

```
[1]: pip install ntk matplotlib seaborn wordcloud
Defaulting to user installation because normal site-packages is not writable
Collecting ntk
  Using cached ntk-3.9.1-py3-none-any.whl.metadata (2.9 kB)
Requirement already satisfied: matplotlib in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (3.9.4)
Collecting seaborn
  Using cached seaborn-0.13.2-py3-none-any.whl.metadata (5.4 kB)
Collecting wordcloud
  Downloading wordcloud-1.9.4-cp39-cp39-macosx_11_0_arm64.whl.metadata (3.4 kB)
Collecting click (from ntk)
  Using cached click-8.1.8-py3-none-any.whl.metadata (2.3 kB)
Requirement already satisfied: joblib in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from ntk) (1.4.2)
Collecting regex<=2021.8.3 (from ntk)
  Downloading regex-2024.11.6-cp39-cp39-macosx_11_0_arm64.whl.metadata (40 kB)
Requirement already satisfied: tqdm in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from ntk) (4.67.1)
Requirement already satisfied: contourpy>=1.0.1 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (1.3.0)
Requirement already satisfied: cycler>=0.10 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (4.55.3)
Requirement already satisfied: kiwisolver>=1.3.3 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (1.4.7)
Requirement already satisfied: numpy>=1.23 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (2.0.2)
Requirement already satisfied: packaging>=20.0 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (11.1.0)
Requirement already satisfied: pyparsing>=3.1 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (3.2.1)
Requirement already satisfied: python-dateutil>=2.7 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (2.9.0.post0)
Requirement already satisfied: importlib-resources>=3.2.0 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from matplotlib) (6.5.2)
Requirement already satisfied: zipp>=3.2 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from importlib-resources>=3.2.0-matplotlib) (3.21.0)
Requirement already satisfied: pytz>=2020.1 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from pandas>=1.2-seaborn) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /Users/madhusudan/Library/Python/3.9/lib/python/site-packages (from pandas>=1.2-seaborn) (2024.2)
Requirement already satisfied: six>=1.5 in /Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/lib/python3.9/site-packages (from python-dateutil>=2.7-matplotlib) (1.15.0)
Using cached ntk-3.9.1-py3-none-any.whl (1.5 MB)
Using cached seaborn-0.13.2-py3-none-any.whl (129 kB)
Downloading wordcloud-1.9.4-cp39-cp39-macosx_11_0_arm64.whl (168 kB)
Downloading regex-2024.11.6-cp39-cp39-macosx_11_0_arm64.whl (284 kB)
Using cached click-8.1.8-py3-none-any.whl (98 kB)
Installing collected packages: regex, click, ntk, wordcloud, seaborn
Successfully installed click-8.1.8 ntk-3.9.1 regex-2024.11.6 seaborn-0.13.2 wordcloud-1.9.4
Note: you may need to restart the kernel to use updated packages.
```

```
In [2]: import os
import nltk
from collections import Counter
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
```

```
%matplotlib inline
```

```
[nltk_data] Downloading package punkt to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

```

Out[2]: True

In [3]: def process_file(file_path):
    """
    Reads a text file, removes Gutenberg boilerplate (header and footer), tokenizes the text,
    converts tokens to lowercase, and filters out non-alphabetic tokens.

    Returns:
        tokens (list): A list of cleaned, alphabetic tokens.
    """
    with open(file_path, "r", encoding="utf-8") as f:
        raw_text = f.read()

    start_marker = "*** START OF THIS PROJECT GUTENBERG EBOOK"
    end_marker = "*** END OF THIS PROJECT GUTENBERG EBOOK"

    start_idx = raw_text.find(start_marker)
    if start_idx != -1:
        text = raw_text[start_idx + len(start_marker):]
    else:
        text = raw_text # If no marker found, use entire text

    end_idx = text.find(end_marker)
    if end_idx != -1:
        text = text[:end_idx]

    text = text.strip()

    tokens = nltk.word_tokenize(text.lower())

    tokens = [token for token in tokens if token.isalpha()]

    return tokens

```

```
Test file tokens (first 20)
```

```
all_files = [f for f in os.listdir(folder) if f.endswith(".txt.txt")]
print(f"Found {len(all_files)} files.")
```

```
aggregated_counter = Counter()
file_stats = [] # This will hold stats for each file
```

```
for filename in all_files:
    file_path = os.path.join(folder, filename)
    tokens = process_file(file_path)

    aggregated_counter.update(tokens)

file_stats.append({
    "filename": filename,
    "num_tokens": len(tokens),
    "unique_tokens": len(set(tokens))
})
```

```
df_stats = pd.DataFrame(file_stats)
print("Per-file statistics (first 5 rows):")
display(df_stats.head())
```

```
Found 2475 files.
Per file statistics (first 5 rows):
```

	filename	num_tokens	unique_tokens
0	4658.txt.txt	161293	15819
1	37009.txt.txt	77551	4991
2	14609.txt.txt	89507	11561
3	5342.txt.txt	87301	6824
4	17.txt.txt	268340	5539

```
In [5]: most_common_all = aggregated_counter.most_common(20)
print("Aggregated Top 20 words:")
print(most_common_all)
```

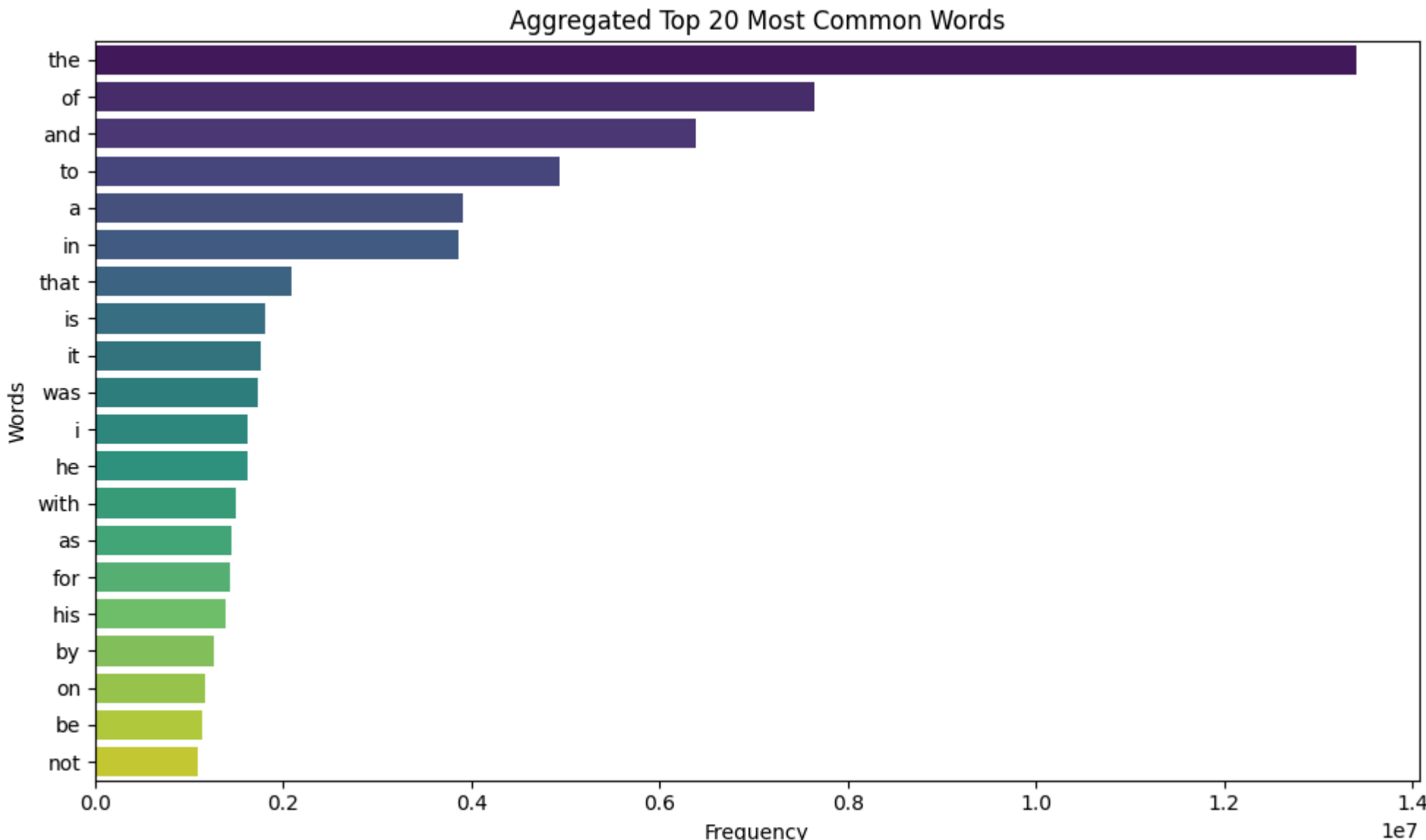
```
words, counts = zip(*most_common_all)
```

```
plt.figure(figsize=(10, 6))
sns.barpplot(x=list(counts), y=list(words), palette="viridis")
plt.title("Aggregated Top 20 Most Common Words")
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()
```

AgarAdapted Top 20 words:
 ('the', 7642031), ('and', 6378766), ('to', 6349662), ('is', 3857791), ('that', 2092342), ('is', 1800929), ('it', 1767951), ('was', 172332), ('I', 1623080), ('he', 1620363), ('with', 1496503), ('as', 1448598), ('for', 1433725), ('his', 1393159), ('by', 1266249), ('on', 1171581), ('be', 1141795), ('not', 3094746).
 /var/folders/jr/vr3v7j0kbbks_ssttqccnqd7/T/tipykenml_22400_2036935065.py:11: FutureWarning:

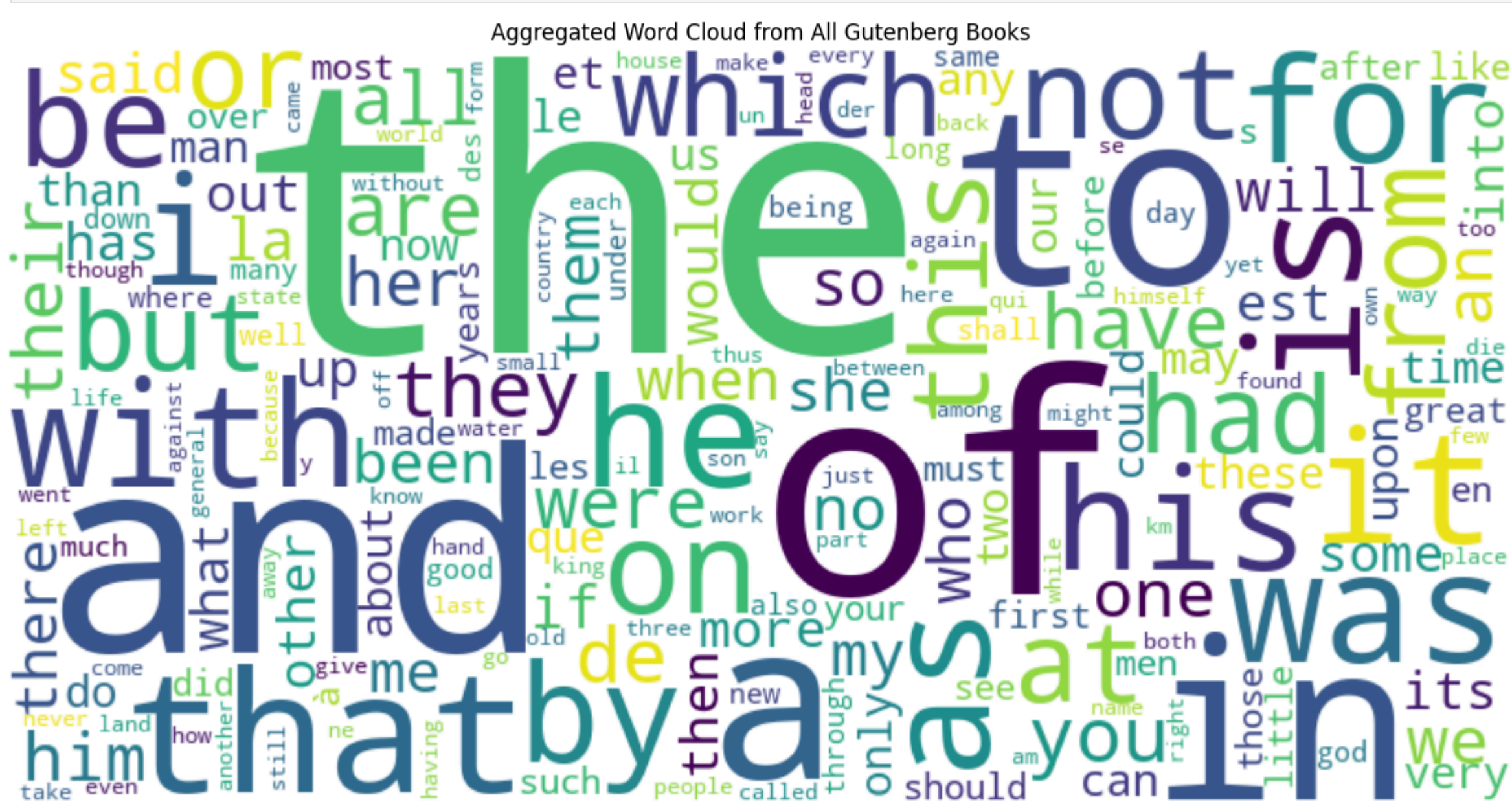
Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'y' variable to 'hue' and set 'legend=False' for the same effect.

```
sns.barplot(x=list(counts), y=list(words), palette="viridis")
```



```
In [6]: wordcloud = WordCloud(width=800, height=400, background_color='white').generate_from_frequencies(aggregated_counter)
```

```
plt.figure(figsize=(15, 8))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.title("Aggregated Word Cloud from All Gutenberg Books",
plt.show()
```



```
In [7]: df_stats['lexical_diversity'] = df_stats['unique_tokens'] / df_stats['num_tokens']
```

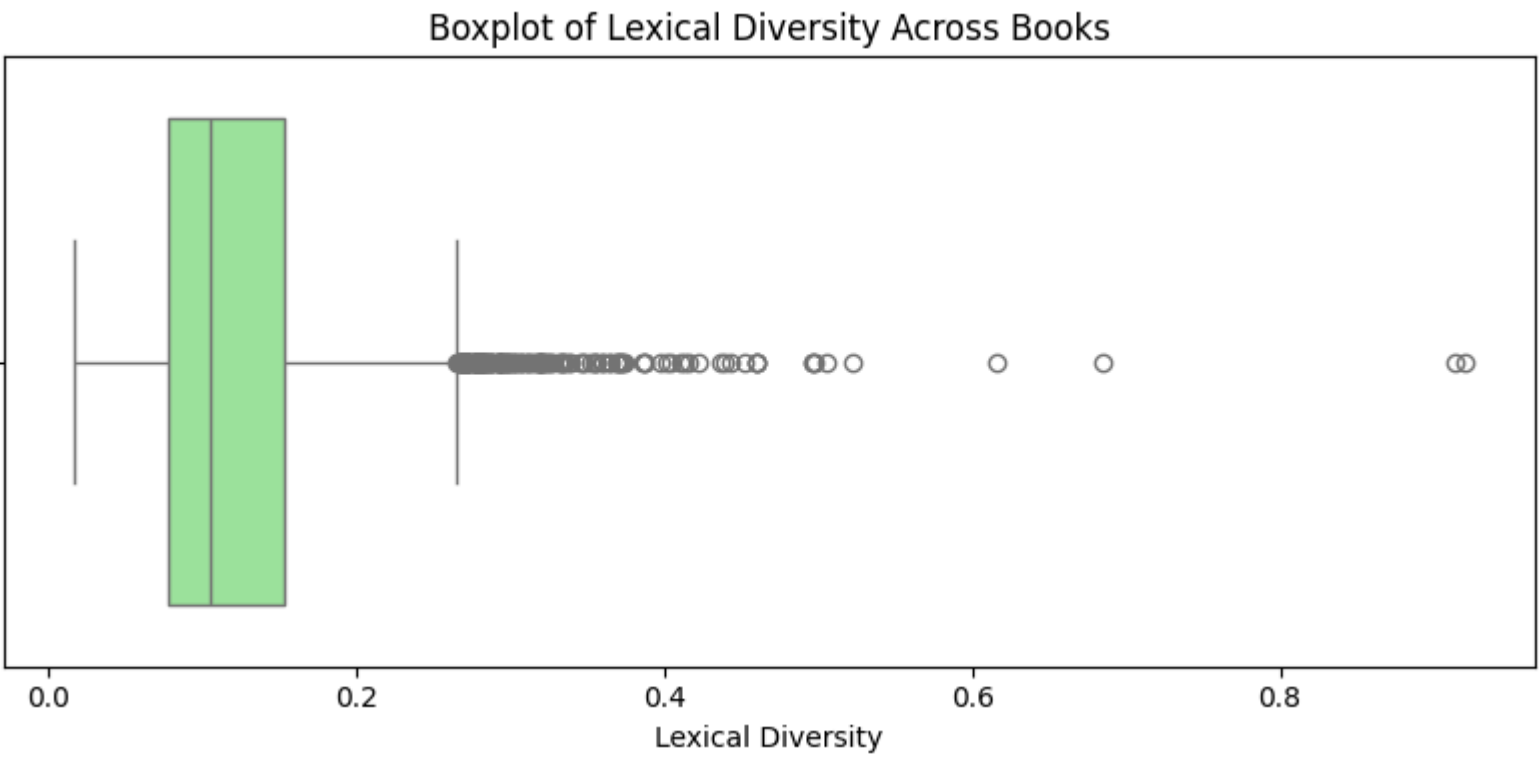
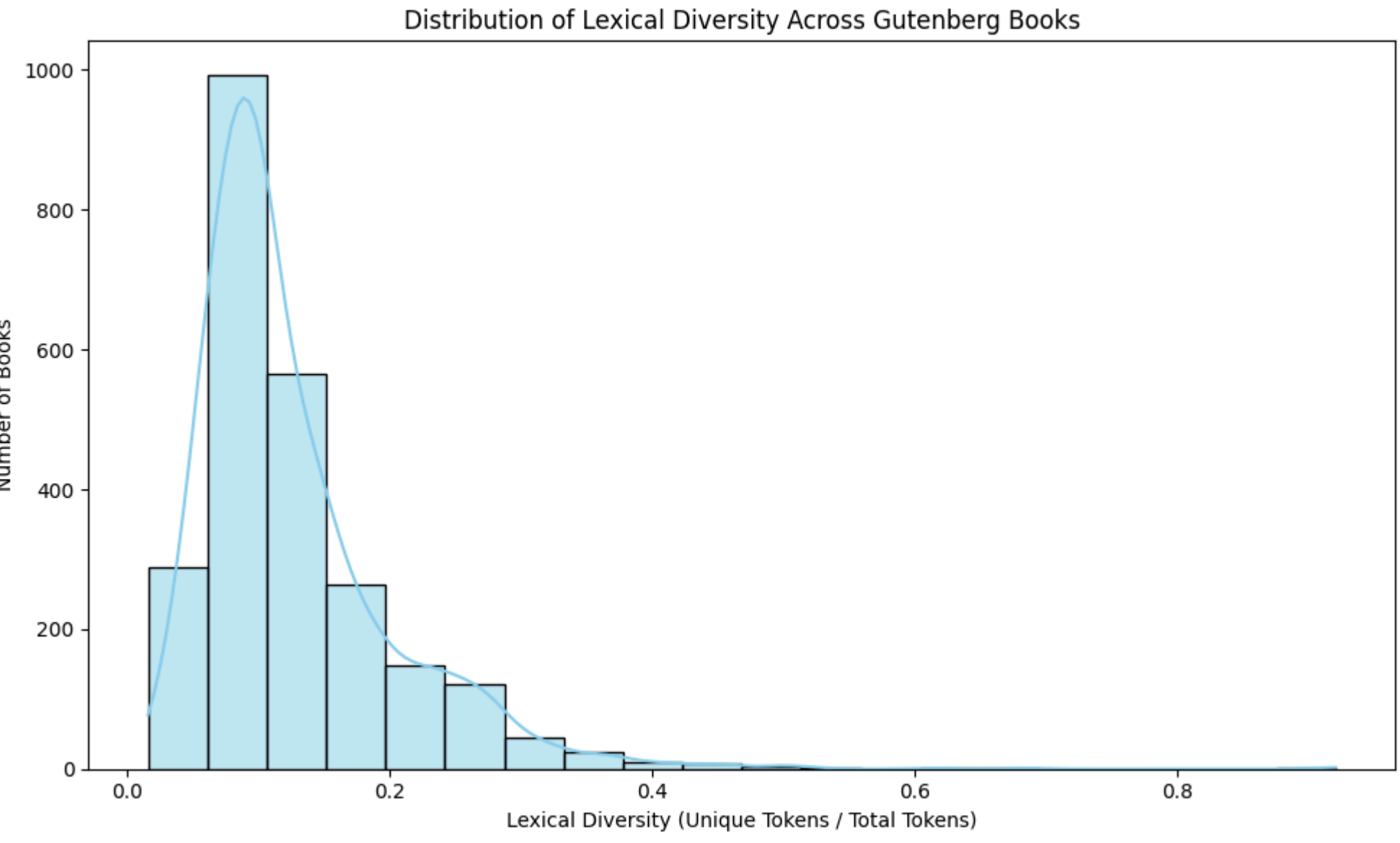
```
display(df_stats.head())
```



```
plt.figure(figsize=(10, 6))
sns.histplot(df_stats['lexical_diversity'], kde=True, bins=20, color='skyblue')
plt.title("Distribution of Lexical Diversity Across Gutenberg Books")
plt.xlabel("Lexical Diversity (Unique Tokens / Total Tokens)")
plt.ylabel("Number of Books")
plt.tight_layout()
plt.show()

plt.figure(figsize=(8, 4))
sns.boxplot(x=df_stats['lexical_diversity'], color='lightgreen')
plt.title("Boxplot of Lexical Diversity Across Books")
plt.xlabel("Lexical Diversity")
plt.tight_layout()
plt.show()
```

	filename	num_tokens	unique_tokens	lexical_diversity
0	4658.txt.txt	161293	15819	0.098076
1	37009.txt.txt	77651	4991	0.064358
2	14609.txt.txt	89507	11561	0.129163
3	5342.txt.txt	87301	6824	0.078166
4	17.txt.txt	268340	5539	0.020642



```
In [8]: all_tokens = []
for filename in all_files:
    file_path = os.path.join(folder, filename)
    tokens = process_file(file_path)
    all_tokens.extend(tokens)

print("Total tokens collected from all files:", len(all_tokens))

Total tokens collected from all files: 209085770
```

```
In [9]: import nltk
nltk.download('averaged_perceptron_tagger')

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
```

Out[9]: True

```
In [10]: print(nltk.data.path)
nltk.download('punkt')

['/Users/mmadhusudan/nltk_data', '/Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/nltk_data', '/Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/share/nltk_data', '/usr/share/nltk_data', '/usr/local/sh
are/nltk_data', '/usr/lib/nltk_data', '/usr/local/lib/nltk_data']

[nltk_data] Downloading package punkt to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

Out[10]: True

```
In [11]: import nltk
nltk.download('averaged_perceptron_tagger')

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
```

Out[11]: True

```
In [12]: print(nltk.data.find('taggers/averaged_perceptron_tagger'))

/Users/mmadhusudan/nltk_data/taggers/averaged_perceptron_tagger
```

```
In [13]: import random

sample_size = min(5000, len(all_tokens))
sample_tokens = random.sample(all_tokens, sample_size)
```

```
In [14]: file_path = os.path.join("Gutenberg_Books", "1.txt.txt")
with open(file_path, "r", encoding="utf-8") as f:
    raw_text = f.read()

start_marker = "*** START OF THIS PROJECT GUTENBERG EBOOK"
end_marker = "*** END OF THIS PROJECT GUTENBERG EBOOK"
start_idx = raw_text.find(start_marker)
if start_idx != -1:
    text = raw_text[start_idx + len(start_marker):]
else:
    text = raw_text

end_idx = text.find(end_marker)
if end_idx != -1:
    cleaned_text = text[:end_idx]
else:
    cleaned_text = text

cleaned_text = cleaned_text.strip()
```

```
In [15]: # --- Sentence-Level Analysis ---

nltk.download('punkt')

sentences = nltk.sent_tokenize(cleaned_text)

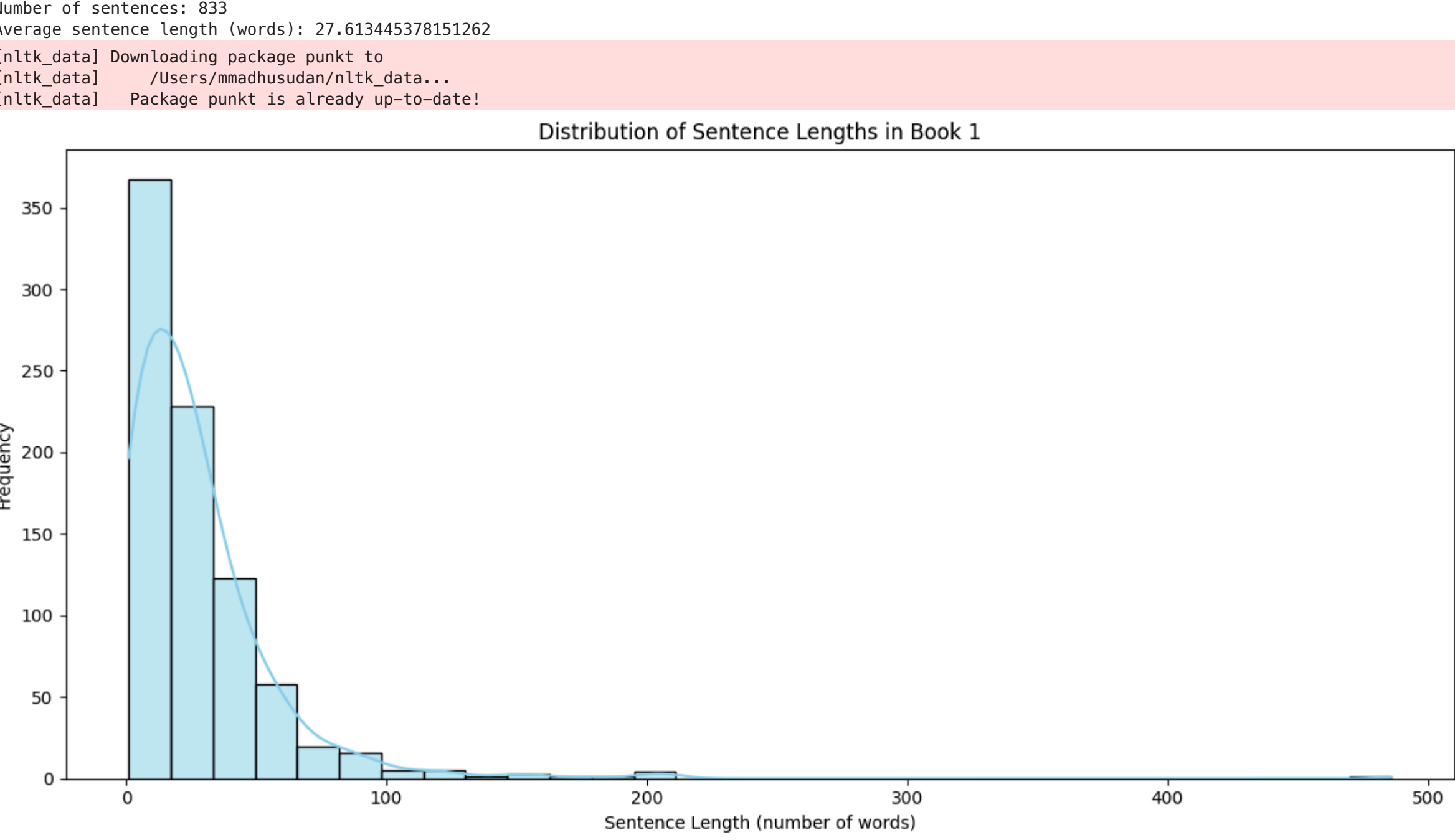
sentence_lengths = [len(nltk.word_tokenize(sentence)) for sentence in sentences]

print("Number of sentences:", len(sentences))
print("Average sentence length (words):", sum(sentence_lengths) / len(sentence_lengths))

plt.figure(figsize=(12, 6))
sns.histplot(sentence_lengths, bins=30, kde=True, color='skyblue')
plt.title("Distribution of Sentence Lengths in Book 1")
plt.xlabel("Sentence Length (number of words)")
plt.ylabel("Frequency")
plt.tight_layout()
plt.show()

Number of sentences: 833
Average sentence length (words): 27.613445378151262

[nltk_data] Downloading package punkt to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```



```
In [24]: !pip install spacy
!python -m spacy download en_core_web_sm
```

Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: spacy in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (3.8.3)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (1.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (1.0.12)
Requirement already satisfied: cytoolz<1.0.0,>=0.2.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.0.11)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (3.0.9)
Requirement already satisfied: thinc<8.4.0,>=8.3.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (8.3.4)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (1.1.3)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.5.1)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.0.10)
Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (0.4.1)
Requirement already satisfied: typer<1.0.0,>=0.3.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (0.15.1)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (4.67.1)
Requirement already satisfied: numpy<=1.19.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.0.2)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.32.3)
Requirement already satisfied: pydantic<1.8.1,>=1.8.1,<3.0.0,>=1.7.4 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (2.10.6)
Requirement already satisfied: Jinja2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (3.1.5)
Requirement already satisfied: setuptools in /Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/lib/python3.9/site-packages (from spacy) (58.0.4)
Requirement already satisfied: packaging<=20.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (24.2)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from spacy) (3.5.0)
Requirement already satisfied: language-data<=1.2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from langcodes<4.0.0,>=3.2.0->spacy) (1.3.0)
Requirement already satisfied: annotated-types<0.6.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pydantic<1.8.1,>=1.8.1,<3.0.0,>=1.7.4->spacy) (0.7.0)
Requirement already satisfied: pydantic-core<=2.27.2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pydantic<1.8.1,>=1.8.1,<3.0.0,>=1.7.4->spacy) (2.27.2)
Requirement already satisfied: typing-extensions<=4.12.2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pydantic<1.8.1,>=1.8.1,<3.0.0,>=1.7.4->spacy) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from requests<3.0.0,>=2.13.0->spacy) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from requests<3.0.0,>=2.13.0->spacy) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from requests<3.0.0,>=2.13.0->spacy) (2.3.0)
Requirement already satisfied: certifi<=2017.4.17 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from requests<3.0.0,>=2.13.0->spacy) (2024.12.14)
Requirement already satisfied: blis<1.3.0,>=1.2.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from thinc<8.4.0,>=8.3.0->spacy) (1.2.0)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from thinc<8.4.0,>=8.3.0->spacy) (0.1.5)
Requirement already satisfied: click<=8.0.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from typer<1.0.0,>=0.3.0->spacy) (8.1.8)
Requirement already satisfied: shellingsham<1.3.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from typer<1.0.0,>=0.3.0->spacy) (1.5.4)
Requirement already satisfied: rich<=10.11.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from typer<1.0.0,>=0.3.0->spacy) (13.9.4)
Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from weasel<0.5.0,>=0.1.0->spacy) (0.20.0)
Requirement already satisfied: smart-open<0.0.0,>=5.2.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from weasel<0.5.0,>=0.1.0->spacy) (7.1.0)
Requirement already satisfied: MarkupSafe<=2.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from Jinja2->spacy) (3.0.2)
Requirement already satisfied: marisa-trie<=1.1.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from language-data<=1.2->langcodes<4.0.0,>=3.2.0->spacy) (1.2.1)
Requirement already satisfied: markdown-it-py<=2.2.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from rich<=10.11.0->typer<1.0.0,>=0.3.0->spacy) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.11.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from rich<=10.11.0->typer<1.0.0,>=0.3.0->spacy) (2.19.1)
Requirement already satisfied: wrapt in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from smart-open<0.0.0,>=5.2.1->weasel<0.5.0,>=0.1.0->spacy) (1.17.2)
Requirement already satisfied: mdurl<=0.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from markdown-it-py<=2.2.0->rich<=10.11.0->typer<1.0.0,>=0.3.0->spacy) (0.1.2)
zsh:1: command not found: python

```
In [25]: !python3 -m spacy download en_core_web_sm

/Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages/urllib3/__init__.py:35: NotOpenSSLWarning: urllib3 v2 only supports OpenSSL 1.1.1+, currently the 'ssl' module is compiled with 'LibreSSL 2.8.3'. See: https://github.com/urllib3/urllib3/issues/3020
warnings.warnl
Defaulting to user installation because normal site-packages is not writeable
Collecting en-core-web-sm==3.8.0
Using cached https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.8.0/en_core_web_sm-3.8.0-py3-none-any.whl (12.8 MB)
✓ Download and installation successful
You can now load the package via spacy.load('en_core_web_sm')
```

```
In [26]: import spacy
from collections import Counter

nlp = spacy.load('en_core_web_sm')

doc = nlp(cleaned_text)

entities = [(ent.text, ent.label_) for ent in doc.entities]

entity_counts = Counter([ent.label_ for ent in doc.entities])

df_entities = pd.DataFrame(entity_counts.items(), columns=['Entity Type', 'Count']).sort_values(by='Count', ascending=False)

print("Named Entity counts:")
display(df_entities)

plt.figure(figsize=(10, 6))
sns.barplot(data=df_entities.head(10), x='Count', y='Entity Type', palette="magma")
plt.title("Top 10 Named Entity Types in Book 1")
plt.xlabel("Frequency")
plt.ylabel("Entity Type")
plt.tight_layout()
plt.show()
```

Named Entity counts:

	Entity Type	Count
4	ORG	429
5	GPE	205
3	PERSON	197
2	DATE	145
1	CARDINAL	114
6	LAW	67
8	WORK_OF_ART	56
0	ORDINAL	33
7	NORP	31
14	PRODUCT	18
15	LOC	16
9	MONEY	14
13	EVENT	8
12	FAC	7
16	TIME	6
11	PERCENT	5
10	QUANTITY	4
17	LANGUAGE	1

/var/folders/?j/rv3w7nj6k6kw_ssltcqpr0000gp/T/ipykernel_22400/556953189.py:24: FutureWarning: Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'y' variable to 'hue' and set 'legend=False' for the same effect.



```
In [27]:
```

Named Entity counts:

	Entity Type	Count
4	ORG	429
5	GPE	205
3	PERSON	197
2	DATE	145
1	CARDINAL	114
6	LAW	67
8	WORK_OF_ART	56
0	ORDINAL	33
7	NORP	31
14	PRODUCT	18
15	LOC	16
9	MONEY	14
13	EVENT	8
12	FAC	7
16	TIME	6
11	PERCENT	5
10	QUANTITY	4
17	LANGUAGE	1

/var/folders/?j/rv3w7nj6k6kw_ssltcqpr0000gp/T/ipykernel_22400/556953189.py:24: FutureWarning: Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'y' variable to 'hue' and set 'legend=False' for the same effect.



```
In [28]: from sklearn.feature_extraction.text import CountVectorizer
from sklearn.decomposition import LatentDirichletAllocation
```

```
def get_cleaned_text(file_path):
    with open(file_path, "r", encoding="utf-8") as f:
        raw_text = f.read()

    start_marker = "=== START OF THIS PROJECT GUTENBERG EBOOK"
    end_marker = "=== END OF THIS PROJECT GUTENBERG EBOOK"

    start_idx = raw_text.find(start_marker)
    if start_idx != -1:
        text = raw_text[start_idx + len(start_marker):]
    else:
        text = raw_text

    end_idx = text.find(end_marker)
    if end_idx != -1:
        text = text[:end_idx]

    # Remove extra whitespace and return
    return text.strip()
```



```
corpus = []
doc_names = [] # Keep track of file names (optional)
for filename in all_files:
    file_path = os.path.join(folder, filename)
    doc = get_cleaned_text(file_path)
    corpus.append(doc)
    doc_names.append(filename)

print(f"Collected {len(corpus)} documents.")

vectorizer = CountVectorizer(stop_words='english', max_df=0.95, min_df=2)
dtm = vectorizer.fit_transform(corpus)
print("DTM shape:", dtm.shape)

n_topics = 5
lda_model = LatentDirichletAllocation(n_components=n_topics, random_state=42)
lda_model.fit(dtm)

def print_top_words(model, feature_names, n_top_words):
    for topic_idx, topic in enumerate(model.components_):
        top_words = [feature_names[i] for i in topic.argsort()[::-n_top_words - 1:-1]]
        print(f"Topic #{topic_idx}: {' '.join(top_words)}")

n_top_words = 10
feature_names = vectorizer.get_feature_names_out()
print("\nTop words per topic:")
print_top_words(lda_model, feature_names, n_top_words)

Collected 2475 documents.
DTM shape: (2475, 526829)

Top words per topic:
Topic #0: water small time great form species large long work used
Topic #1: said man time great men did like day good little
Topic #2: la et le les il que en des qui est
Topic #3: die der en que la und el den se los
Topic #4: est 000 km years na total male female population rate
```

```
In [30]: !pip install networkx

Defaulting to user installation because normal site-packages is not writeable
Collecting networkx
  Downloading networkx-3.2.1-py3-none-any.whl.metadata (5.2 kB)
  Downloading networkx-3.2.1-py3-none-any.whl (1.6 MB)
    1.6/1.6 MB 5.2 MB/s eta 0:00:00a 0:00:01
Installing collected packages: networkx
Successfully installed networkx-3.2.1
```

```
In [41]: import networkx as nx
import matplotlib.pyplot as plt
import matplotlib.cm as cm
import numpy as np

window_size = 5 # Sliding window size
co_occurrence = {}

for i, token in enumerate(tokens):
    for j in range(i+1, min(i+window_size, len(tokens))):
        pair = tuple(sorted([token, tokens[j]]))
        co_occurrence[pair] = co_occurrence.get(pair, 0) + 1

G = nx.Graph()
threshold = 5
for pair, weight in co_occurrence.items():
    if weight >= threshold:
        G.add_edge(pair[0], pair[1], weight=weight)

print(f"Graph has {G.number_of_nodes()} nodes and {G.number_of_edges()} edges.")

deg centrality = nx.degree_centrality(G)
degrees = dict(G.degree())
node_color = [deg_centrality[node] for node in G.nodes()]
node_size = [degrees[node] * 100 for node in G.nodes()]
edge_weights = [G[u][v]['weight'] for u, v in G.edges()]
edge_width = [w / 2 for w in edge_weights]

fig, ax = plt.subplots(figsize=(15, 15))
pos = nx.spring_layout(G, k=0.15, seed=42)

nodes = nx.draw_networkx_nodes(
    G, pos, ax=ax,
    node_size=node_size,
    node_color=node_color,
    cmap=cm.viridis,
    alpha=0.9
)

sm = plt.cm.ScalarMappable(cmap=cm.viridis, norm=plt.Normalize(vmin=min(node_color), vmax=max(node_color)))
sm.set_array(np.array(node_color))
fig.colorbar(sm, ax=ax, label="Degree Centrality")

edges = nx.draw_networkx_edges(
    G, pos, ax=ax,
    width=edge_width,
    edge_color=edge_weights,
    edge_cmap=cm.plasma,
    alpha=0.7
)

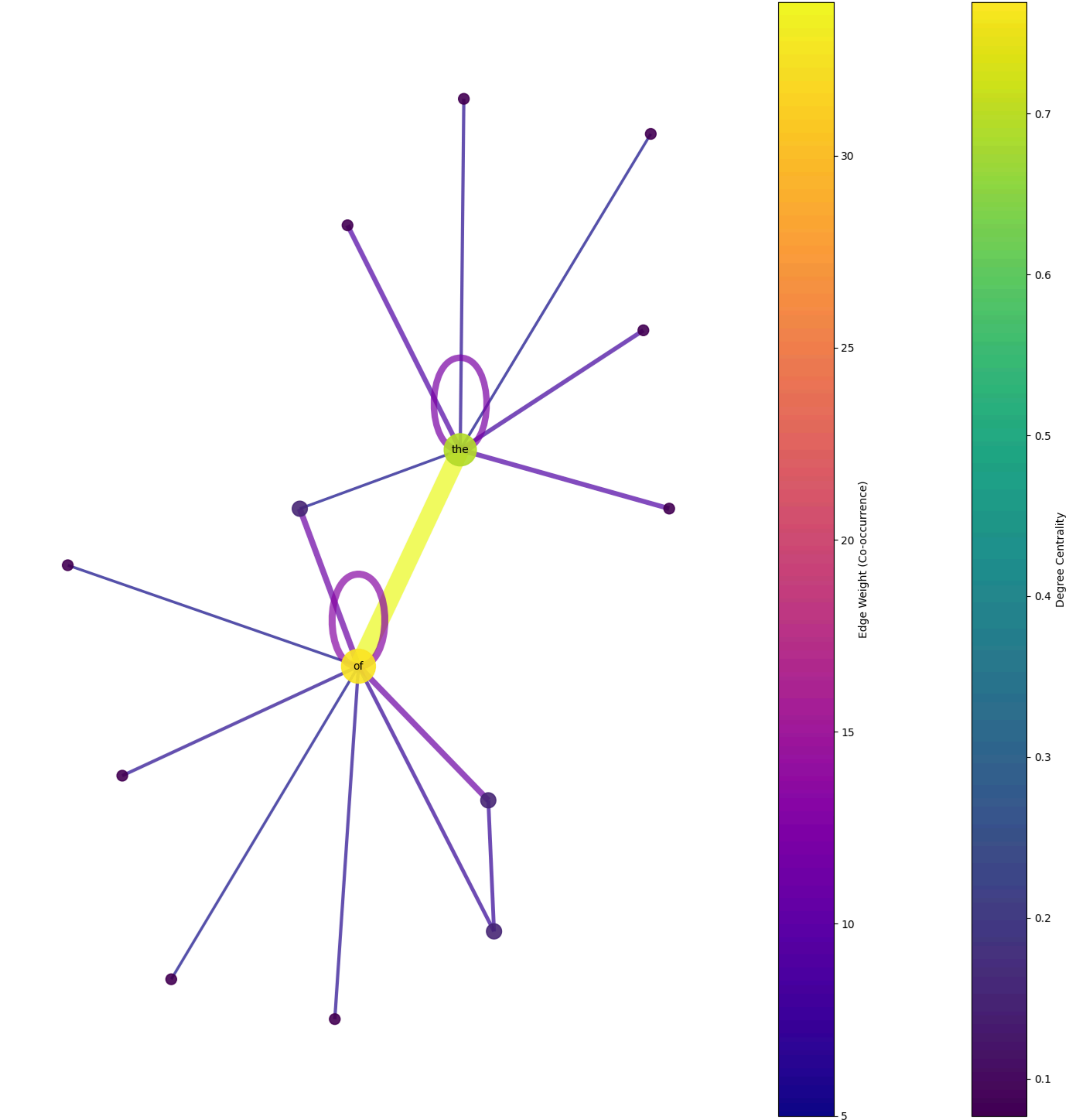
ecolor = np.array(edge_weights)
sm2 = plt.cm.ScalarMappable(cmap=cm.plasma, norm=plt.Normalize(vmin=min(ecolor), vmax=max(ecolor)))
sm2.set_array(ecolor)
fig.colorbar(sm2, ax=ax, label="Edge Weight (Co-occurrence)")

centrality_values = np.array(list(deg_centrality.values()))
threshold_label = np.percentile(centrality_values, 90)
high_central_nodes = {node: node for node in G.nodes() if deg_centrality[node] >= threshold_label}
nx.draw_networkx_labels(G, pos, labels=high_central_nodes, font_size=10, font_color='black', ax=ax)

ax.set_title("Informative Word Co-occurrence Network")
ax.axis("off")
plt.tight_layout()
plt.show()
```

Graph has 14 nodes and 17 edges.

Informative Word Co-occurrence Network



```
In [42]: import nltk

nltk.download('vader_lexicon')
from nltk.sentiment.vader import SentimentIntensityAnalyzer

sentences = nltk.sent_tokenize(cleaned_text)

sia = SentimentIntensityAnalyzer()

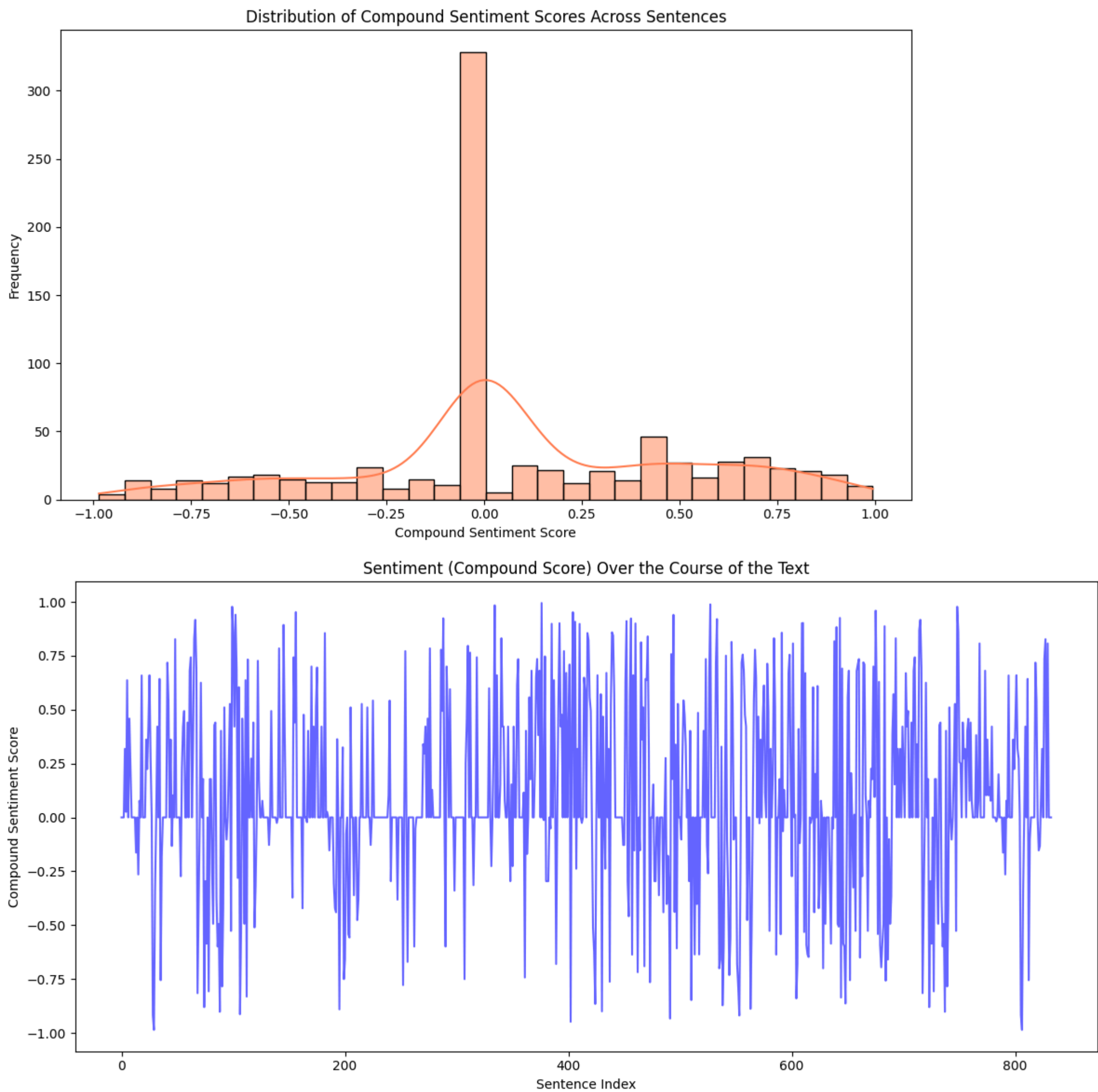
sentiment_scores = [sia.polarity_scores(sentence)['compound'] for sentence in sentences]

print("Number of sentences:", len(sentences))
print("Average compound sentiment score:", sum(sentiment_scores)/len(sentiment_scores))

plt.figure(figsize=(10, 6))
sns.histplot(sentiment_scores, bins=30, kde=True, color='coral')
plt.title("Distribution of Compound Sentiment Scores Across Sentences")
plt.xlabel("Compound Sentiment Score")
plt.ylabel("Frequency")
plt.tight_layout()
plt.show()

plt.figure(figsize=(12, 6))
plt.plot(sentiment_scores, color='blue', alpha=0.6)
plt.title("Sentiment (Compound Score) Over the Course of the Text")
plt.xlabel("Sentence Index")
plt.ylabel("Compound Sentiment Score")
plt.tight_layout()
plt.show()

Number of sentences: 833
Average compound sentiment score: 0.08476578631452589
[nltk_data] Downloading package vader_lexicon to
[nltk_data] /Users/madhusudan/nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
```



```
In [43]: from nltk.sentiment.vader import SentimentIntensityAnalyzer
import os
import nltk
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

nltk.download('punkt')
nltk.download('vader_lexicon')

def get_cleaned_text(file_path):
    with open(file_path, "r", encoding="utf-8") as f:
        raw_text = f.read()
        start_marker = "=== START OF THIS PROJECT GUTENBERG EBOOK"
        end_marker = "=== END OF THIS PROJECT GUTENBERG EBOOK"
        start_idx = raw_text.find(start_marker)
        if start_idx != -1:
            text = raw_text[start_idx + len(start_marker):]
        else:
            text = raw_text
        end_idx = text.find(end_marker)
        if end_idx != -1:
            text = text[:end_idx]
        return text.strip()

def compute_file_metrics(file_path):
    text = get_cleaned_text(file_path)
    # Word-level analysis:
    tokens = nltk.word_tokenize(text.lower())
    tokens = [t for s in tokens if t.isalpha()]
    num_tokens = len(tokens)
    num_unique = len(set(tokens))
    lexical_diversity = num_unique / num_tokens if num_tokens > 0 else 0

    # Sentence-level analysis:
    sentences = nltk.sent_tokenize(text)
    num_sentences = len(sentences)
    sentence_lengths = [len(nltk.word_tokenize(s)) for s in sentences]
    avg_sentence_length = sum(sentence_lengths)/num_sentences if num_sentences > 0 else 0

    # Sentiment analysis using VADER:
    sia = SentimentIntensityAnalyzer()
    sentiment_scores = [sia.polarity_scores(s)['compound'] for s in sentences]
    avg_sentiment = sum(sentiment_scores)/num_sentences if num_sentences > 0 else 0

    return {
        'file': os.path.basename(file_path),
        'num_tokens': num_tokens,
        'num_unique': num_unique,
        'lexical_diversity': lexical_diversity,
        'num_sentences': num_sentences,
        'avg_sentence_length': avg_sentence_length,
        'avg_sentiment': avg_sentiment
    }

folder = "Gutenberg_Books"
all_files = [f for f in os.listdir(folder) if f.endswith(".txt.txt")]

metrics_list = []
for filename in all_files:
    file_path = os.path.join(folder, filename)
    try:
        metrics = compute_file_metrics(file_path)
        metrics_list.append(metrics)
    except Exception as e:
        print(f"Error processing {filename}: {e}")

df_metrics = pd.DataFrame(metrics_list)
print("Per-file Metrics:")
display(df_metrics)

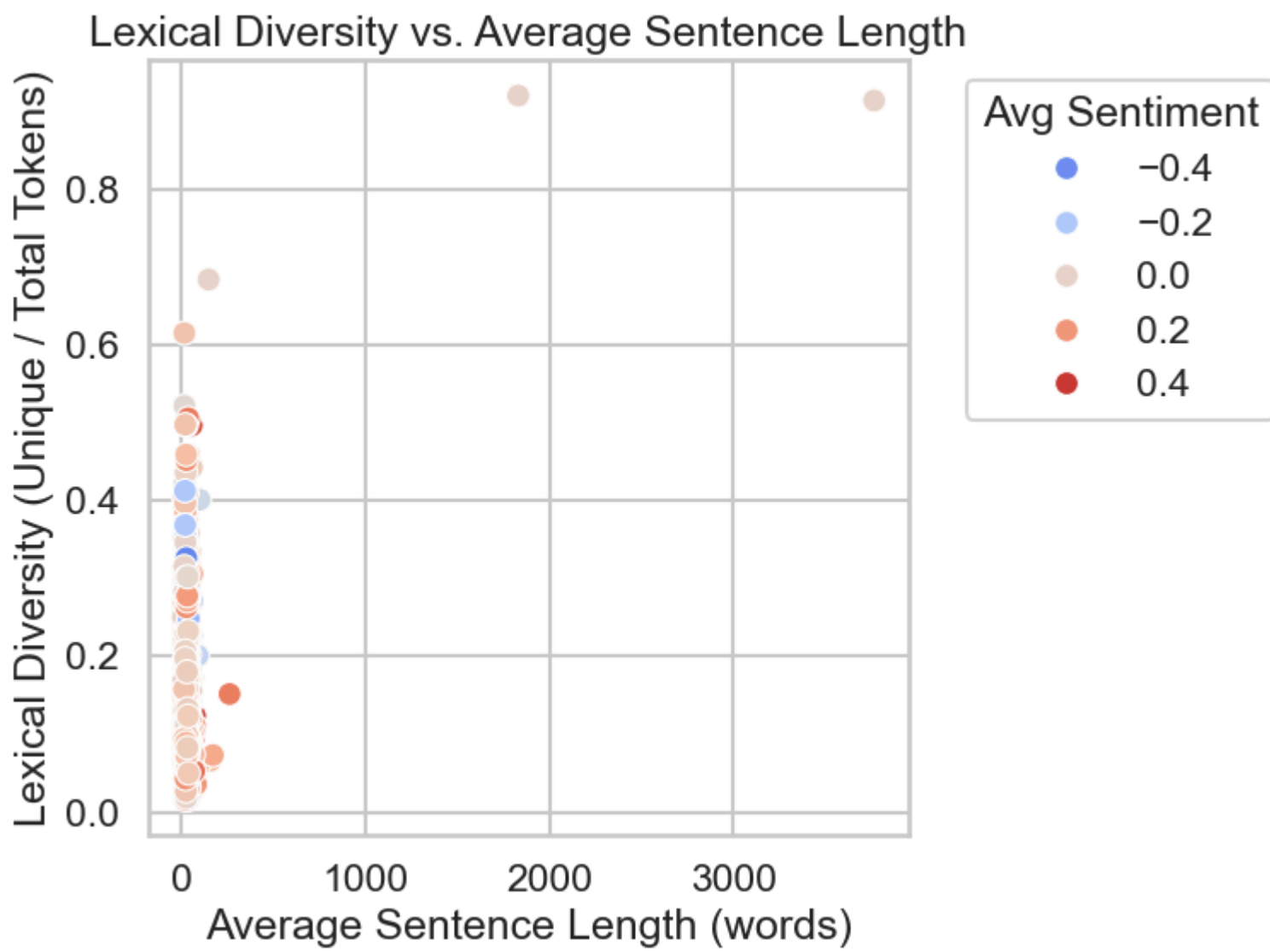
sns.set(style="whitegrid", context="talk")
plt.figure(figsize=(8,6))
sns.scatterplot(data=df_metrics, x='avg_sentence_length', y='lexical_diversity',
                hue='avg_sentiment', palette='coolwarm', s=100)
plt.title("Lexical Diversity vs. Average Sentence Length")
plt.xlabel("Average Sentence Length (words)")
plt.ylabel("Lexical Diversity (Unique / Total Tokens)")
plt.legend(title="Avg Sentiment", bbox_to_anchor=(1.05, 1), loc=2)
plt.tight_layout()
plt.show()

plt.figure(figsize=(8,6))
sns.scatterplot(data=df_metrics, x='num_sentences', y='num_tokens',
                hue='lexical_diversity', palette='viridis', s=100)
plt.title("Number of Sentences vs. Total Tokens")
plt.xlabel("Number of Sentences")
plt.ylabel("Total Tokens")
plt.legend(title="Lexical Diversity", bbox_to_anchor=(1.05, 1), loc=2)
plt.tight_layout()
plt.show()

[nltk_data] Downloading package punkt to
[nltk_data] /Users/madhusudan/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package vader_lexicon to
[nltk_data] /Users/madhusudan/nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
Per-file Metrics:
```

	file	num_tokens	num_unique	lexical_diversity	num_sentences	avg_sentence_length	avg_sentiment
0	4658.txt.txt	161293	15819	0.098076	8086	26.728172	0.040154
1	37009.txt.txt	77551	4991	0.064358	3295	27.069803	0.119815
2	14609.txt.txt	89507	11561	0.129163	4501	25.416130	0.013161
3	5342.txt.txt	87301	6824	0.078166	6239	17.282898	0.054297
4	17.txt.txt	268340	5539	0.020642	7676	40.448801	0.062850
...
2470	55836.txt.txt	57012	7572	0.132814	2200	31.482727	0.023423
2471	1452.txt.txt	132776	10999	0.082839	5061	30.670223	0.032598
2472	10061.txt.txt	3575	1081	0.302378	143	31.643357	-0.012249
2473	8395.txt.txt	18416	2276	0.123588	626	34.236422	0.035615
2474	31011.txt.txt	462	230	0.497835	33	18.515152	0.074618

2475 rows x 7 columns



Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: pip in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (25.0)
Requirement already satisfied: setuptools in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (75.8.0)
Requirement already satisfied: wheel in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (0.45.1)
Defaulting to user installation because normal site-packages is not writeable
Collecting gensim==4.3.0
 Downloading gensim-4.3.0.tar.gz (23.3 MB)
 23.3/23.3 MB 2.3 MB/s eta 0:00:0000:0100:01
 Preparing metadata (setup.py) ... done
Requirement already satisfied: numpy<=1.18.5 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from gensim==4.3.0) (1.26.4)
Requirement already satisfied: scipy<=1.7.0 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from gensim==4.3.0) (1.13.1)
Requirement already satisfied: smart_open<=1.8.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from gensim==4.3.0) (7.1.0)
Collecting FuzzyTM<=0.4.0 (from gensim==4.3.0)
 Downloading FuzzyTM-2.0.9-py3-none-any.whl.metadata (7.9 kB)
Requirement already satisfied: pandas in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from FuzzyTM<=0.4.0->gensim==4.3.0) (2.2.3)
Collecting pyfume (from FuzzyTM<=0.4.0->gensim==4.3.0)
 Downloading pyFUME-0.3.4-py3-none-any.whl.metadata (9.7 kB)
Requirement already satisfied: wrapt in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from smart_open<=1.8.1->gensim==4.3.0) (1.17.2)
Requirement already satisfied: python-dateutil<=2.8.2 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pandas->FuzzyTM<=0.4.0->gensim==4.3.0) (2.9.0.post0)
Requirement already satisfied: pytz<=2020.1 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pandas->FuzzyTM<=0.4.0->gensim==4.3.0) (2024.2)
Requirement already satisfied: tzdata<=2022.7 in /Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages (from pandas->FuzzyTM<=0.4.0->gensim==4.3.0) (2024.2)
Collecting scipy<=1.7.0 (from gensim==4.3.0)
 Downloading scipy-1.10.1-cp39-cp39-macosx_12_0_arm64.whl.metadata (53 kB)
Collecting numpy<=1.18.5 (from gensim==4.3.0)
 Downloading numpy-1.24.4-cp39-cp39-macosx_11_0_arm64.whl.metadata (5.6 kB)
Collecting simpful<=2.12.0 (from pyfume->FuzzyTM<=0.4.0->gensim==4.3.0)
 Downloading simpful-2.12.0-py3-none-any.whl.metadata (4.8 kB)
Collecting fst-ps<=1.8.1 (from pyfume->FuzzyTM<=0.4.0->gensim==4.3.0)
 Downloading fst-psy-1.8.1.tar.gz (10 kB)
 Preparing metadata (setup.py) ... done
Collecting pandas (from FuzzyTM<=0.4.0->gensim==4.3.0)
 Downloading pandas-1.5.3-cp39-cp39-macosx_11_0_arm64.whl.metadata (11 kB)
Collecting miniful (from fst-ps<=1.8.1-pyfume->FuzzyTM<=0.4.0->gensim==4.3.0)
 Downloading miniful-0.0.6.tar.gz (2.8 kB)
 Preparing metadata (setup.py) ... done
Requirement already satisfied: six<=1.5 in /Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/lib/python3.9/site-packages (from python-dateutil<=2.8.2->pandas->FuzzyTM<=0.4.0->gensim==4.3.0) (1.15.0)
Requirement already satisfied: sio<=1.5 in /Library/Developer/CommandLineTools/Library/Frameworks/Python3.framework/Versions/3.9/lib/python3.9/site-packages (from python-dateutil<=2.8.2->pandas->FuzzyTM<=0.4.0->gensim==4.3.0) (1.15.0)
 Downloading FuzzyTM-2.0.9-py3-none-any.whl (31 kB)
 Downloading pyFUME-0.3.4-py3-none-any.whl (60 kB)
 Downloading numpy-1.24.4-cp39-cp39-macosx_11_0_arm64.whl (13.9 MB)
 13.9/13.9 MB 3.0 MB/s eta 0:00:0000:0100:01
 Downloading scipy-1.10.1-cp39-cp39-macosx_12_0_arm64.whl (28.9 MB)
 28.9/28.9 MB 6.1 MB/s eta 0:00:0000:0100:01
 Downloading pandas-1.5.3-cp39-cp39-macosx_11_0_arm64.whl (11.0 MB)
 11.0/11.0 MB 6.7 MB/s eta 0:00:00 0:00:01
 Downloading simpful-2.12.0-py3-none-any.whl (24 kB)
Building wheels for collected packages: gensim, fst-ps, miniful
 Building wheel for gensim (setup.py) ... done
 Created wheel for gensim: filename=gensim-4.3.0-cp39-cp39-macosx_10_9_universal2.whl size=24457501 sha256=bd82df7d5d2124a4856c85b33fe883236359b773ef620d627b39b6d70071d13c
 Stored in directory: /Users/mmadhusudan/Library/Caches/pip/wheels/f4/68/4d/7be8c2e7a9e0bd4d8892e33aea529c5c77a1f94a362290191
 Building wheel for fst-ps (setup.py) ... done
 Created wheel for fst-ps: filename=fst-psy-1.8.1-py3-none-any.whl size=20478 sha256=10c887fe00944bd599ac995626e85d57135d5d9da3404236c859de4e213d078
 Stored in directory: /Users/mmadhusudan/Library/Caches/pip/wheels/99/66/48/d7ce0c6927f6abf167bbcd0e537affc7b92c03632f78028411
 Building wheel for miniful (setup.py) ... done
 Created wheel for miniful: filename=miniful-0.0.6-py3-none-any.whl size=3554 sha256=5cb485ba076b077a00aa957e03bc5bc242b2dc31a296ac1f513f4765789318e
 Stored in directory: /Users/mmadhusudan/Library/Caches/pip/wheels/d9/c7/71/db1d4646d963b34c530667501d3d6f34c0825eafae2f0f2cb
Successfully built gensim fst-ps miniful
Installing collected packages: numpy, scipy, pandas, simpful, miniful, fst-ps, pyfume, FuzzyTM, gensim
 Attempting uninstall: numpy
 Found existing installation: numpy 1.26.4
 Uninstalling numpy-1.26.4:
 Successfully uninstalled numpy-1.26.4
 Attempting uninstall: scipy
 Found existing installation: scipy 1.13.1
 Uninstalling scipy-1.13.1:
 Successfully uninstalled scipy-1.13.1
 Attempting uninstall: pandas
 Found existing installation: pandas 2.2.3
 Uninstalling pandas-2.2.3:
 Successfully uninstalled pandas-2.2.3
 Attempting uninstall: gensim
 Found existing installation: gensim 4.3.3
 Uninstalling gensim-4.3.3:
 Successfully uninstalled gensim-4.3.3
Successfully installed FuzzyTM-2.0.9 fst-ps-1.8.1 gensim-4.3.0 miniful-0.0.6 numpy-1.24.4 pandas-1.5.3 pyfume-0.3.4 scipy-1.10.1 simpful-2.12.0

In [71]:

```
!python3 -m spacy download en_core_web_md
```

9442.04s - pydevd: Sending message related to process being replaced timed-out after 5 seconds
/Users/mmadhusudan/Library/Python/3.9/lib/python/site-packages/urllib3/_init_.py:35: NotOpenSSLWarning: urllib3 v2 only supports OpenSSL 1.1.1+, currently the 'ssl' module is compiled with 'LibreSSL 2.8.3'. See: https://github.com/urllib3/urllib3/issues/3020
 warnings.warnl
Defaulting to user installation because normal site-packages is not writeable
Collecting en-core-web-md==3.8.0
 Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_md-3.8.0/en_core_web_md-3.8.0-py3-none-any.whl (33.5 MB)
 33.5/33.5 MB 8.8 MB/s eta 0:00:0000:0100:01
Installing collected packages: en-core-web-md
Successfully installed en-core-web-md-3.8.0
 Download and installation successful
You can now load the package via spacy.load('en_core_web_md')

In [72]:

```
import spacy
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
import numpy as np
import nltk
from collections import Counter

nltk.download('punkt')

nlp = spacy.load("en_core_web_md")

doc = nlp(cleaned_text)

tokens = [token.text.lower() for token in doc if token.is_alpha and not token.is_stop and token.has_vector]

freq = Counter(tokens)

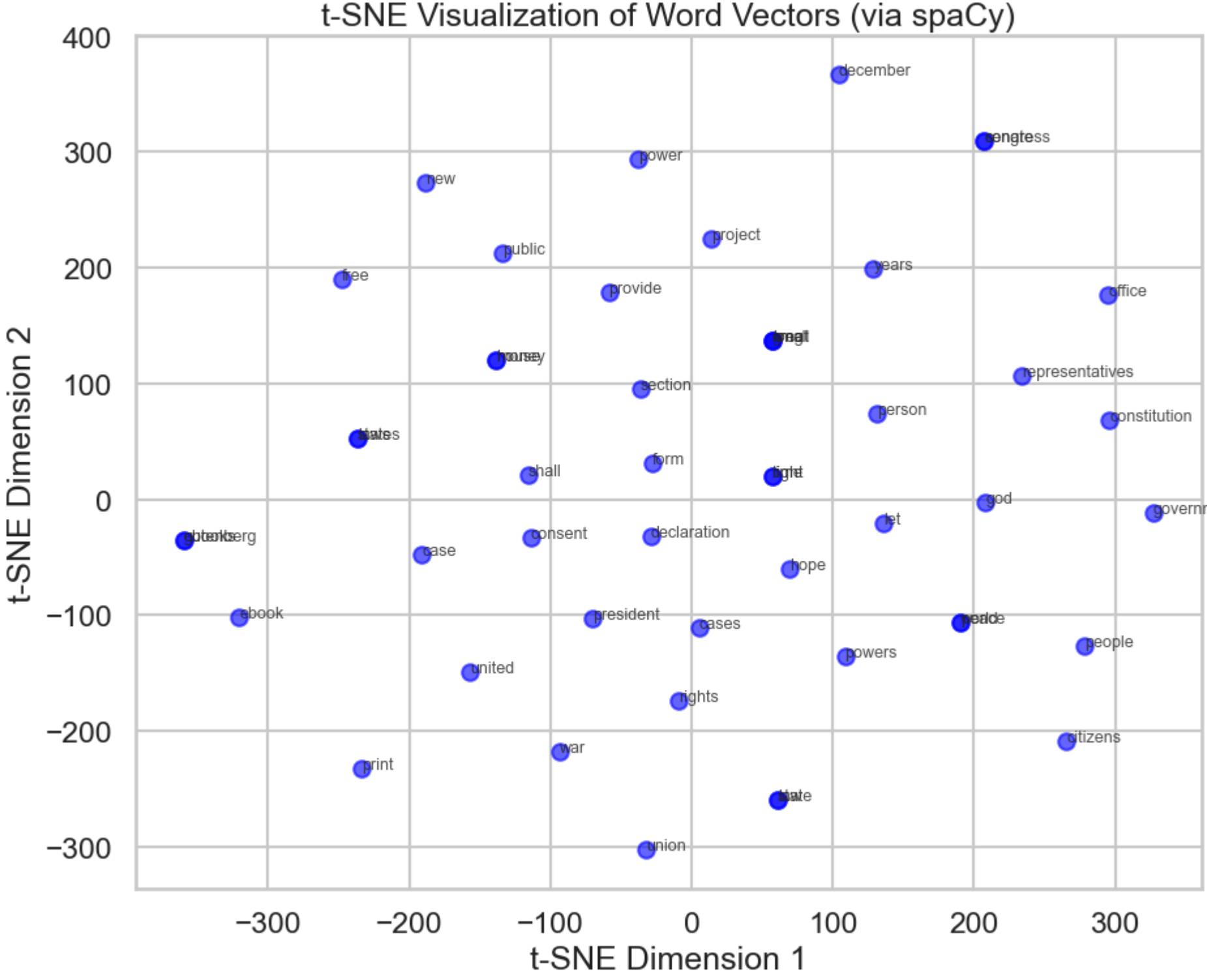
most_common_tokens = [word for word, count in freq.most_common(50)]
print(f"Most common tokens:", most_common_tokens)

word_vectors = np.array([nlp.vocab[word].vector for word in most_common_tokens])

tsne = TSNE(n_components=2, random_state=42)
tsne_results = tsne.fit_transform(word_vectors)

plt.figure(figsize=(10, 8))
plt.scatter(tsne_results[:, 0], tsne_results[:, 1], color='blue', alpha=0.6)
for i, word in enumerate(most_common_tokens):
    plt.annotate(word, (tsne_results[i, 0], tsne_results[i, 1]), fontsize=9, alpha=0.8)
plt.title("t-SNE Visualization of Word Vectors (via spaCy)")
plt.xlabel("t-SNE Dimension 1")
plt.ylabel("t-SNE Dimension 2")
plt.tight_layout()
plt.show()
```

[nltk_data] Downloading package punkt to
[nltk_data] /Users/mmadhusudan/nltk_data...
[nltk_data] Package punkt is already up-to-date!
Most common tokens: ['shall', 'states', 'project', 'united', 'gutenberg', 'state', 'people', 'time', 'law', 'constitution', 'laws', 'congress', 'government', 'president', 'right', 'new', 'war', 'public', 'house', 'union', 'free', 'power', 'ebook', 'person', 'let', 'section', 'years', 'cases', 'world', 'form', 'case', 'god', 'office', 'peace', 'hope', 'small', 'money', 'right
s', 'december', 'great', 'citizens', 'ebooks', 'print', 'powers', 'consent', 'representatives', 'senate', 'long', 'provide', 'declaration']



In [73]:

```
import os
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

def get_cleaned_text(file_path):
    with open(file_path, "r", encoding="utf-8") as f:
        raw_text = f.read()
        # Use common markers to strip Gutenberg header and footer.
        start_marker = "==== START OF THIS PROJECT GUTENBERG EBOOK"
        end_marker = "==== END OF THIS PROJECT GUTENBERG EBOOK"
        start_idx = raw_text.find(start_marker)
        if start_idx != -1:
            text = raw_text[start_idx + len(start_marker):]
        else:
            text = raw_text
        end_idx = text.find(end_marker)
        if end_idx != -1:
            text = text[:end_idx]
        return text.strip()

folder = "Gutenberg Books"
all_files = [f for f in os.listdir(folder) if f.endswith(".txt.txt")]

corpus = []
doc_names = []
for filename in all_files:
    file_path = os.path.join(folder, filename)
    try:
        text = get_cleaned_text(file_path)
        corpus.append(text)
        doc_names.append(filename)
    except Exception as e:
        print(f"Error processing {filename}: {e}")

print(f"Collected {len(corpus)} documents.")

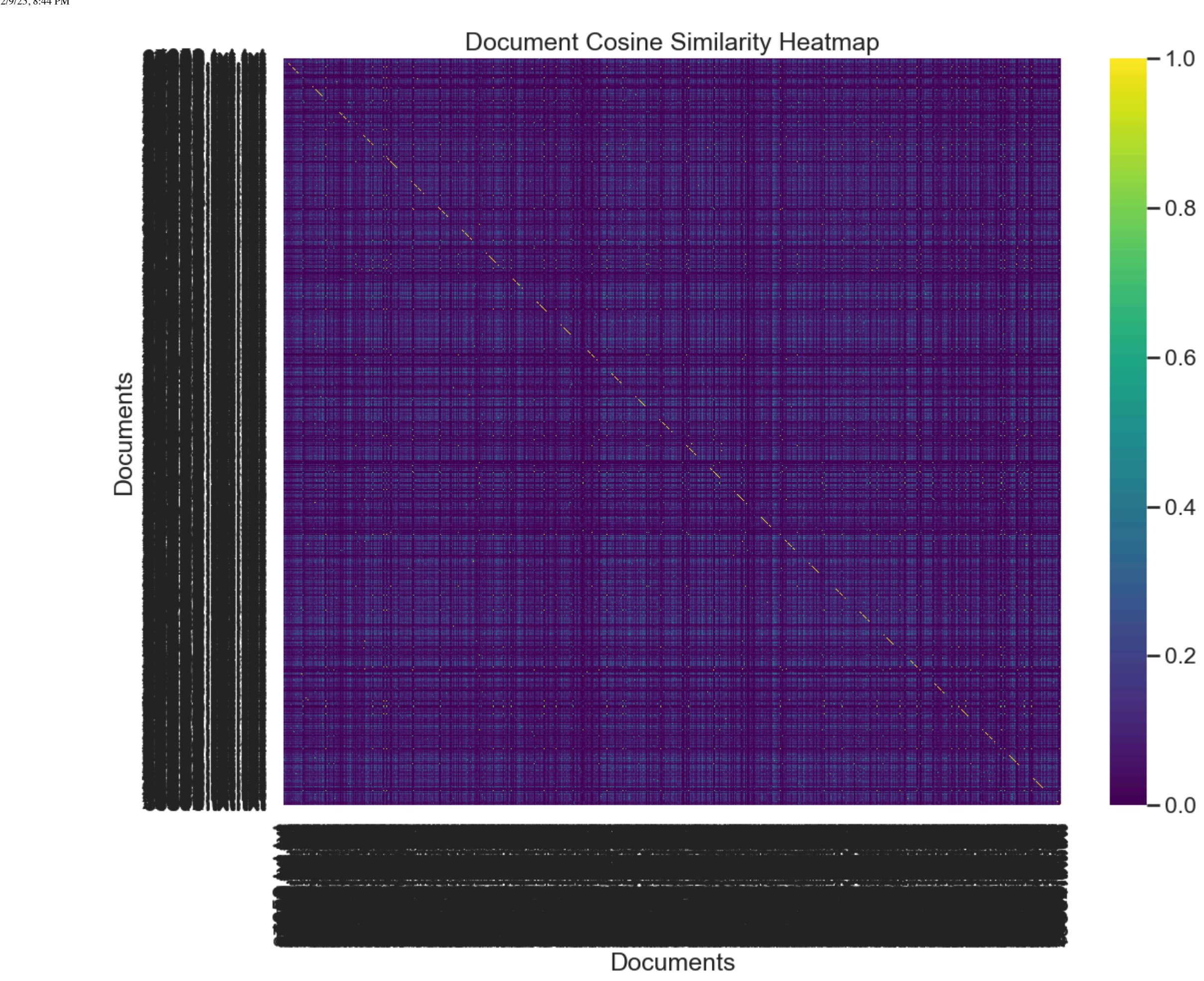
vectorizer = TfidfVectorizer(stop_words='english', max_df=0.95, min_df=2)
tfidf_matrix = vectorizer.fit_transform(corpus)
print(f"TF-IDF matrix shape: {tfidf_matrix.shape}")

cos_sim_matrix = cosine_similarity(tfidf_matrix)

df_sim = pd.DataFrame(cos_sim_matrix, index=doc_names, columns=doc_names)

# --- Visualization: Document Similarity Heatmap ---
plt.figure(figsize=(12, 10))
sns.heatmap(df_sim, cmap='viridis', xticklabels=True, yticklabels=True)
plt.title("Document Cosine Similarity Heatmap")
plt.xlabel("Documents")
plt.ylabel("Documents")
plt.tight_layout()
plt.show()
```

Collected 2475 documents.
TF-IDF matrix shape: (2475, 526829)



```
In [74]: import nltk
from nltk.collocations import BigramCollocationFinder, BigramAssocMeasures
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

all_tokens = []
for doc in corpus:
    # Tokenize each document into words, convert to lowercase, and filter out non-alphabetic tokens.
    tokens = nltk.word_tokenize(doc.lower())
    tokens = [token for token in tokens if token.isalpha()]
    all_tokens.extend(tokens)

print(f"Total tokens aggregated from the corpus: {len(all_tokens)}")

bigram_measures = BigramAssocMeasures()
finder = BigramCollocationFinder.from_words(all_tokens)

finder.apply_freq_filter(3)

scored_bigrams = finder.score_ngrams(bigram_measures.pmi)
# Convert to DataFrame for easy handling
df_bigrams = pd.DataFrame(scored_bigrams, columns=["bigram", "PMI"])
df_bigrams.sort_values(by="PMI", ascending=False).head(20)
print("Top 20 bigrams by PMI:")
display(df_bigrams)

df_bigrams["bigram_str"] = df_bigrams["bigram"].apply(lambda x: " ".join(x))

plt.figure(figsize=(10, 6))
sns.barplot(data=df_bigrams, x="PMI", y="bigram_str", palette="Blues_d")
plt.xlabel("PMI Score")
plt.ylabel("Bigram")
plt.title("Top 20 Bigrams by PMI in the Corpus")
plt.tight_layout()
plt.show()
```

Total tokens aggregated from the corpus: 209085770
Top 20 bigrams by PMI:

	bigram	PMI
0	(abgeruehrter, kugelhopf)	26.054557
104	(khandu, wangchuk)	26.054557
132	(napao, wetikoo)	26.054557
133	(ndeh, ntumazah)	26.054557
134	(nerbia, espartafiarlo)	26.054557
135	(nikica, valentic)	26.054557
136	(nurzhan, subkhanberdin)	26.054557
137	(nuzas, rocabertis)	26.054557
138	(ochthodromus, wilsonius)	26.054557
139	(odjo, tankpinon)	26.054557
140	(ojasta, allikkoon)	26.054557
141	(olaudah, equiano)	26.054557
142	(orhan, ucok)	26.054557
143	(otinielu, tausi)	26.054557
144	(oudom, khattiya)	26.054557
145	(palafoxes, nuzas)	26.054557
146	(papeis, avulsos)	26.054557
147	(pastissons, giraumous)	26.054557
148	(paucás, horas)	26.054557
149	(paucis, annis)	26.054557

```
/var/folders/7j/rv3w77nj6kb6kw_ss1tcqpk0000gp/T/ipykernel_22400/2160407253.py:37: FutureWarning:
Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'y' variable to 'hue' and set 'legend=False' for the same effect.
sns.barplot(data=df_bigrams, x="PMI", y="bigram_str", palette="Blues_d")
```

