```
In [1]:  1  import sys
         2  print(sys.path)
```

['C:\\Users\\harsh\\OneDrive\\Desktop\\UNT Subjects\\NLP', 'C:\\Users\\harsh\\anaconda3\\envs\\tf-gpu\\python38.zi
p', 'C:\\Users\\harsh\\anaconda3\\envs\\tf-gpu\\DLLs', 'C:\\Users\\harsh\\anaconda3\\envs\\tf-gpu\\lib', 'C:\\Users
\\harsh\\anaconda3\\envs\\tf-gpu', '', 'C:\\Users\\harsh\\AppData\\Roaming\\Python\\Python38\\site-packages', 'C:\\U
sers\\harsh\\anaconda3\\envs\\tf-gpu\\lib\\site-packages', 'C:\\Users\\harsh\\anaconda3\\envs\\tf-gpu\\lib\\site-pac
kages\\win32', 'C:\\Users\\harsh\\anaconda3\\envs\\tf-gpu\\lib\\site-packages\\win32\\lib', 'C:\\Users\\harsh\\anaco
nda3\\envs\\tf-gpu\\lib\\site-packages\\Pythonwin']

Importing the required libraries

```
In [2]:  1  import nltk
         2  from nltk.tokenize import word_tokenize
         3  from nltk.util import ngrams
         4  from nltk.probability import ConditionalFreqDist
         5  import docx2txt
         6  import docx
         7  import random
         8  import re
         9  import matplotlib.pyplot as plt
```

```
In [3]:    1
           2  # Open the DOCX file
           3  doc = docx.Document('5 papers related to blockchain.docx')
           4  text = '\n'.join([paragraph.text for paragraph in doc.paragraphs])
           5
           6  # Remove unnecessary punctuations,spaces brackets
           7  text = re.sub(r'\[\d+\]', '', text)
           8  # Remove underscores and hyphens from the sides of words
           9  text = re.sub(r'[-_]', '', text)
          10  # Remove numbers used for points
          11  text = re.sub(r'\d+\.', '', text)
          12  # remove "
          13  text = re.sub(r'"', '', text)
          14  # remove '
          15  text = re.sub(r"'", '', text)
          16  # remove special characters
          17  text = re.sub(r'[+-.,!@#$%^&<>?/\{}()*_=:;|]', '', text)
          18  # remove multiple spaces
          19  text = re.sub(r'\n\s*\n', '\n', text)
          20  # remove numbers
          21  text = re.sub(r'\s\d+\s', ' ', text)
          22  # remove new line, tabs
          23  text = re.sub(r'\n|\t', ' ', text)
          24  # Remove bullets
          25  text = re.sub(r'^[\s\u2022\u2023\u25E6\u2043]*', '', text, flags=re.MULTILINE)
          26  # remove multiple spaces
          27  text = re.sub(r' +', ' ', text)
          28  text = text.lower()
          29
          30  # Write the cleaned text to a TXT file
          31  with open('input.txt', 'w',encoding="utf-8") as file:
          32      file.write(text)
          33
```

File reading starts from here

```
In [4]:   1  filename = "input.txt"
          2  with open(filename,encoding="utf-8") as f:
          3      content = f.readlines()
          4  content = [ line for line in content if line != '\n' ]
          5  content = [ line.lower() for line in content ]
```

```
In [5]:   1  tokens = [t for l in content for t in l.split() ]
          2  # tokens
```

```
In [6]:   1  tl = len(tokens)
          2  tl
```

Out[6]: 9703

Calculating Diversity score of entire data set
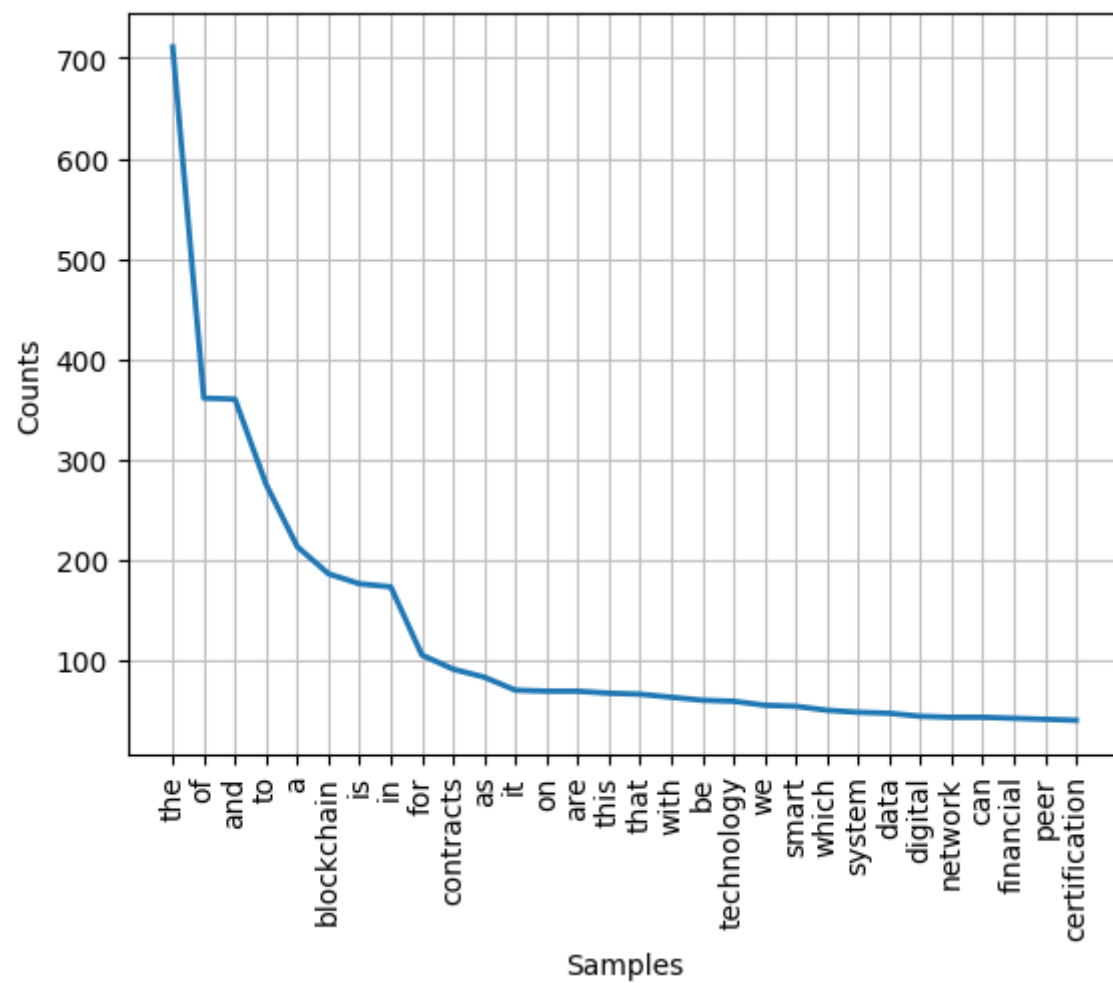
```
In [7]:   1  # Calculate  unique words
          2  unique_words = set(tokens)
          3  num_unique_words = len(unique_words)
          4
          5  diversity_score = num_unique_words / len(tokens)
          6  print(diversity_score)
```

0.20807997526538183

```
In [8]:   1  freq = nltk.FreqDist(tokens)
          2  probs = {k: v/tl for (k,v) in freq.items()}
          3  pvals = list(probs.values())
          4  cumprobs = {k: sum(pvals[0:ix+1]) for ix, (k,v) in enumerate(probs.items())}
```

Frequency plot for the generated tokens
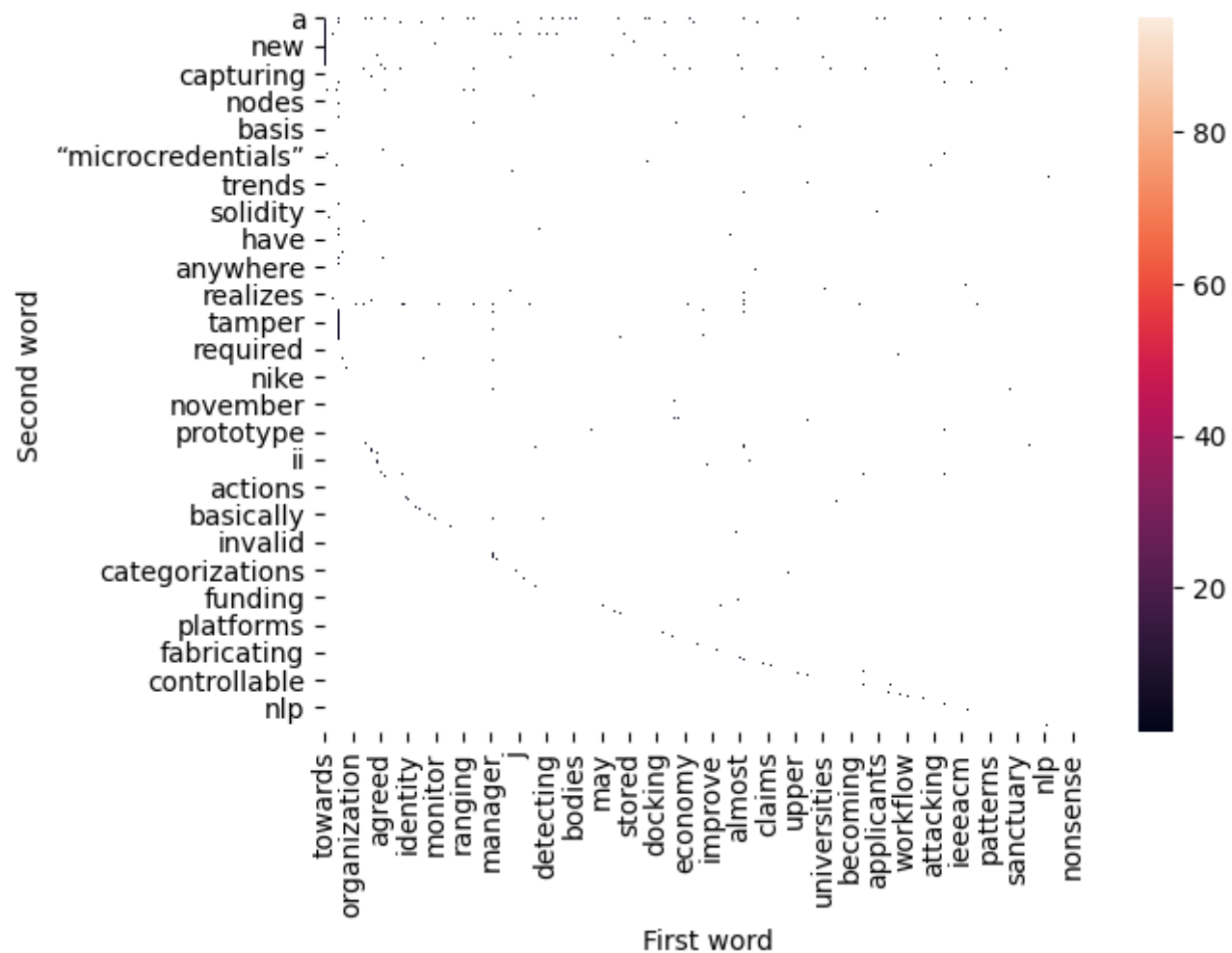
```
1 freq.plot(30, cumulative=False)
2 plt.show()
```

```
1 cfd_bigram = nltk.ConditionalFreqDist(nltk.bigrams(tokens))
```

Heat map for bigrams with conditional Frequency distribution

```python
import seaborn as sns
import pandas as pd

df = pd.DataFrame(cfd_bigram)

# Create the heatmap
sns.heatmap(df)
plt.xlabel("First word")
plt.ylabel("Second word")
plt.show()
```

```
1  cfd_bigram.items()
```

1, 'supervision': 1, 'future': 1, 'main': 1, ...})), ('technical', FreqDist({'principles': 2, 'requirements': 2, 'limitations': 1, 'side': 1, 'characteristics': 1, 'defects': 1, 'means': 1, 'issues': 1, 'problems': 1})), ('limitations', FreqDist({'however': 1})), ('however', FreqDist({'the': 6, 'recognition': 1, 'as': 1, 'peer': 1, 'in': 1, 'with': 1, 'emerging': 1, 'by': 1, 'we': 1})), ('as', FreqDist({'a': 17, 'the': 9, 'well': 6, 'peer': 3, 'an': 2, 'blockchain': 2, 'follows': 2, 'chain': 2, 'security': 1, 'healthcare': 1, ...})), ('or', FreqDist({'a': 1, 'additions': 1, 'consortium': 1, 'individual': 1, 'the': 1, 'quality': 1, 'minimizing': 1, 'who': 1, 'authority': 1, 'bitcoin': 1, ...})), ('software', FreqDist({'component': 2, 'engineering': 2})), ('component', FreqDist({'requires': 1, 'in': 1, 'of': 1})), ('requires', FreqDist({'a': 1, 'more': 1, 'regulatory': 1})), ('comprehensive', FreqDist({'understanding': 1, 'evaluation': 1})), ('characterization', FreqDist({'of': 1})), ('principles', FreqDist({'and': 2, 'of': 1})), ('characteristics', FreqDist({'the': 3, 'of': 3, 'it': 1, 'such': 1, 'athe': 1, 'and': 1})), ('latter', FreqDist({'introduces': 2, 'is': 1})), ('introduces', FreqDist({'an': 2, 'several': 1, 'the': 1})), ('an', FreqDist({'important': 4, 'open': 3, 'architecture': 2, 'attempt': 2, 'uncertainty': 1, 'organization': 1, 'identical': 1, 'ongoing': 1, 'urgent': 1, 'isolated': 1, ...})), ('uncertainty', FreqDist({'for': 1})), ('organization', FreqDist({'to': 1, 'this': 1})), ('decide', FreqDist({'which': 1})), ('which', FreqDist({'is': 13, 'blockchain': 4, 'are': 4, 'we': 3, 'includes': 2, 'will': 2, 'might': 2, 'aims': 1, 'type': 1, 'indexes': 1, ...})), ('protocol', FreqDist({'to': 3, 'best': 1, 'is': 1, 'handles': 1, 'under': 1, 'in': 1, 'architecture': 1, 'can': 1})), ('best', FreqDist({'meets': 1})), ('meets', FreqDist({'its': 1})), ('needs', FreqDist({'to': 4, 'and': 1, 'safety': 1, 'the': 1, 'constant': 1, 'of': 1})), ('demands', FreqDist({'in': 1})), ('in', FreqDist({'the': 53, 'a': 13, 'this': 7, 'digital': 7, 'order': 6, 'terms': 5, 'particular': 5, 'addition': 4, 'many': 3, 'its': 3, ...})), ('general', FreqDist({'there': 1, 'inconsistency': 1, 'but': 1})), ('there', FreqDist({'are': 12, 'is': 6, 'was':

Calculating perplexity score of bigrams

```
In [13]:   1  import math
           2
           3  log_prob_sum = 0
           4  num_tokens = len(tokens)
           5  for i in range(num_tokens - 1):
           6      w0 = tokens[i]
           7      w1 = tokens[i+1]
           8      bigram_count = cfd_bigram[w0][w1]
           9      bigram_prob = (bigram_count + 1) / (cfd_bigram[w0].N() + len(cfd_bigram[w0]))
          10      log_prob_sum += math.log2(bigram_prob)
          11
          12  perplexity = math.pow(2, -log_prob_sum / num_tokens)
          13  print(perplexity)
          14
```

10.965655218617805

Calculating diversity score of bigrams

```
In [14]:   1  num_unique_bigrams = len(set(nltk.bigrams(tokens)))
           2  diversity = num_unique_bigrams / len(list(nltk.bigrams(tokens)))
           3  print(diversity)
           4
```

0.7062461348175634

```
In [15]:   1  def generate_bigram_sentence(data):
           2      # randomly select a starting word
           3      word = random.choice(list(data.keys()))
           4      sentence = [word]
           5
           6      # generate the rest of the sentence using bigrams
           7      for i in range(9):
           8          if word not in data:
           9              break
          10          freq_dist = data[word]
          11          next_word = freq_dist.most_common(1)[0][0]
          12          sentence.append(next_word)
          13          word = next_word
          14
          15      # join the sentence and return
          16      return " ".join(sentence)
          17
          18  for i in range(5):
          19      sentence = generate_bigram_sentence(cfd_bigram)
          20      print("Sentence : ",i+1)
          21      print(sentence)
          22
```

```
Sentence :  1
architecture for the blockchain technology and the blockchain technology and
Sentence :  2
reads data and the blockchain technology and the blockchain technology
Sentence :  3
modules is a blockchain technology and the blockchain technology and
Sentence :  4
though the blockchain technology and the blockchain technology and the
Sentence :  5
faults in the blockchain technology and the blockchain technology and
```

```
In [16]:   1  trigrams = nltk.trigrams(tokens)
```

```
In [17]:   1  condition_pairs = (((w0, w1), w2) for w0, w1, w2 in trigrams)
           2  cfd_trigram = nltk.ConditionalFreqDist(condition_pairs)
```

Calculating perplexity score of trigrams

```
In [18]:   1  trigram_prob = 0.0
           2  for w0_w1, freq_w2 in cfd_trigram.items():
           3      if '' in w0_w1:
           4          continue
           5      for w2 in freq_w2:
           6          trigram_count = cfd_trigram[w0_w1][w2]
           7          bigram_count = cfd_bigram[w0_w1[0]][w0_w1[1]]
           8          trigram_prob *= (trigram_count + 1) / (bigram_count + len(freq))
           9
          10
          11  if trigram_prob == 0.0:
          12      perplexity = float('inf')
          13  else:
          14      perplexity = math.pow(2, -1 * (math.log2(trigram_prob)))
          15
          16  print(perplexity)
          17
```
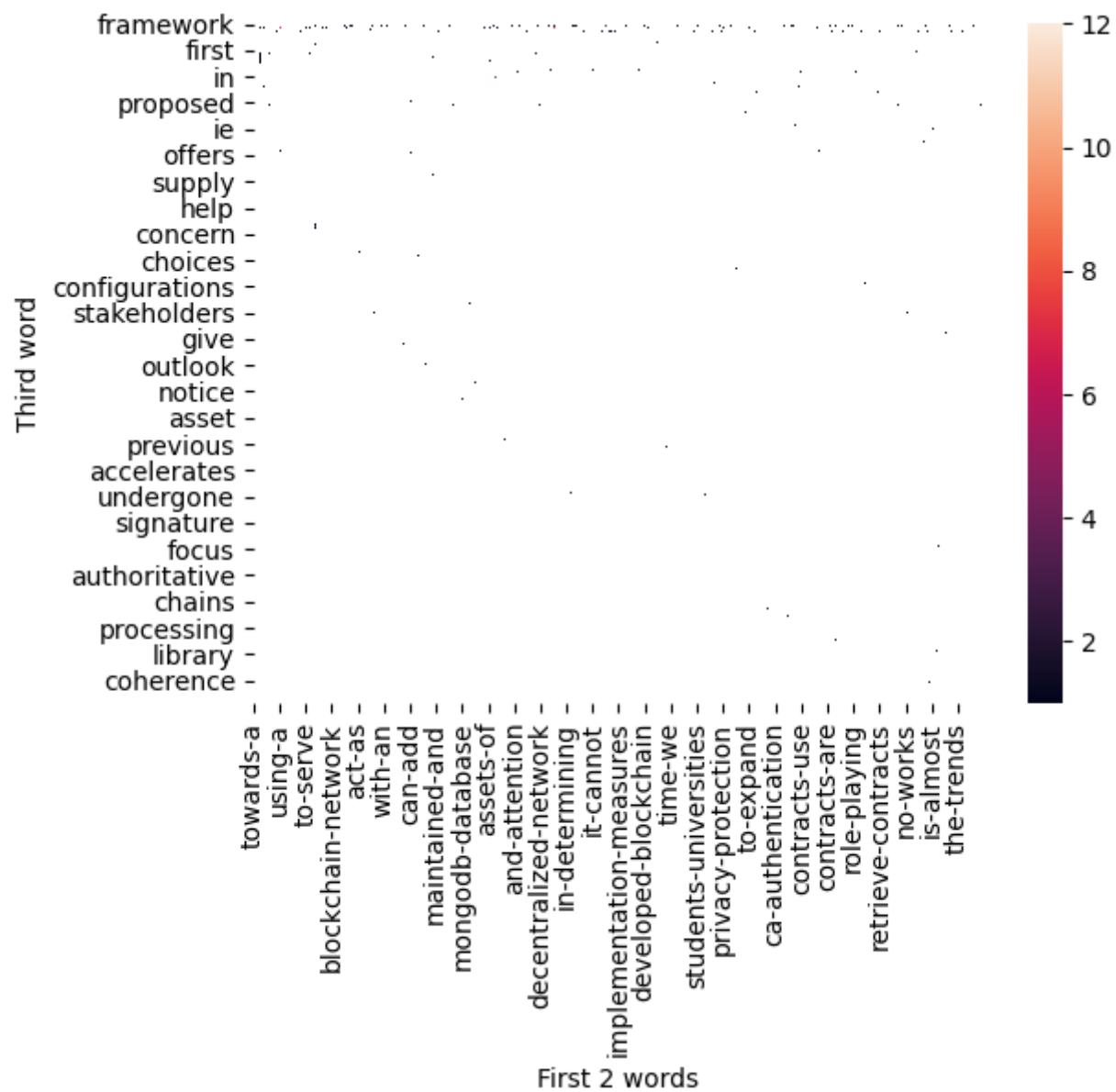
inf

Calculating diversity score of trigrams

```
In [19]:   1  trigrams = list(nltk.trigrams(tokens))
           2  diversity_score = len(set(trigrams)) / len(trigrams)
           3  print(diversity_score)
           4
```

0.920832903824348

```
In [20]:  1  df = pd.DataFrame(cfd_trigram)
          2  sns.heatmap(df)
          3  plt.xlabel("First 2 words")
          4  plt.ylabel("Third word")
          5  plt.show()
```

```
In [21]:    1  cfd_trigram.items()
```

e': 1})), (('evaluating', 'the'), FreqDist({'behaviour': 1, 'topic': 1})), (('the', 'behaviour'), FreqDist({'of':
1})), (('behaviour', 'of'), FreqDist({'blockchain': 1})), (('protocols', 'under'), FreqDist({'different': 1})),
(('under', 'different'), FreqDist({'test': 1})), (('different', 'test'), FreqDist({'scenarios': 1})), (('test', 's
cenarios'), FreqDist({'a': 1})), (('scenarios', 'a'), FreqDist({'distributed': 1})), (('a', 'distributed'), FreqDi
st({'ledger': 3, 'network': 1, 'ca': 1})), (('ledger', 'is'), FreqDist({'often': 1})), (('is', 'often'), FreqDist
({'described': 1})), (('often', 'described'), FreqDist({'as': 1})), (('described', 'as'), FreqDist({'a': 1})),
(('a', 'shared'), FreqDist({'distributed': 1, 'governance': 1})), (('shared', 'distributed'), FreqDist({'databas
e': 1})), (('distributed', 'database'), FreqDist({'which': 1, 'the': 1})), (('database', 'which'), FreqDist({'is':
1})), (('which', 'is'), FreqDist({'a': 4, 'accessed': 1, 'considered': 1, 'ideal': 1, 'an': 1, 'feasible': 1, 'joi
ntly': 1, 'impressive': 1, 'the': 1, 'also': 1})), (('is', 'accessed'), FreqDist({'and': 1})), (('accessed', 'an
d'), FreqDist({'maintained': 1})), (('and', 'maintained'), FreqDist({'by': 1})), (('maintained', 'by'), FreqDist
({'the': 2, 'a': 1})), (('by', 'a'), FreqDist({'set': 1})), (('of', 'independent'), FreqDist({'possibly': 1})),
(('independent', 'possibly'), FreqDist({'untrusted': 1})), (('possibly', 'untrusted'), FreqDist({'participants':
1})), (('untrusted', 'participants'), FreqDist({'ie': 1})), (('participants', 'ie'), FreqDist({'nodes': 1})), (('i
e', 'nodes'), FreqDist({'each': 1})), (('nodes', 'each'), FreqDist({'participant': 1, 'network': 1})), (('each',
'participant'), FreqDist({'owns': 1})), (('participant', 'owns'), FreqDist({'an': 1})), (('owns', 'an'), FreqDist
({'identical': 1})), (('an', 'identical'), FreqDist({'copy': 1})), (('identical', 'copy'), FreqDist({'of': 1})),
(('copy', 'of'), FreqDist({'the': 1})), (('of', 'the'), FreqDist({'blockchain': 12, 'application': 5, 'nodes': 4,
'financial': 4, 'network': 3, 'most': 3, 'system': 3, 'database': 2, 'state': 2, 'deployment': 2, ...})), (('the',
'database'), FreqDist({'of': 1, 'in': 1, 'and': 1, 'is': 1, 'at': 1})), (('database', 'of'), FreqDist({'transactio

```python
def generate_trigram_sentence():
    sentence = []
    # Pick the first word at random
    w0 = random.choice(list(cfd_trigram.keys()))[0]
    sentence.extend([w0])

    # Pick the second word from bigram
    if w0 in cfd_trigram:
        w1 = random.choice(list(cfd_trigram[w0].keys()))
    else:
        # If w0 not present, pick random word from corpus
        w1 = random.choice(tokens)
    sentence.extend([w1])

    # Use those two words to pick the most frequent third word
    if w0 in cfd_trigram and w1 in cfd_trigram[w0]:
        w2 = cfd_trigram[w0][w1].most_common(1)[0][0]
    else:
        # If w0 or w1 not present, picking random word from corpus
        w2 = random.choice(tokens)
    sentence.extend([w2])

    # Iterate until 10 words
    for i in range(7):
        w0, w1, w2 = w1, w2, None
        if w0 in cfd_trigram and w1 in cfd_trigram[w0]:
            w2 = cfd_trigram[w0][w1].most_common(1)[0][0]
        if w2 is None:
            w2 = random.choice(tokens)
        sentence.extend([w2])

    return ' '.join(sentence)

for i in range(5):
    print('Sentence : ',i+1)
    print(generate_trigram_sentence())
```

```
Sentence :  1
final developers feasibility institutes the the regulatory digital the at
Sentence :  2
of blockchain generate with not system aspects peer security and
Sentence :  3
brains a development systemic feasible this application controlled a number
Sentence :  4
time questions protocol affect a the industrial database a structure
Sentence :  5
framework popular tremendous incorporating aspects in can to tolerance consensus
```

In [23]:
```python
1  quadgrams = nltk.ngrams(tokens,4)
```

In [24]:
```python
1  condition_pairs = (((w0, w1, w2), w3) for w0, w1, w2, w3 in quadgrams)
2  cfd_quadgram = nltk.ConditionalFreqDist(condition_pairs)
```

Calculating perplexity score of quadgrams

In [25]:
```python
1  quadgram_prob = 1.0
2  for w0, w1, w2, w3 in quadgrams:
3      trigram_count = cfd_quadgram[(w0, w1, w2)][w3]
4      bigram_count = cfd_trigram[(w0, w1)][w2]
5      quadgram_prob *= (trigram_count + 1) / (bigram_count + len(freq))
6
7  perplexity = math.pow(2, -1 * (math.log2(quadgram_prob)))
8  print(perplexity)
```
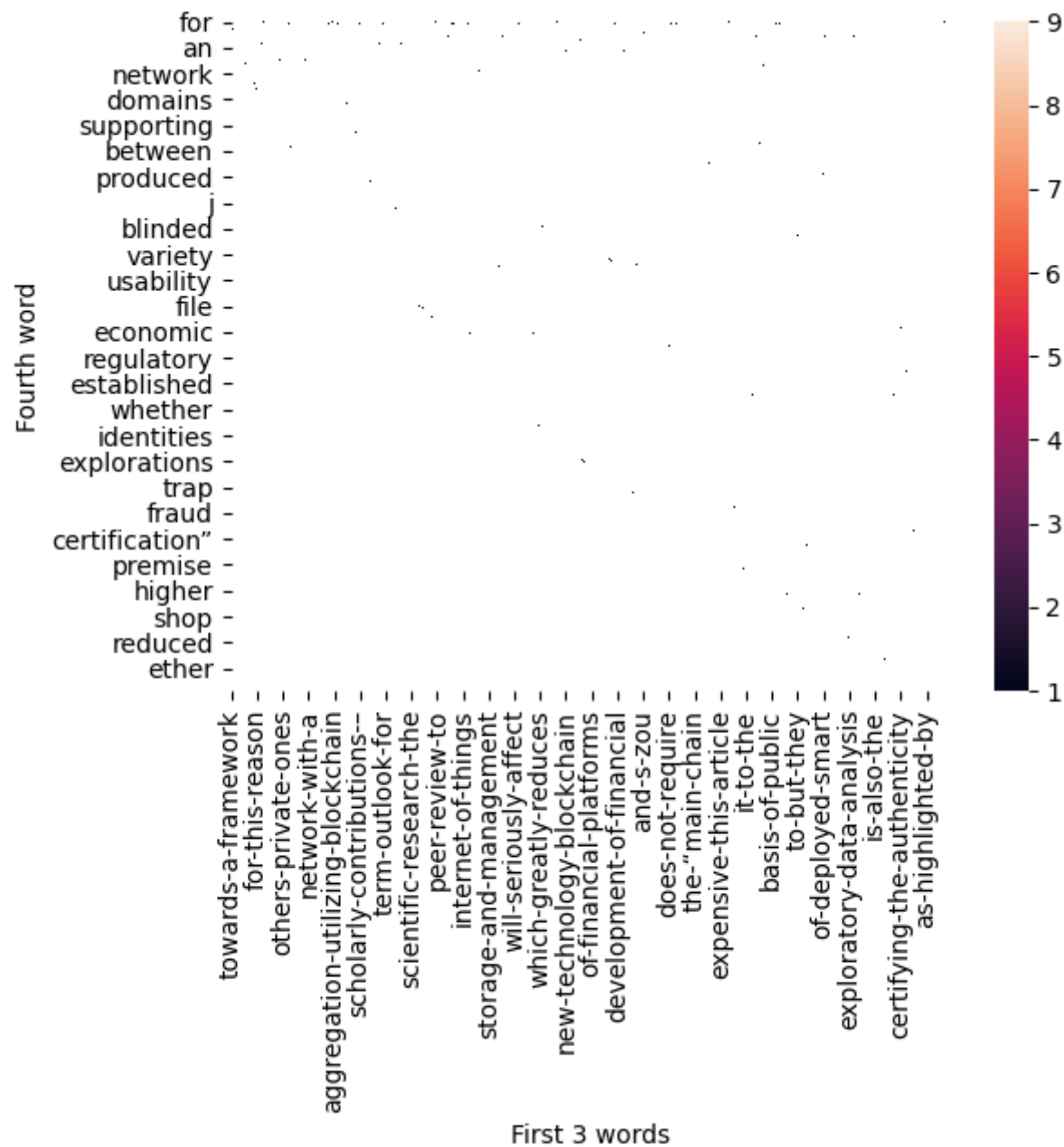
```
1.0
```

Calculating diversity score of quadgrams

```python
diversity_score = 0.0
freq_quadgrams = nltk.FreqDist(quadgrams)
unique_quadgrams = set(quadgrams)
if len(freq_quadgrams) > 0 :
    diversity_score = len(unique_quadgrams) / len(freq_quadgrams)
print(diversity_score)
```

```
0.0
```

```python
df = pd.DataFrame(cfd_quadgram)

# Create the heatmap
sns.heatmap(df)
plt.xlabel("First 3 words")
plt.ylabel("Fourth word")
plt.show()
```

```
In [28]:   1  cfd_quadgram.items()
```

'challenges': 1})), (('ledger', 'technology', 'dlt'), FreqDist({'appears': 1, 'was': 1})), (('technology', 'dlt', 'appears'), FreqDist({'to': 1})), (('dlt', 'appears', 'to'), FreqDist({'be': 1})), (('appears', 'to', 'be'), FreqDist({'at': 1})), (('to', 'be', 'at'), FreqDist({'a': 1})), (('be', 'at', 'a'), FreqDist({'worldwide': 1})), (('at', 'a', 'worldwide'), FreqDist({'threshold': 1})), (('a', 'worldwide', 'threshold'), FreqDist({'of': 1})), (('worldwide', 'threshold', 'of'), FreqDist({'acceptance': 1})), (('threshold', 'of', 'acceptance'), FreqDist({'and': 1})), (('of', 'acceptance', 'and'), FreqDist({'adoption': 1})), (('acceptance', 'and', 'adoption'), FreqDist({'since': 1})), (('and', 'adoption', 'since'), FreqDist({'their': 1})), (('adoption', 'since', 'their'), FreqDist({'inception': 1})), (('since', 'their', 'inception'), FreqDist({'several': 1})), (('their', 'inception', 'several'), FreqDist({'innovative': 1})), (('inception', 'several', 'innovative'), FreqDist({'projects': 1})), (('several', 'innovative', 'projects'), FreqDist({'have': 1})), (('innovative', 'projects', 'have'), FreqDist({'been': 1})), (('projects', 'have', 'been'), FreqDist({'proposing': 1})), (('have', 'been', 'proposing'), FreqDist({'solutions': 1})), (('been', 'proposing', 'solutions'), FreqDist({'to': 1})), (('proposing', 'solutions', 'to'), FreqDist({'the': 1})), (('solutions', 'to', 'the'), FreqDist({'blockchain': 1})), (('to', 'the', 'blockchain'), FreqDist({'trilemma': 1})), (('the', 'blockchain', 'trilemma'), FreqDist({'improving': 1})), (('blockchain', 'trilemma', 'improving'), FreqDist({'blockchain': 1})), (('trilemma', 'improving', 'blockchain'), FreqDist({'features': 1})), (('improving', 'blockchain', 'features'), FreqDist({'and': 1})), (('blockchain', 'features', 'and'), FreqDist({'its': 1})), (('features', 'and', 'its'), FreqDist({'technical': 1})), (('and', 'its', 'technical'), FreqDist({'limitations': 1})), (('its', 'technical', 'limitations'), FreqDist({'however': 1})), (('technical', 'limitations', 'however'), FreqDist({'the': 1})), (('limitations', 'however', 'the'), FreqDist({'adoption': 1})), (('however', 'the', 'adoption'), FreqDist({'of': 1})), (('the', 'adoption', 'of'), FreqDist({'blockchain': 1})), (('adoption', 'of', 'blockcha

```
In [29]:   1  for i in range(5):
           2      # select random word
           3      print('Sentence : ',i+1)
           4      w0, w1, w2 = random.choice(list(cfd_quadgram))
           5
           6      sentence = [w0, w1, w2]
           7
           8      # Using the conditional frequencies to generate the sentence
           9      for j in range(7):
          10          w3_freqdist = cfd_quadgram[(w0, w1, w2)]
          11          w3 = w3_freqdist.max()
          12          sentence.append(w3)
          13          w0, w1, w2 = w1, w2, w3
          14
          15      print(" ".join(sentence))
```

```
Sentence :  1
technology and application development white paper blockchain reference framework and
Sentence :  2
produced by the network as well as a connectivity manager
Sentence :  3
with enough cryptocurrency to pay for the user gas costs
Sentence :  4
will not be conducive to the development trend of blockchain
Sentence :  5
have also demonstrated what can be achieved while pushing the
```

In [ ]:   1