

Kết quả

```
0. hardhat - Why are there so many async and await in smart contract tests (question mark) - Ethereum ...
0. hardhat - Why are there so many async and await in smart contract tests (question mark) - Ethereum Stack Exchange.txt
pipeline > stackoverflow > 0. hardhat - Why are there so many async and await in smart contract tests (question mark) - Ethereum St
1 hardhat - why are there so many async and await in smart contract tests? - Ethereum St
2 ##### Hoàng Sang #####
3 The text is written in the language: en
4 ##### Hoàng Sang #####
5 Phân tử cho văn bản tiếng Việt: (Bởi vì truy vấn trong trang stackoverflow có khá ít r
6 ['Phân', 'tử', 'cho', 'văn bản', 'tiếng', 'Việt', ':', '(', 'Bởi', 'vì', 'truy vấn', '
7 ##### Hoàng Sang #####
8 This is the original text:
9
10 What's the use of async and await and why are there so many of them in smart contract t
11
12
13 ##### Hoàng Sang #####
14 This is the text that has been Unicode Normalization:
15
16
17 b'\nWhat\xe2\x80\x99s the use of async and await and why are there so many of them in
18
19 ##### Hoàng Sang #####
20 Count and sorted stopWords:
21 {'of': 2, 'and': 2, 's': 1, 'the': 1, 'why': 1, 'are': 1, 'there': 1, 'so': 1, 'them':
22 ##### Hoàng Sang #####
23 STEMMING:
24 What: what
25 ', '
26 s: s
27 the: the
28 use: use
29 of: of
30 async: async
31 and: and
32 await: await
33 and: and
34 why: whi
35 are: are
36 there: there
37 so: so
38 many: mani
39 of: of
40 them: them
41 in: in
42 smart: smart
43 contact: contact
44 tests: test
45 ?: ?
46
47 ##### Hoàng Sang #####
48 LEMMATIZATION:
49 What: What
50 ', '
51 s: s
52 the: the
53 use: use
54 of: of
55 async: async
56 and: and
57 await: await
58 and: and
59 why: why
60 are: are
61 there: there
62 so: so
63 many: many
64 of: of
65 them: them
66 in: in
67 smart: smart
68 contact: contact
69 tests: test
70 ?: ?
71
```

Code

```
Pipeline_PhamHoangSang_19IT113.py x 0. hardhat - Why are there so many async and await in smart contract
pipeline > Pipeline_PhamHoangSang_19IT113.py > Crawler > run
1
2 import codecs
3 from urllib.parse import urljoin
4 import logging
5 from bs4 import BeautifulSoup
6 from nltk.tokenize import word_tokenize, sent_tokenize
7 import requests
8 from nltk.corpus import stopwords
9 from nltk.stem import PorterStemmer, WordNetLemmatizer
10 from langdetect import detect
11 import underthesea
12
13 logging.basicConfig(
14     format='%(asctime)s %(levelname)s: %(message)s',
15     level=logging.INFO)
16
17
18 class Crawler:
19
20     def __init__(self):
21         self.visited_urls = []
22         self.urls_to_visit = []
23         self.number_of_article = 0
24         self.filtered_words = []
25         self.stop_words = []
26         self.divider = "\n##### Hoàng Sang #####\n"
27
28     def create_text_file(self, soup):
29         # Get text
30         title = soup.find('h1', {'class': 'fs-headline1'})
31         content = soup.find('div', {'class': 'js-post-body'})
32         # Write text to file with encoding is utf-8
33         file = codecs.open("D:\\E23.1\\NaturalLanguageProcessing\\Ex\\pipeline\\stack",
34                             "a", "utf-8")
35         file.write(str(self.number_of_article) + ". " +
36                     str(soup.title.string)
37                     .replace("|", "").replace("(", "")
38                     .replace("!", "").replace("/", "")
39                     .replace("?", " (question mark)").replace(",", "")) + '\n')
40
41         file.write(str(soup.title.string))
42
43         file.write(self.divider)
44         file.write("The text is written in the language: ")
45         file.write(detect(soup.title.string))
46
47         file.write(self.divider)
48         sentence = "Phân tử cho văn bản tiếng Việt: (Bởi vì truy vấn trong trang stack)"
49         file.write(sentence + "\n")
50         file.write(str(underthesea.word_tokenize(sentence)))
51
52         file.write(self.divider)
53         file.write("This is the original text:\n")
54         file.write(str(content.text) + "\n")
55
56         # split sentences with sent_tokenize
57         file.write(self.divider)
58         file.write("This is the text that has been Unicode Normalization: \n\n")
59         for paragraph in sent_tokenize(content.text):
60             # Unicode Normalization with .encode("utf-8")
61             file.write(str(paragraph.encode("utf-8")) + "\n")
62             self.remove_stopwords(word_tokenize(paragraph))
63
64         file.write(self.divider)
65         file.write("Count and sorted stopWords: \n")
66         file.write(str(self.count_stop_words()))
67
68         file.write(self.divider)
69         file.write("STEMMING: \n")
70         for paragraph in sent_tokenize(content.text):
71             for word in word_tokenize(paragraph):
72                 file.write(word + ": " + str(self.to_stemming(word)) + "\n")
73         file.write(self.divider)
```

```
Pipeline_PhamHoangSang_19IT113.py x 0. hardhat - Why are there so many async and await in smart contract
pipeline > Pipeline_PhamHoangSang_19IT113.py > Crawler > run
74     file.write("LEMMATIZATION: \n")
75     for paragraph in sent_tokenize(content.text):
76         for word in word_tokenize(paragraph):
77             file.write(word + ": " + str(self.to_lemmatizer(word)) + "\n")
78
79     file.close()
80     self.number_of_article += 1
81
82     def to_lemmatizer(self, word):
83         lemmatizer = WordNetLemmatizer()
84         return lemmatizer.lemmatize(word)
85
86     def to_stemming(self, word):
87         stemmer = PorterStemmer()
88         return stemmer.stem(word)
89
90     def remove_stopwords(self, wordTokenize):
91         for word in wordTokenize:
92             if word not in stopwords.words('english'):
93                 self.filtered_words.append(word)
94             else:
95                 self.stop_words.append(word)
96
97     def count_stop_words(self):
98         count = {}
99         for word in self.stop_words:
100             count[word] = self.stop_words.count(word)
101         sorted_count_words = {k: v for k, v in sorted(
102             count.items(), key=lambda item: item[1], reverse=True)}
103         return sorted_count_words
104
105     def get_linked_urls(self, url, soup):
106         for link in soup.find_all('a', {'class': 'question-hyperlink'}):
107             path = link.get('href')
108             if path and path.startswith('/'):
109                 path = urljoin(url, path)
110             yield path
111
112     def add_url_to_visit(self, url):
113         if url not in self.visited_urls and url not in self.urls_to_visit:
114             self.urls_to_visit.append(url)
115
116     def crawl(self, url):
117         html = requests.get(url)
118         soup = BeautifulSoup(html.content, 'html.parser')
119         self.create_text_file(soup)
120         self.stop_words = []
121         self.filtered_words = []
122         for url in self.get_linked_urls(url, soup):
123             self.add_url_to_visit(url)
124         print("Link not visited: ", len(self.urls_to_visit))
125         print("Link visited: ", len(self.visited_urls))
126
127     def run(self):
128         for num_page in range(0, 10):
129             num_page = num_page + 1
130             print("Page: ", num_page)
131             page = requests.get(
132                 "https://stackoverflow.com/questions?tab=newest&pag" + str(num_page))
133             soup = BeautifulSoup(page.content, 'html.parser')
134             for url in self.get_linked_urls("https://stackoverflow.com/questions/", soup):
135                 self.add_url_to_visit(url)
136         while (len(self.urls_to_visit) > 0 & self.number_of_article < 1000):
137             url = self.urls_to_visit.pop(0)
138             logging.info(f'Crawling: {url}')
139             try:
140                 self.crawl(url)
141             except Exception:
142                 logging.exception(f'Failed to crawl: {url}')
143             finally:
144                 self.visited_urls.append(url)
145
146
147 0 2 hrs 50 mins Ln 128, Col 36 Spaces: 4 UTF-8 CRLF Python 3.10.6 64-bit Editor: ready
```