and requirements.txt with the `pip install requirement.txt` command in the terminal. Some packages require specific versions. If you need to revert to a previous version, use the terminal command `pip install PackageName==Version` to revert to a previous version.
3. Test the correct installation of the StanfordCoreNLP using the document text files from remote repository by following the ReadMe instructions. Progress to transcript extraction after the successful execution of the StanfordCoreNLP. Otherwise, review the above installation process before progressing.

1.2 Transcript Extraction

This section informs extracting earnings call transcripts from Capital IQ.

1. Review the `firm_id` column in the `1.firm_score.xlsx` sheet from Option-1 to identify the companies related to the 258 earnings call transcripts.
2. Navigate to Capital IQ, selecting the companies tab, followed by the transcripts link.
3. In the search criteria company search bar, type in and select each of the unique companies listed in the `firm_id` column mentioned above. The selected entities will update to list sixteen companies as Telecom Corp of New Zealand Ltd changed rebranded to Spark new Zealand Limited.
4. Change the time frame from 01/01/2009 to 01/10/2021 to ensure you include all 258 transcripts listed in the `1.firm_score.xlsx` spreadsheet and select search in middle-right of the webpage.
5. Select all transcripts on the page by ticking the top tick box middle-left of the webpage.
6. Click the options dropdown middle-left of the webpage, and select Download in .Zip file to download all selected documents into a .Zip file.
7. Scroll to the bottom of the page to select the next subset of transcripts.
8. Repeat steps five through seven to download all transcripts in .Zip files.
9. Create a new local directory titled 'transcripts', unzip all .Zip files, moving all transcripts to this newly created directory.
10. Review the transcripts in the 'transcripts' directory. The filenames align with the filename column in the `1.firm_score.xlsx` spreadsheet. There will be multiple transcripts with the same name e.g., Air New Zealand Limited - ShareholderAnalyst Call.pdf. Consult filename and calltime columns in `1.firm_score.xlsx` to identify the correct transcripts according to date e.g., 201510 is October 2015, deleting the incorrect duplicates. After, the subset of 258 transcripts will exist amongst the full set in the transcript directory.

1.3 Implementation

This section highlights the code to process earnings call transcripts, execute the StanfordCoreNLP and compare the results. The implementation was partitioned into three Python functions within the finance-751-cmcd398.py script (1.5.3). This section provides a high level overview of the code with further details described in the code comments. Transcripts have a common structure. The first three pages are front-matter. The last page is the legal disclaimer. Some transcripts don't have Q&A sections while others have multiple. The transcripts without Q&A sections are isolated and excluded during processing. Transcripts with multiple Q&A sections are manually condensed prior to processing but excluded during comparison.

1.3.1 Variables

The definition of several variables and arrays take place prior to implementation.

1. Set strings describing the relative paths for the 1.firm_score.xlsx file, transcript directory, selected transcript directory to move 258 transcripts of interest, transcript directory for processed transcripts after removing Q&A sections, documents.txt file, documents_ids.txt, and processed text directory.
2. Review each transcript in the 1.firm_score.xlsx filename column to record the page number for the first page of the Q&A section, appending each value to the end of an array. If no Q&A section exists, record a value of 4. **The preservation of order is imperative with the position of the page number matching the position of the filename in the filename list from 1.firm_score.xlsx.**
3. Set an array listing the set of company ids from the 1.firm_score.xlsx spreadsheet aligning with an array listing the cumulative position of the final transcript corresponding to the company id. For example, Air New Zealand (ANZ) has 11 transcripts. Auckland International Airport (AIA) has 13 transcripts. Therefore, ANZ and AIA have values of 11 and 24, respectively, in the cumulative position array.
4. Set strings describing the relative paths for output files, results spreadsheet, and firm scores outputs from the StanfordCoreNLP.
5. Set binary variables (TRUE or FALSE) to control the execution of the below functions.

1.3.2 Prepare_documents.py

This function isolates the Q&A sections of each transcript, converts each transcript to a line in a text file, and returns the document text file and identification. The following sequence of functions are nested within, called on in the order below.

1. **get_transcripts** extracts a list of filenames from the 1.firm_score.xlsx spreadsheet, transferring transcripts of interest to the transcripts selected directory.
2. **remove_transcript_metadata** deploys the pdfrw package to extract each page of the Q&A section per transcript, using the array denoting the starting page number for the Q&A section, creating a processed transcript stored in the transcripts processed directory.
3. **create_ids** creates various forms of identification in data frames for comparison while excluding transcripts without Q&A sections.
4. **create_documents_text** deploys the pdfminer package to convert each processed transcript into a single line of text, appending each line to the document.txt file to use as an input for the StanfordCoreNLP.

1.3.3 Perform_stanford_nlp.py

This function executes each one of the five Python functions integral to StanfordCoreNLP in the following order. The provision of two separate dictionaries (NZD/AUS and US) informs analysis.

1. **parse.py** to parse the raw documents.
2. **clean_and_train.py** to clean, remove stopwords, and named entities in the parsed documents text file.
3. **create_dict.py** to create the expanded dictionary.
4. **score.py** to score the document. This implementation uses the TF-IDF weights used in the article.
5. **aggregate_firms.py** to aggregate the scores to the firm-time level.

Complete steps one, two and three. Next, replace the `expanded_dictionary.csv` in the `dict` directory with the AUS/NZD trained dictionary. It is possible to manually edit these dictionaries in attempts to improve scores. However, the provided dictionaries trained to ascertain the original scores. Therefore, the provided dictionaries were left unchanged in replicating scores. Next, Run `score.py` and `aggregate_firms.py`, saving the `scores_TFIDF.csv` as an `xlsx` file to the `comparisons` directory. Repeat steps four and five with the US dictionary.

1.3.4 Compare_results.py

This function combines a formatted `1.firm_scores.xlsx` document with the TF-IDF output scores from `perform_stanford_nlp.py` by merging data frames on document identification in order to make comparisons. `Compare_results.py` must be repeated for both dictionaries. After, combine both comparison spreadsheets to compare results from both sets of dictionaries, deleting duplicate values.

1.4 Comparison

This section compares our replication of the measures for corporate culture across the five values (Innovation, Integrity, Quality, Respect, Teamwork) using NZ/AUS and US dictionaries. We acknowledge the provided scores have slightly shorter document lengths, likely from different pdf to text conversion methodologies. Our analysis detected a few abnormalities in the aforementioned subset of transcripts, omitting the majority of Q&A sections (1.5.1), in addition to a subset of transcripts not including Q&A sections but trained on presentation sections. The author's emphasize the presentation sections in transcripts are likely not a true reflection of company culture as edited by corporate lawyers and PR personal. Subsequently, we exclude these transactions during processing.

1.4.1 Accuracy Measures

Absolute and percentage differences between our replication and the provided results are displayed in the `751-comparison.xlsx` workbook. However, we utilize the following equations to measure the accuracy of our replication across companies, values, and total results.

$$\text{Individual} = 1 - \frac{\sum_i |\text{New}_{i,j,k} - \text{Old}_{i,j,k}|}{\sum_i \text{Old}_{i,j,k}} \forall j, k \quad (1) \quad \text{Total} = 1 - \frac{\sum_i \sum_j \sum_k |\text{New}_{i,j,k} - \text{Old}_{i,j,k}|}{\sum_i \sum_j \sum_k \text{Old}_{i,j,k}} \quad (2)$$

$$\text{Company} = 1 - \frac{\sum_i \sum_k |\text{New}_{i,j,k} - \text{Old}_{i,j,k}|}{\sum_i \sum_k \text{Old}_{i,j,k}} \forall j \quad (3) \quad \text{Value} = 1 - \frac{\sum_i \sum_j |\text{New}_{i,j,k} - \text{Old}_{i,j,k}|}{\sum_i \sum_j \text{Old}_{i,j,k}} \forall k \quad (4)$$

$$i \in \{1, \dots, N\} \quad (5)$$

$$j \in \{\text{Air New Zealand}, \dots, \text{Westpac Banking Corporation}\} \quad (6)$$

$$k \in \{\text{Innovation}, \text{Integrity}, \text{Quality}, \text{Respect}, \text{Teamwork}\} \quad (7)$$

1.4.2 Results

Individual, **Company**, **Value**, and **Total** measure the accuracy of our replication for a specific company and value, company across all values, value across all companies, and across all values and companies respectively. The accuracy results are displayed in a matrix (1.5.2). There are a few abnormalities. The value Teamwork for Goodman Property Trust is NA as both values in the original results are zero. Our replication for Teamwork using the NZD/AUS dictionary, and Respect using the US dictionary, deviate relatively from provided figures in our replication. The later driven by material differences in Infratil's replication (-89%) and 80% accuracy for Westpac Banking Corporation. The remaining results from the Respect value using the US dictionary are above 80%. However, all Teamwork results using the NZD/AUS dictionary are above 83%, not raising cause for concern. The **Company** level of accuracy is above 90% for all Company IDs. Each **Value** level of accuracy is above 90% except for the Respect value measured by the US dictionary (80%). Finally, the **Total** level of accuracy is 93%. Discrepancies may be caused by small differences in documents lengths, or abnormalities when parsing documents using StanfordCoreNLP. In summary, our results are highly accurate and satisfactory across Company IDs and Values, providing supporting evidence our replication is successful.

References

Li, K., Mai, F., Shen, R., & Yan, X. (2021). Measuring corporate culture using machine learning. *The Review of Financial Studies*, 34(7), 3265–3315.

1.5 Appendix

1.5.1 Transcripts with Multiple Q&A Sections

The following transcripts have multiple Q&A sections. The sections are consolidated into one section by deleting the presentation material in between the Q&A sections. However, they are excluded from comparison calculation as Helen only uses the last Q&A section. We took the perspective the last section alone does not proxy for the entire Q&A sections in the transcript. Therefore, not suitable for measuring corporate culture given the document lengths.

- Australia and New Zealand Banking Group Limited - ShareholderAnalyst Call.pdf,
- Bank of Queensland Ltd. - ShareholderAnalyst Call.pdf
- Commonwealth Bank of Australia - ShareholderAnalyst Call.pdf
- Infratil Limited - AnalystInvestor Day.pdf
- Infratil Ltd. - AnalystInvestor Day.pdf
- National Australia Bank Limited - ShareholderAnalyst Call.pdf

1.5.2 Result Matrix

Firm	Document Length	Innovation (ANZ)	Integrity (ANZ)	Quality (ANZ)	Respect (ANZ)	Teamwork (ANZ)	Innovation (US)	Integrity (US)	Quality (US)	Respect (US)	Teamwork (US)	Company
Air New Zealand Limited	96%	95%	91%	97%	96%	92%	93%	94%	95%	90%	93%	97%
Auckland International Airport Limited	96%	95%	94%	95%	96%	87%	94%	91%	93%	92%	93%	98%
Australia New Zealand Banking Group Limited	96%	93%	94%	93%	94%	93%	94%	90%	94%	84%	95%	96%
Bank of Queensland Limited	96%	94%	93%	95%	93%	83%	95%	91%	93%	93%	89%	97%
Bendigo and Adelaide Bank Limited	96%	93%	95%	94%	94%	94%	94%	94%	94%	89%	95%	97%
Commonwealth Bank of Australia	96%	92%	92%	94%	92%	93%	93%	92%	92%	88%	93%	96%
Contact Energy Ltd	94%	90%	93%	93%	94%	93%	91%	93%	91%	88%	91%	94%
Fisher Paykel Healthcare Corporation Limited	95%	93%	92%	94%	92%	84%	94%	92%	93%	84%	93%	97%
Fletcher Building Ltd	96%	93%	94%	95%	94%	96%	92%	93%	94%	95%	95%	96%
Goodman Property Trust	96%	93%	91%	93%	97%	95%	94%	90%	90%	89%	#N/A	97%
Infratil Limited	95%	92%	94%	93%	94%	95%	94%	90%	92%	-89%	93%	94%
Kiwi Income Property Trust	95%	96%	95%	95%	81%	95%	96%	87%	93%	96%	95%	98%
National Australia Bank Limited	96%	94%	94%	95%	94%	94%	94%	93%	94%	89%	94%	97%
Spark New Zealand Limited	95%	94%	96%	93%	96%	88%	94%	96%	92%	89%	95%	96%
Telecom Corp of New Zealand Ltd	97%	95%	92%	94%	96%	84%	93%	88%	94%	85%	95%	97%
Vector Limited	94%	92%	93%	92%	92%	90%	91%	92%	92%	87%	92%	95%
Westpac Banking Corporation	96%	96%	94%	94%	94%	95%	94%	91%	94%	79%	93%	98%
Value	96%	94%	93%	94%	94%	91%	94%	92%	93%	80%	93%	93%

Figure 1: Results Matrix

1.5.3 Python

```

1 # Descriptions
2 # This script/function implements the StanfordNLP to score corporate culture,
  replicating the production of inputs in Option 1 as outputs
3 # Inputs for Option 1 include:
4 # 1. Firm_score.xlsx contains five scores estimated with two different dictionaries
  for all calls. Scores ended with 1 (for example, integrity1) are estimated with
  the dictionary trained on the 258 call transcripts included in this sample.
  Scores ended with 2 (for example, integrity2) are estimated with the dictionary
  from the original paper (Table IA3 in the Internet Appendix). Other variables
  include document_id (used in your coding), filename (file name used by CapitalIQ)
  , firm_id (firm name) and call time (year and month of the call).
5 # 2. Expanded_dict1.csv is the culture dictionary trained with the 258 call
  transcripts (the new dictionary).
6 # 3. Expanded_dict2.csv is the culture dictionary from the original paper (the
  original dictionary).
7 # 4. Word_contributin_TFIDF1.csv (Word_contributin_TFIDF2.csv) contains word
  contribution based on TFIDF score estimated with the new dictionary (the original
  dictionary).
8 # 5. The Li, Mai, Shen and Yan (2021) paper and the Internet Appendix of this paper.
9 #
10 # Author: Connor McDowall
11 # Date: 25th August 2021
12
13 # Imports
14 # Transcript Processing Modules
15 import pandas as pd
16 from pathlib import Path
17 import shutil as sh
18 from pdfw import PdfReader, PdfWriter
19 import pdfminer as pdm
20 from pdfminer.converter import TextConverter
21 from pdfminer.layout import LAParams
22 from pdfminer.pdfdocument import PDFDocument
23 from pdfminer.pdfinterp import PDFResourceManager, PDFPageInterpreter
24 from pdfminer.pdfpage import PDFPage
25 from pdfminer.pdfparser import PDFParser
26 import io
27 import datefinder as dtf
28 # General Python Modules
29 import datetime
30 import functools
31 import logging
32 import sys
33 import math
34 import os
35 import pickle
36 import gensim
37 import itertools
38 from pprint import pprint
39 from collections import Counter, defaultdict, OrderedDict
40 from tqdm.auto import tqdm
41 from typing import Dict, List, Optional, Set
42 from multiprocessing import Pool
43 from operator import itemgetter
44 from tqdm import tqdm as tqdm
45
46 # StanfordNLP Specific Functions
47 from culture import culture_models, file_util, preprocess, culture_dictionary,
  preprocess_parallel
48 from stanfordnlp.server import CoreNLPClient
49 import global_options
50 import parse
51 import clean_and_train
52 import create_dict
53 import score
54 import aggregate_firms
55
56 # Functions
57 def get_transcripts(firm_score_excel, transcript_directory, transcript_selected):
58     """Locates and isolates transcripts for processing
59
60     Args:
61         firm_score_excel (xlsx): Excel file containing the initial list of transcripts
62         transcript_directory (str): Source of all transcripts
63         transcript_selected (str): Destination for transcripts of interest
64
65     Returns:
66         transcript_list (list): List of transcript filenames
67         calltimes (list): List of calltimes
68     """
69     # Get list of filenames
70     firms_df = pd.read_excel(firm_score_excel)
71     firms_df = firms_df.dropna()

```

```

72 firms_df.columns = firms_df.iloc[0]
73 firms_df = firms_df.drop(2)
74 firms_df = firms_df.reset_index(drop=True)
75 transcript_list = firms_df['filename'].tolist()
76 # Get list of calltimes for the firm ID
77 calltimes = firms_df['calltime'].tolist()
78 # Copy file into selection if exists
79 files_found = 0
80 files_to_find = len(transcript_list)
81 missing_files_list = []
82 for filename in transcript_list:
83     transcript_x = Path(transcript_directory + '/' + filename)
84     if transcript_x.is_file():
85         transcript_y = Path(transcript_selected + '/' + filename)
86         sh.copy(transcript_x, transcript_y)
87         files_found = files_found + 1
88     else:
89         missing_files_list.append(filename)
90 missing_files = files_to_find - files_found
91 if missing_files > 0:
92     print('You are missing the following transcripts...')
93     print(missing_files_list)
94 else:
95     print('All transcripts found')
96 return transcript_list, calltimes
97
98 def create_ids(transcript_list, qa_num, company_ids_set, company_ids_order,
99 documents_ids_text, calltimes):
100     """Creates document identification, updates transcript list to only include
101     transcript lists
102     with Question and Answer Sections, and creates dataframe to compare results.
103     Args:
104         transcript_list (list): List of transcript filenames
105         qa_num (list): List of page numbers denoting the start of question and answer
106         sections
107         company_ids_set (list): List of company names
108         company_ids_order (list): List of numbers referencing number of file relating
109         to one company
110         documents_ids_text (str): Directory to store document id list as a text file
111         calltimes (list): List of calltimes
112     Returns:
113         updated_transcript_list (list): List of updated filenames
114         updated_document_ids (list): List of updated document ids
115         updated_firm_id (list): List of updated firm ids
116         output_df (dataframe): Dataframe with document information
117     """
118     # Initial lists
119     document_ids = []
120     firm_id = []
121     # Updated lists
122     updated_document_ids = []
123     updated_firm_id = []
124     updated_transcript_list = []
125     updated_calltimes = []
126     # Assigns document id
127     idx = 0
128     for i in range(len(transcript_list)):
129         document_ids.append(str(i + 1) + '.F')
130         if i < company_ids_order[idx]:
131             firm_id.append(company_ids_set[idx])
132         else:
133             idx = idx + 1
134             firm_id.append(company_ids_set[idx])
135     # Updates lists to remove entries with no question and answer sections
136     for j in range(len(qa_num)):
137         if qa_num[j] != 4:
138             updated_document_ids.append(document_ids[j])
139             updated_firm_id.append(firm_id[j])
140             updated_transcript_list.append(transcript_list[j])
141             updated_calltimes.append(calltimes[j])
142     # Creates document_id text file
143     with open(documents_ids_text, "w") as file:
144         # Clear the file
145         file.truncate(0)
146         for element in updated_document_ids:
147             file.write(element + "\n")
148         file.close()
149     # Creates a dataframe with updated transcript list
150     output_df = pd.DataFrame(list(zip(updated_document_ids, updated_transcript_list,
151                                     updated_firm_id)),
152                             columns=['document_id', 'filename', 'firm_id'])

```

```

150 # Creates id2firm csv
151 for i in range(len(updated_calltimes)):
152     val = updated_calltimes[i]
153     new_val = int(str(val)[:4])
154     updated_calltimes[i] = new_val
155 id2firms_df = pd.DataFrame(list(zip(updated_document_ids, updated_firm_id,
156                                     updated_calltimes)),
157                             columns=['document_id', 'firm_id', 'time'])
158 print(id2firms_df.head())
159 id2firms_df.to_csv('data/input/id2firms.csv')
160 return updated_transcript_list, updated_document_ids, updated_firm_id, output_df
161
162 def remove_transcript_metadata(transcript_list, qa_num, transcript_selected,
163                               transcript_processed):
164     """Removes front matter, table of contents, call participants, and copyright
165     disclaimer
166     to process transcripts to a format suitable for combination. This is possible as
167     the
168     format is consistent for all earnings call transcripts.
169
170     Args:
171         transcript_list (list): List of transcript filenames
172         qa_num (list): List of page numbers denoting the start of question and answer
173         sections
174         transcript_selected (str): String of selected transcript directory
175         transcript_processed (str): String of processed transcript directory
176
177     """
178     # Count for
179     i = 0
180     # Create copy, remove pages, and move to processed directory
181     for filename in transcript_list:
182         # Defines pdfs
183         input_pdf = Path(transcript_selected + '/' + filename)
184         output_pdf = Path(transcript_processed + '/' + filename)
185         # Defines objects
186         reader_input = PdfReader(input_pdf)
187         writer_output = PdfWriter()
188         for page_x in range(len(reader_input.pages)):
189             # Adds pages excluding sections prior to Q&A section and legal disclaimer
190             if page_x >= qa_num[i]-1 and page_x < (len(reader_input.pages)-1):
191                 writer_output.addpage(reader_input.pages[page_x])
192             writer_output.write(output_pdf)
193             i = i + 1
194     return
195
196 def create_documents_text(transcript_list, transcript_processed, text_processed,
197                           documents_text):
198     """Creates documents.txt file for the Stanford NLP
199
200     Args:
201         transcript_list(str): List of processed transcripts
202         transcript_processed (str): String of processed transcript directory
203         text_processed (str): Directory to store text file
204         documents_text (str): Directory for documents.txt file
205
206     Returns:
207         documents_test_list (list): Returns a list of processed transcript document
208         strings
209
210     """
211     # Adapted from https://towardsdatascience.com/pdf-text-extraction-in-python-5
212     # b6ab9e92dd
213     # Erase object contents to reset the textfile
214     with open(documents_text, "r+") as file:
215         file.truncate(0)
216         file.close()
217     # Creates empty list
218     documents_test_list = []
219     # Begin extracting files
220     for file_name in transcript_list:
221         file_pdf = Path(transcript_processed + '/' + file_name)
222         file_text = io.StringIO()
223         with open(file_pdf, 'rb') as in_file:
224             parser = PDFParser(in_file)
225             doc = PDFDocument(parser)
226             rsrcmgr = PDFResourceManager()
227             device = TextConverter(rsrcmgr, file_text, laparams=LAParams())
228             interpreter = PDFPageInterpreter(rsrcmgr, device)
229             for page in PDFPage.create_pages(doc):
230                 interpreter.process_page(page)
231             # Extract text to and remove characters
232             textname = Path(text_processed + '/output.txt')
233             with open(textname, "w") as file:
234                 file.write(file_text.getvalue())

```



```

225         file.close()
226         # Print the lines
227         with open(textname, "r+") as file:
228             line = file.read().replace("\n", " ")
229             file.truncate(0)
230             file.close()
231         # Write line to the documents file
232         with open(documents_text, "a") as file:
233             file.write(line)
234             if file_name != transcript_list[-1]:
235                 file.write("\n")
236             file.close()
237         # Create list of texts and dates
238         documents_test_list.append(line)
239     return documents_test_list
240
241 def prepare_documents(firm_score_xlsx, transcript_directory, transcript_selected,
242                      transcript_processed, text_processed, documents_text, documents_ids_text, qa_num,
243                      company_ids_set, company_ids_order):
244     """ Isolate transcripts of interest, process Q&A sections, and create document
245     files
246     Args:
247         firm_score_xlsx (xlsx): Excel file containing the initial list of transcripts
248         transcript_directory (str): Source of all transcripts
249         transcript_selected (str): Destination for transcripts of interest
250         transcript_processed (str): Directory for processed transcripts
251         text_processed (str): Directory to store text file
252         documents_text (str): Directory for documents.txt file
253         documents_ids_text (str): Directory to store document id list as a text file
254         qa_num (list): List of page numbers denoting the start of question and answer
255         sections
256         company_ids_set (list): List of company names
257         company_ids_order (list): List of numbers referencing number of file relating
258         to one company
259     Returns:
260         documents_test_list (list): Returns a list of processed transcript document
261         strings
262         document_ids (list): List of document ids
263         firm_id (list): List of firm ids
264         output_df (df): Dataframe with document information
265     """
266     # Prepares the documentation
267     # Get list of transcripts
268     transcript_list, calltimes = get_transcripts(firm_score_xlsx, transcript_directory,
269         , transcript_selected)
270     # Isolates Q&A sections while removing legal disclaimers
271     remove_transcript_metadata(transcript_list, qa_num, transcript_selected,
272         transcript_processed)
273     # Creates supplementary identification (Changed here to remove files without text
274     files)
275     transcript_list, document_ids, firm_id, output_df = create_ids(transcript_list,
276         qa_num, company_ids_set, company_ids_order, documents_ids_text, calltimes)
277     # Creates the documents.txt file, documents ids, firm_ids, and dataframe of
278     outputs
279     documents_test_list = create_documents_text(transcript_list, transcript_processed,
280         text_processed, documents_text)
281     # Saves csv for comparison
282     dataframe_file = Path('data/input/results.csv')
283     output_df.to_csv(dataframe_file)
284     return documents_test_list, document_ids, firm_id, output_df
285
286 def perform_stanford_nlp():
287     """Executes Stanford NLP algorithm on processed documentation via
288     """
289     print("Implementing Stanford NLP...")
290     # Creates variables and directories in global options
291     exec(open("global_options.py").read())
292     # Step 1: Use 'python parse.py' to use Stanford CoreNLP to parse the raw
293     documents.
294     exec(open("parse.py").read())
295     # Step 2: Use 'python clean_and_train.py' to clean, remove stopwords, and named
296     entities in parsed 'documents.txt'
297     exec(open("clean_and_train.py").read())
298     # Step 3: Use 'python create_dict.py' to create the expanded dictionary.
299     exec(open("create_dict.py").read())
300     # Step 4: Use 'python score.py' to score the documents.
301     exec(open("score.py").read())
302     # Step 5: Use 'python aggregate_firms.py' to aggregate the scores to the firm-
303     time level.
304     exec(open("aggregate_firms.py").read())
305     return

```

```

293
294 def compare_results(results,output_scores):
295     """Creates comparison excel sheets with helens results
296
297     Args:
298         results (str): Directory to the document id files
299         output_scores (str): Directory for scoring sheets
300     """
301     # Load in the results
302     output_df = pd.read_csv(results)
303     # Set directories
304     tf = 'firm_scores_TF.csv'
305     tfidf = 'firm_scores_TFIDF.csv'
306     wfidf = 'firm_scores_WFIDF.csv'
307     helen_results = 'outputs/scores/firm_score_helen.xlsx'
308     firm_scores_tf = Path(output_scores+'/'+tf)
309     firm_scores_tfidf = Path(output_scores+'/'+tfidf)
310     firm_scores_wfidf = Path(output_scores+'/'+wfidf)
311     helen_results = Path(helen_results)
312     # Read csv and excel files
313     firm_scores_tf_df = pd.read_csv(firm_scores_tf)
314     firm_scores_tfidf_df = pd.read_csv(firm_scores_tfidf)
315     firm_scores_wfidf_df = pd.read_csv(firm_scores_wfidf)
316     helen_results = pd.read_excel(helen_results)
317     # Merge results with dataframes for comparison
318     target_df = firm_scores_tfidf_df
319     user_results_df = pd.merge(output_df, target_df,how = 'left',on = output_df.
index)
320     comparison_df = pd.merge(user_results_df,helen_results,how = 'left',on = ['
document_id'])
321     print('Please enter a filename')
322     filename = input()
323     # Save comparison csv
324     file_string = 'outputs/comparisons'++'/'+filename+'.xlsx'
325     comparison_df.to_excel(file_string)
326     return
327
328 # Inputs - established all the directories for the locations
329 # Inputs for processing
330 firm_score_xlsx = 'data/input/option-1/1.firm_score.xlsx'
331 transcript_directory = 'data/input/transcripts'
332 transcript_selected = 'data/raw/selected_transcripts'
333 transcript_processed = 'data/processed/processed_transcripts'
334 text_processed = 'data/processed/processed_text'
335 documents_text = 'data/input/documents.txt'
336 documents_ids_text = 'data/input/document_ids.txt'
337 # Creates array of pages numbers indicating the start of the Q&A section for each PDF
338 # Note: This is labourous but necessary. Values of 4 indicate no Q&A section in the
document,
339 # starting at the presentation section
340 air_nz_num=[8,10,10,7,8,10,8,11,8,8,8]
341 aia_num = [4,4,12,12,12,9,10,15,11,10,10,10] # Changed to 4
342 anz_num =
[14,6,10,11,11,13,11,13,11,10,11,13,13,8,7,8,7,10,11,12,13,11,12,10,11,11,13,12,11,12,8]
343 bql_num =
[24,12,10,11,11,12,11,11,14,11,9,16,12,13,16,14,13,12,13,10,10,11,12,12,12,13,8]
344 bab_num = [4,10,10,12,11,10,10,10,14,15,10,10,10,12,10,10,12,12,15,15,9,10]
345 cba_num =
[5,11,11,12,12,11,12,11,10,10,4,11,11,10,11,11,11,10,12,12,11,12,12,12,6,6,6,7,8]
# Changed to 4 (29)
346 ce_num = [8,4] # Changed to 4
347 fph_num = [9,8,9,8,9,8,8,7,8,8,8]
348 fbu_num = [12,10,11,10,9,10,10,9,10,11,9]
349 gpt_num = [10,9]
350 il_num = [15,15,15,15,13,14,13,12,13,15,15,16,13]
351 kip_num = [11,10]
352 nab_num =
[12,4,4,13,12,14,10,18,15,9,10,11,11,12,13,13,10,12,15,10,9,10,10,12,11,9,8,15,7,7,6]
# Changed to 4 (31)
353 spk_num = [16,12,12]
354 tnz_num = [15,14,11,16,14,12,13,9,16]
355 vec_num = [9,12,9,9,10,9,9,8,12,11,10,9]
356 wpc_num =
[12,19,12,14,14,13,13,11,12,11,11,12,11,11,16,14,12,12,12,11,11,12,10,11,7,8,7,7,7]
357 # Combines the arrays
358 qa_num = [air_nz_num,
359 aia_num,
360 anz_num,
361 bql_num,
362 bab_num,
363 cba_num,
364 ce_num,

```

```

365         fph_num,
366         fbu_num,
367         gpt_num,
368         il_num,
369         kip_num,
370         nab_num,
371         spk_num,
372         tnz_num,
373         vec_num,
374         wpc_num]
375
376 qa_num = air_nz_num+aia_num+anz_num+bql_num+bab_num+cba_num+ce_num+fph_num+fbu_num+
          gpt_num+il_num+kip_num+nab_num+spk_num+tnz_num+ vec_num + wpc_num
377 # Sets list for company ids
378 company_ids_set = ['Air New Zealand Limited','Auckland International Airport Limited'
          , 'Australia New Zealand Banking Group Limited', 'Bank of Queensland Limited',
          Bendigo and Adelaide Bank Limited', 'Commonwealth Bank of Australia', 'Contact
          Energy Ltd', 'Fisher Paykel Healthcare Corporation Limited', 'Fletcher Building Ltd
          ', 'Goodman Property Trust', 'Infratil Limited', 'Kiwi Income Property Trust',
          National Australia Bank Limited', 'Spark New Zealand Limited', 'Telecom Corp of New
          Zealand Ltd', 'Vector Limited', 'Westpac Banking Corporation']
379 company_ids_order = [11,24,55,82,104,133,135,146,157,159,172,174,205,208,217,229,258]
380 # Inputs for comparison
381 output_scores = 'outputs/scores'
382 results = 'data/input/results.csv'
383 output_word_contributions = 'outputs/scores/word_contributions'
384 firm_scores_tf = 'outputs/scores/firm_scores_TF.csv'
385 firm_scores_tfidf = 'outputs/scores/firm_scores_TFIDF.csv'
386 firm_scores_wfidf = 'outputs/scores/firm_scores_WFIDF.csv'
387 #####
388 # Function Calls
389 # Set binary variables to control function calls
390 transcript_preparation = False
391 stanford_nlp_implementation = False
392 results_comparison = True
393 # Executes functions based on binary variables
394 if transcript_preparation == True:
395     # Prepare the documents
396     print("Preparing documents...")
397     documents_test_list, document_ids, firm_id, output_df = prepare_documents(
          firm_score_xlsx, transcript_directory, transcript_selected, transcript_processed,
          text_processed, documents_text, documents_ids_text, qa_num, company_ids_set,
          company_ids_order)
398 if stanford_nlp_implementation == True:
399     # Implements Stanford NLP
400     perform_stanford_nlp()
401 if results_comparison == True:
402     print('Comparing results...')
403     compare_results(results,output_scores)
404 # Note: Australia and New Zealand Banking Group Limited - ShareholderAnalyst Call.pdf
          , Bank of Queensland Ltd. - ShareholderAnalyst Call.pdf
405 # Commonwealth Bank of Australia - ShareholderAnalyst Call.pdf, Infratil Limited -
          AnalystInvestor Day.pdf, Infratil Ltd. - AnalystInvestor Day.pdf
406 # National Australia Bank Limited - ShareholderAnalyst Call.pdf

```