

ПРОГРАММА ДЛЯ ЭВМ

**Утилита для скачивания и анализа сообщений с
форума velomania.ru**

Фрагменты исходного текста программы

Листов 4

Правообладатель-автор: Трофимов Владислав Александрович

© Трофимов Владислав Александрович, 2016

г. Санкт-Петербург
2016 г.

```

import re
import csv

import requests
from bs4 import BeautifulSoup

__author__ = 'Stranger'

thread_url = 'http://forum.velomania.ru/showthread.php?t=74626'
image_label = '**IMAGE**'
video_label = '**VIDEO**'
output_file = 'thread1.tsv'

def extract_data_from_url(url):
    html = requests.get(url).content
    soup = BeautifulSoup(html, 'html.parser')
    posts = soup.findAll('div', {'class': 'postdetails'})

    extracted = []

    for post in posts:
        author = get_author(post)
        content = get_content(post)
        contains_video = video_label in content
        contains_images = image_label in content
        extracted.append((url, author, content, contains_images, contains_video))

    return extracted

def get_content(post):
    raw_content = post.find('div', {'class': 'content'}).prettify()
    pre_content = remove_tags(raw_content)
    pre_content = linearize_and_remove_trash(pre_content)
    pre_content = replace_images_and_video_with_label(pre_content)
    pre_content = linearize(pre_content)
    # pre_content = highlight_links(pre_content)
    content = remove_nonprintable_chars(pre_content)
    return content

def get_author(post):
    author_block = post.find('a', {'class': 'username'})
    author = author_block.find('strong').contents[0]
    if not isinstance(author, str):

```



```
return pages

def print_to_csv(data):
    with open(output_file, 'w', newline='') as csvfile:
        csvwriter = csv.writer(csvfile, delimiter='\t')
        for row in data:
            csvwriter.writerow(list(row))

def main():
    data = []
    pages = get_all_pages_url()

    for url in pages:
        print('Parsing page ' + url + ' of ' + str(len(pages)), end='\r')
        data += extract_data_from_url(url)

    print_to_csv(data)

if __name__ == '__main__':
    main()
```

Всего пронумеровано и прошнуровано
4 листов фрагментов исходного
текста программы

Правообладатель

_____ / Трофимов В.А.
подпись