

**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ - ĐHQG HN**  
**KHOA TRÍ TUỆ NHÂN TẠO**

---



**BÁO CÁO**  
**BÀI TẬP CUỐI KÌ**

**Giảng viên: TS. ĐẶNG TRẦN BÌNH**

**Trợ giảng: CN. Nguyễn Văn Phi**

**Lớp: K67 – AI1**

**Sinh viên: Trần Văn Dy\_22022523**

# Đề bài: Phân tích tương tác và nội dung của một/ nhiều tài khoản Facebook

## Giới thiệu

Trong thực tế hiện nay việc khai thác dữ liệu từ Facebook là một việc vô cùng cần thiết và có ứng dụng rất rộng rãi. Nhiều công ty truyền thông tại Việt Nam đã tích lũy kinh nghiệm trong việc thu thập dữ liệu từ Facebook để đánh giá tác động của chương trình quảng cáo và tiếp thị. Các doanh nghiệp này thường thực hiện việc cào dữ liệu để tổng hợp thông tin và đo lường mức độ ảnh hưởng của họ. Một số công ty khác thì bán giải pháp trích xuất thông tin từ Facebook với đa dạng các nội dung có thể thu thập được. Hoặc đơn giản nhất với sinh viên hoặc một số shop nhỏ họ thực hiện cào dữ liệu để theo dõi tình trạng phát triển của fanpage, xu hướng phát triển trong thời gian gần đây.

Việc cào dữ liệu FaceBook là một ý tưởng rất hay và cần thiết để nghiên cứu . Chính vì vậy trong bài tập này em đã thực hiện cào dữ liệu Facebook về để nghiên cứu và phân tích dựa vào các trường thông tin đã thu thập được. Do chưa có nhiều kinh nghiệm về làm báo cáo cũng như những hạn chế về mặt kiến thức , trong bài báo cáo chắc chắn cũng không thể tránh khỏi những thiếu sót . Em mong nhận được sự đóng góp , phê bình từ phía thầy để báo cáo của em được hoàn thiện hơn .

Em xin chân thành cảm ơn thầy!

## Nội dung

### Phần 1: Thu thập dữ liệu .....

#### I. Yêu cầu cài đặt.....

1. Thư viện dùng .....
2. File cookies.txt.....

#### II. Cào dữ liệu

1. Thực hành .....
2. Tổ chức lại dữ liệu thô .....
3. Lưu file .....

### Phần 2: Tiền xử lí dữ liệu

#### I. Bài viết.....

#### II. Bình luận.....

### Phần 3: Phân tích

#### I. Bài viết

1. Bài đăng theo các thứ .....
2. Các mốc thời gian trong ngày đăng bài .....
3. Phân bố các lượt tương tác theo giờ .....
4. Reaction .....
5. Mối quan hệ các loại tương tác và độ dài bài viết.....
6. Hình ảnh và video .....
7. Lượng tương tác trung bình.....
8. Nội dung bài viết .....

#### II. Bình luận

1. Phân tích comment .....
2. Phân tích commenter .....

## Phần 1: Thu thập dữ liệu

### I, Yêu cầu về cài đặt

#### 1. Thư viện dùng

Cài đặt thư viện facebook\_scraper để cào thông tin các bài viết

```
%pip install facebook_scraper pandas numpy

Collecting facebook_scraper
  Downloading facebook_scraper-0.2.59-py3-none-any.whl (45 kB)
    0.0/45.5 kB ? eta -:-:--
    45.5/45.5 kB 1.4 MB/s eta 0:00:00
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (1.5.3)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (1.23.5)
Collecting dateparser<2.0.0,>=1.0.0 (from facebook_scraper)
  Downloading dateparser-1.2.0-py2.py3-none-any.whl (294 kB)
    295.0/295.0 kB 7.7 MB/s eta 0:00:00
Collecting demjson3<4.0.0,>=3.0.5 (from facebook_scraper)
  Downloading demjson3-3.0.6.tar.gz (131 kB)
    131.5/131.5 kB 21.9 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting requests-html<0.11.0,>=0.10.0 (from facebook_scraper)
  Downloading requests_html-0.10.0-py3-none-any.whl (13 kB)
Requirement already satisfied: python-dateutil<=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz<=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2023.3.post1)
Requirement already satisfied: regex<=2019.02.19,!=2021.8.27 in /usr/local/lib/python3.10/dist-packages (from dateparser<2.0.0,>=1.0.0->facebook_scraper) (2023.6.3)
Requirement already satisfied: six<=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil<=2.8.1->pandas) (1.16.0)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from requests-html<0.11.0,>=0.10.0->facebook_scraper) (2.31.0)
Collecting pyquery (from requests-html<0.11.0,>=0.10.0->facebook_scraper)
  Downloading pyquery-2.0.0-py3-none-any.whl (22 kB)
Collecting fake-useragent (from requests-html<0.11.0,>=0.10.0->facebook_scraper)
  Downloading fake_useragent-1.4.0-py3-none-any.whl (15 kB)
...
Found existing installation: urllib3 2.0.7
Uninstalling urllib3-2.0.7:
  Successfully uninstalled urllib3-2.0.7
Successfully installed bs4-0.0.1 cssselect-1.2.0 dateparser-1.2.0 demjson3-3.0.6 facebook_scraper-0.2.59 fake-useragent-1.4.0 parse-1.20.0 pyee-8.2.2 pyppeteer-1.0.2 pyquery-2.0.0 requ
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

#### 2. File cookies.txt

Trong quá trình crawl data, việc sử dụng cookies có một số ứng dụng:

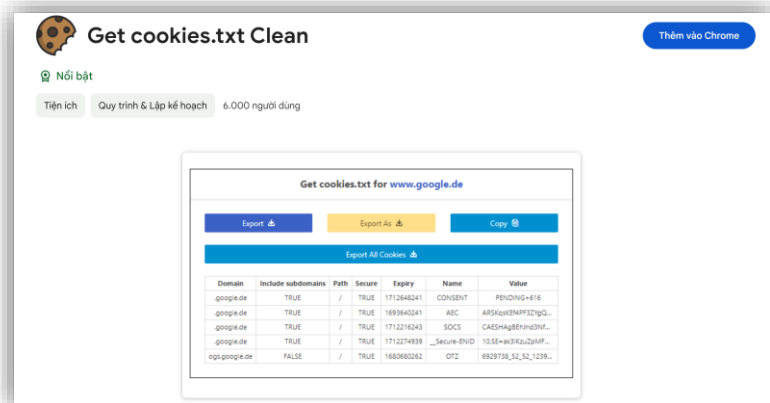
- Thu thập dữ liệu công khai: thông tin người dùng, bài đăng, hình ảnh và các thông tin công cộng khác
- Tích hợp với FB API: xác thực và truy cập vào dữ liệu API yêu cầu

Các bước thực hiện:

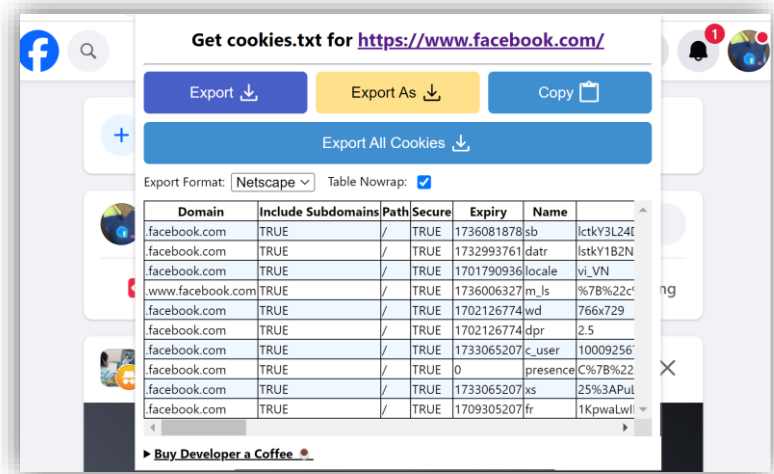
Bước 1: Thêm tiện ích

get cookies

[Tại đây](#)



## Bước 2: Mở facebook và download cookies



## II, Cào dữ liệu

### 1. Thực hành

Ở bài này em cào dữ liệu của fanpage Hai Chiều trên google colab

```
[ ] FANPAGE_LINK = "2chieu"
FOLDER_PATH = "/content/drive/MyDrive/Crawl"
COOKIE_PATH = "/content/drive/MyDrive/Crawl/cookies.txt"

PAGES_NUMBER = 15 # Number of pages to crawl

[ ] post_list = []
for post in get_posts(FANPAGE_LINK,
                      options={"comments": True, "reactions": True, "allow_extra_requests": True},
                      extra_info=True, pages=PAGES_NUMBER, cookies=COOKIE_PATH):
    print(post)
    post_list.append(post)

warnings.warn(
    ERROR:facebook_scraper.facebook_scraper:Exception while requesting URL: https://m.facebook.com/photo.php?fbid=944801937004120&id=1000442296794018&set=a.210658970418424&eav=Afbxn7Pd
Exception: HTTPError('500 Server Error: Internal Server Error for url: https://m.facebook.com/photo.php?fbid=944801937004120&id=1000442296794018&set=a.210658970418424&eav=Afbxn7Pd
Traceback (most recent call last):
  File "/usr/local/lib/python3.10/dist-packages/facebook_scraper/facebook_scraper.py", line 880, in get
    response.raise_for_status()
  File "/usr/local/lib/python3.10/dist-packages/requests/models.py", line 1021, in raise_for_status
    raise HTTPError(http_error_msg, response=self)
requests.exceptions.HTTPError: 500 Server Error: Internal Server Error for url: https://m.facebook.com/photo.php?fbid=944801937004120&id=1000442296794018&set=a.210658970418424&eav=Afbxn7Pd
ERROR:facebook_scraper.extractors:500 Server Error: Internal Server Error for url: https://m.facebook.com/photo.php?fbid=944801937004120&id=1000442296794018&set=a.210658970418424&eav=Afbxn7Pd
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_video didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_video_thumbnail didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_video_id didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_video_meta didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_factcheck didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_share_information didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_listing didn't return anything
WARNING:facebook_scraper.extractors:[944802137004100] Extract method extract_with didn't return anything
ERROR:facebook_scraper.extractors:Unable to parse comment <Element 'div' class=('_2a_i', '_2a_l') data-store={'token':"944802137004100_1500544807410270"} id='1500544807410270' data-store={'token':"944802137004100_1020135569253460"} id='1020135569253460' data-unique-id='1020135569253460' data-uniqu
ERROR:facebook_scraper.extractors:Unable to parse comment <Element 'div' class=('_2a_i', '_2a_l') data-store={'token':"944802137004100_288598836958597"} id='288598836958597' data-unique-id='288598836958597' data-uniqu
ERROR:facebook_scraper.extractors:Unable to parse comment <Element 'div' class=('_2a_i', '_2a_l') data-store={'token':"944802137004100_1038349844110198"} id='1038349844110198' data-unique-id='1038349844110198' data-uniqu
ERROR:facebook_scraper.extractors:Unable to parse comment <Element 'div' class=('_2a_i', '_2a_l') data-store={'token':"944802137004100_873761391068159"} id='873761391068159' data-unique-id='873761391068159' data-uniqu
ERROR:facebook_scraper.extractors:Unable to parse comment <Element 'div' class=('_2a_i', '_2a_l') data-store={'token':"944802137004100_729037605925106"} id='729037605925106' data-unique-id='729037605925106' data-uniqu
```

...

Số lượng bài là 150

```
[ ] len(post_list)
```

150

# Bài đầu tiên

```
{
  'post_list': [
    {
      'post_id': '944802137004100',
      'text': 'Cuốn quá hong dứt ra được \U0001fae2',
      'post_text': 'Cuốn quá hong dứt ra được \U0001fae2',
      'shared_text': '',
      'original_text': None,
      'time': datetime.datetime(2023, 11, 28, 11, 12, 18),
      'timestamp': '1701169938',
      'image': None,
      'image_lowquality': 'https://scontent-ord5-2.xx.fbcdn.net/v/t39.30808-6/405182200_944801930317454_9185539492962934583_n.jpg?sto=c0_dst-jpg_e15_p320x320_g658_nc_cat=106&ccb=1-78_nc_sidsab73678f9e5y0j0jdC798_nc_phc=hf7n0-InYAX9-ZH1i8_nc_ht=scontent-ord5-2.xx&oh=0e_AfC_oAFoICwAtP0R0m1PDa14EwK3HrcsNIUB7NkshZouZ0&oe=656E0FB5',
      'images': [],
      'images_description': [],
      'images_lowquality': 'https://scontent-ord5-2.xx.fbcdn.net/v/t39.30808-6/405182200_944801930317454_9185539492962934583_n.jpg?sto=c0_dst-jpg_e15_p320x320_g658_nc_cat=106&ccb=1-78_nc_sidsab73678f9e5y0j0jdC798_nc_phc=hf7n0-InYAX9-ZH1i8_nc_ht=scontent-ord5-2.xx&oh=0e_AfC_oAFoICwAtP0R0m1PDa14EwK3HrcsNIUB7NkshZouZ0&oe=656E0FB5',
      'image_lowquality_description': '["Cổ thể là hình minh họa về văn bản cho biết 'TÔI KHI CHƠI 1 BÀI NHẠC KÌ CỰC CÙNG LÀ TÔI CÀ NGÂY HÂM ĐỒ: Làn đầu tiên Grav thành ló trong mĩ tơn/ 7/2u hieushieu suzustudio SUZU STUDIO"]',
      'video': None,
      'video_duration_seconds': None,
      'video_height': None,
      'video_id': None,
      'video_quality': None,
      'video_size_MB': None,
      'video_thumbnail': None,
      'video_watches': None,
      'video_width': None,
      'likes': None,
      'comments': 2202,
      'shares': 671,
      'post_url': 'https://facebook.com/zchieu/posts/944802137004100',
      'link': None,
      'links': [
        {
          'link': '/story.php?story_fb_id=pfbiD03320x04b4t8fyefneJogrKcMRhtCeKzfzC2wq65mTwc3LE9FgtYp3qW77FChR1&id=100044229679401&eav=AfBfQsV9-PFnFrYgV1avBKXrKXbkAovsR-96KlFC6fLMB17H6p2m2x5HYHJAC5K&m_entstream_sourceurl=link&refid=178_ft_encrypted_tracking_data.0AYW8RCgltJ-bzqr_gDw0MqcnD6t1VERpAQ0wS1IpYN0SQ4WuRhwHw5m4V6u0x8BG6Aqhv2v1ChYvYmQ_jnV-uru09TaTDzmfzHw4Y0UYUW385FmR1R-B-DlXuUftJumknt2xk7LMtyxL110e1xc5g0mK3thwFT0F50wQ1pBgwN58C-RPLQA1sUa0BLXt7RGJgEyhCBV5311L00Z2DcwDYMVG7BuIT1EYJ30mNG30R2tJ325BUP7q-WqJwHmWVRXZG1Fq-U50UtzFngly1-1',
        }
      ]
    }
  ]
}
```

• • •

## 2. Tổ chức dữ liệu thô

[illegible]

...

## 10 post đầu tiên:

id	post_id	text	post_text	shared_text	original_text	time	timestamp	image	image_lowquality	image	width	page_id	shares	image_id	image_id	was_live	fetch_time	video_id	video_id	header
0	94480137004100	Cười quá hơng đi ra được	Cười quá hơng đi ra được		None	2023-11-11 12:18	1701169958	None	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	94480137004100	94480137004100	False	2023-11-30 11:16:06.607672	NaN	NaN	NaN
0	944221083723892	em mới hai tuổi, mà cười ngây đến phớt lệ, 10...	em mới hai tuổi, mà cười ngây đến phớt lệ, 10...		None	2023-11-11 13:23	1701086375	https://m.facebook.com/photos/view_full_size/?...	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	944221083723892	944221083723892	False	2023-11-30 11:16:36.348809	NaN	NaN	NaN
0	943155540502090	Trong vườn sảnh "khảm pha độc từ đến tượng bà...	Trong vườn sảnh "khảm pha độc từ đến tượng bà...		None	2023-11-12 30:50	1700915450	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	None		False	2023-11-30 11:17:42.631546	NaN	NaN	NaN	
0	942579587186333	Khi nào thì nhà không còn là môi để trú ẩn?	Khi nào thì nhà không còn là môi để trú ẩn?		None	2023-11-04 28:38	1700880318	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	None		False	2023-11-30 11:18:05.629222	NaN	NaN	NaN	
0	942641393888641	Dưới đó là đồng người, là mặt đường hay là cái...	Dưới đó là đồng người, là mặt đường hay là cái...		None	2023-11-12 30:20	1700879820	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	942641393888641	942641393888641	False	2023-11-30 11:18:18.103905	NaN	NaN	NaN	
0	941979430610704	Hôm trước, mình ngồi xem phim với anh đồng nghiệp...	Hôm trước, mình ngồi xem phim với anh đồng nghiệp...		None	2023-11-04 11:18	1700715078	https://content-cdf-2.xx.fbcdn.net/v/139.306...	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	941979430610704	941979430610704	False	2023-11-30 11:18:18.587385	NaN	NaN	NaN
0	940839840364988	4 tuổi rồi đã biết lấy tay kim nước mắt, thì ...	4 tuổi rồi đã biết lấy tay kim nước mắt, thì ...		None	2023-11-09 01:01	1700557261	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	940839840364988	940839840364988	False	2023-11-30 11:18:21.290764	NaN	NaN	NaN	
0	939784010830440	Ở trường, những bạn học không chỉ cần sức khỏe...	Ở trường, những bạn học không chỉ cần sức khỏe...		None	2023-11-13 01:39	1700398899	https://content-cdf-2.xx.fbcdn.net/v/139.306...	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	939784010830440	939784010830440	False	2023-11-30 11:19:48.549092	NaN	NaN	NaN
0	938410671948340	Không có chiến là không được mà được	Không có chiến là không được mà được		None	2023-11-09 03:40	1700181400	https://m.facebook.com/photos/view_full_size/?...	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	938410671948340	938410671948340	False	2023-11-30 11:20:28.287444	NaN	NaN	NaN
0	934231571392290	Cười đó là nụ tươi tươi, mà đôi mắt thì biết...	Cười đó là nụ tươi tươi, mà đôi mắt thì biết...		None	2023-11-04 30:00	1699763400	https://content-cdf-2.xx.fbcdn.net/v/139.306...		None	246384732960019	None	934231571392290	934231571392290	False	2023-11-30 11:21:19.007080	NaN	NaN	NaN	

## 3. Lưu file

Lưu file dưới 2 dạng :

- Csv: để phân tích bài viết (nội dung, tương tác,...)
- Npy: để phân tích bình luận và người bình luận

```
path1=FOLDER_PATH + FANPAGE_LINK + ".npy"
np.save(path1, arr) # .npy extension is added if not given
print(path1)
}

/content/drive/MyDrive/Craw2chieu.npy

# To df
path=FOLDER_PATH + FANPAGE_LINK + "new.csv"
post_df_full.to_csv(path, index=False)
print(path)
}

/content/drive/MyDrive/Craw2chieunew.csv
```

## Phần 2: Tiền xử lý dữ liệu

### I. Bài viết

#### Dataframe trước

	post_id	text	post_text	shared_text	original_text	time	timestamp	image	image_lowquality	images	...	with	page_id
0	944802137004100	Cuốn quá hong dứt ra được	Cuốn quá hong dứt ra được	NaN	NaN	2023-11-28 11:12:18	1701169938	NaN	https://scontent-ord5-2.xx.fbcdn.net/v/t39.308...	[]	--	NaN	246384732965019
1	944225863728392	elm mới hai tuổi, mà suốt ngày dính phở-ớt, tồ...	elm mới hai tuổi, mà suốt ngày dính phở-ớt, tồ...	NaN	NaN	2023-11-27 11:59:35	1701086375	https://m.facebook.com/photo/view_full_size/?f...	https://scontent-ord5-2.xx.fbcdn.net/v/t39.308...	[https://m.facebook.com/photo/view_full_size/...	--	NaN	246384732965019
2	943155540502093	Trong cuốn sách "Khám phá đầu trẻ bên trong bạ...	Trong cuốn sách "Khám phá đầu trẻ bên trong bạ...	NaN	NaN	2023-11-25 12:30:50	1700915450	NaN	https://scontent-ord5-2.xx.fbcdn.net/v/t15.525...	[]	--	NaN	246384732965019
3	942979587186355	Khi nào thì nhà không còn là nơi để trở về? \nK...	Khi nào thì nhà không còn là nơi để trở về? \nK...	NaN	NaN	2023-11-25 04:28:38	1700886518	NaN	https://scontent-ord5-2.xx.fbcdn.net/v/t15.525...	[]	--	NaN	246384732965019
4	942643393886641	Dưới đó là dòng người, là mặt đường hay là cái...	Dưới đó là dòng người, là mặt đường hay là cái...	NaN	NaN	2023-11-24 12:30:20	1700829020	NaN	https://scontent-ord5-2.xx.fbcdn.net/v/t39.308...	[]	--	NaN	246384732965019
...	...	...	...	...	...	...	...	...	...	...	...	...	...

#### Thông tin

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 54 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   post_id                               150 non-null    int64
1   text                                  149 non-null    object
2   post_text                             148 non-null    object
3   shared_text                           4 non-null      object
4   original_text                         0 non-null      float64
5   time                                  150 non-null    object
6   timestamp                             150 non-null    int64
7   image                                 122 non-null    object
8   image_lowquality                      149 non-null    object
9   images                                150 non-null    object
10  images_description                    150 non-null    object
11  images_lowquality                     150 non-null    object
12  images_lowquality_description          150 non-null    object
13  video                                 3 non-null      object
14  video_duration_seconds                 0 non-null      float64
15  video_height                           0 non-null      float64
16  video_id                               3 non-null      float64
17  video_quality                          0 non-null      float64
18  video_size_MB                          0 non-null      float64
19  video_thumbnail                       3 non-null      object
...
52  videos                                6 non-null      object
53  header                                4 non-null      object
dtypes: bool(3), float64(18), int64(7), object(26)
memory usage: 60.3+ KB
```



Nhìn qua chúng ta thấy nhiều cột bị thiếu dữ liệu( bài đăng không có nội dung, hạn chế quyền truy cập của token,..) vì thế để thuận tiện cho bước phân tích ta sẽ bỏ những cột không cần thiết, và cài lại giá trị những hàng hoặc cột nếu cần thiết

Dataframe sau

	post_id	text	post_text	time	image	video	comments	shares	reactions	reaction_count
0	944802137004100	Cuốn quá hong dứt ra được 🤔	Cuốn quá hong dứt ra được 🤔	2023-11-28 11:12:18	NaN	NaN	2202	671	('thích': 4462, 'yêu thích': 82, 'haha': 7902,...	14595
1	944225883728392	elm mới hai tuổi, mà suốt ngày dính phở-ớt, tó...	elm mới hai tuổi, mà suốt ngày dính phở-ớt, tó...	2023-11-27 11:59:35	<a href="https://m.facebook.com/photo/view_full_size/?f...">https://m.facebook.com/photo/view_full_size/?f...</a>	NaN	446	148	('thích': 9850, 'yêu thích': 1206, 'haha': 521...	15157
2	943155540502093	Trong cuốn sách "Khám phá đứa trẻ bên trong bạ...	Trong cuốn sách "Khám phá đứa trẻ bên trong bạ...	2023-11-25 12:30:50	NaN	<a href="https://scontent-ord5-2.xx.fbcdn.net/v/t42.179...">https://scontent-ord5-2.xx.fbcdn.net/v/t42.179...</a>	82	70	('thích': 946, 'yêu thích': 230, 'wow': 3, 'th...	1234
3	942979587186355	Khi nào thì nhà không còn là nơi để trở về? \nK...	Khi nào thì nhà không còn là nơi để trở về? \nK...	2023-11-25 04:28:38	NaN	<a href="https://scontent-ord5-2.xx.fbcdn.net/v/t42.179...">https://scontent-ord5-2.xx.fbcdn.net/v/t42.179...</a>	15	33	('thích': 383, 'yêu thích': 210, 'haha': 2, 'w...	661
4	942643393886641	Dưới đó là dòng người, là mặt đường hay là cái...	Dưới đó là dòng người, là mặt đường hay là cái...	2023-11-24 12:30:20	NaN	NaN	30	83	('thích': 1002, 'yêu thích': 348, 'haha': 3, '...	1546
...	...	...	...	...	...	...	...	...	...	...

Thông tin

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   post_id         150 non-null   int64
1   text            150 non-null   object
2   post_text       150 non-null   object
3   time            150 non-null   object
4   image           150 non-null   int32
5   video           150 non-null   int32
6   comments        150 non-null   int64
7   shares          150 non-null   int64
8   reactions       150 non-null   object
9   reaction_count  150 non-null   int64
dtypes: int32(2), int64(4), object(4)
memory usage: 10.7+ KB
```

## II. Bình luận

Bình luận chưa xử lí là 1 phần trong file HaiChieu.npy

```
{ 'comment_id': '365993212562670', 'comment_url': 'https://facebook.com/365993212562670', 'commenter_id': '100043997325521', 'commenter_url': 'https://facebook.com/thanganhconem2eavvAfBx',
  'comment_id': '755440759961813', 'comment_url': 'https://facebook.com/755440759961813', 'commenter_id': '100082965051760', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=100082',
  'comment_id': '883063806457448', 'comment_url': 'https://facebook.com/883063806457448', 'commenter_id': '100066800510946', 'commenter_url': 'https://facebook.com/tram4520eavvAfafk_yfu',
  'comment_id': '783450291731155', 'comment_url': 'https://facebook.com/783450291731155', 'commenter_id': '100022189274901', 'commenter_url': 'https://facebook.com/hien.trandieu.5764?af',
  'comment_id': '1283450549013879', 'comment_url': 'https://facebook.com/1283450549013879', 'commenter_id': '100091124224730', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=1000',
  'comment_id': '7738791419488573', 'comment_url': 'https://facebook.com/7738791419488573', 'commenter_id': '100079552307144', 'commenter_url': 'https://facebook.com/HuynhCuti2eavvAfYCO',
  'comment_id': '1129697324679857', 'comment_url': 'https://facebook.com/1129697324679857', 'commenter_id': '61551250733345', 'commenter_url': 'https://facebook.com/nguong22082eavvAfaf9',
  'comment_id': '1375848600109978', 'comment_url': 'https://facebook.com/1375848600109978', 'commenter_id': '100074399505527', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=1000',
  'comment_id': '896439598773387', 'comment_url': 'https://facebook.com/896439598773387', 'commenter_id': '100075142973415', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=100075',
  'comment_id': '2442989502554514', 'comment_url': 'https://facebook.com/2442989502554514', 'commenter_id': '100012827495073', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=1000',
  'comment_id': '6902245030888462', 'comment_url': 'https://facebook.com/6902245030888462', 'commenter_id': '100073785095036', 'commenter_url': 'https://facebook.com/lyvne.092eavvAfYmVvIX',
  'comment_id': '6535400453254440', 'comment_url': 'https://facebook.com/6535400453254440', 'commenter_id': '100092903115114', 'commenter_url': 'https://facebook.com/fling38th12eavvAfZ',
  'comment_id': '1307719116600750', 'comment_url': 'https://facebook.com/1307719116600750', 'commenter_id': '100082602784461', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=1000',
  'comment_id': '310015538651954', 'comment_url': 'https://facebook.com/310015538651954', 'commenter_id': '100006631743247', 'commenter_url': 'https://facebook.com/gemstone27952eavvAfYU',
  'comment_id': '663962988940545', 'comment_url': 'https://facebook.com/663962988940545', 'commenter_id': '100041859378844', 'commenter_url': 'https://facebook.com/th.ngen2eavvAfBf9LlX',
  'comment_id': '881160939407139', 'comment_url': 'https://facebook.com/881160939407139', 'commenter_id': '100037587129863', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=100037',
  'comment_id': '729033465910843', 'comment_url': 'https://facebook.com/729033465910843', 'commenter_id': '100012813404450', 'commenter_url': 'https://facebook.com/tamhu.1e.90812eavvAfY',
  'comment_id': '895914201993203', 'comment_url': 'https://facebook.com/895914201993203', 'commenter_id': '100071878036637', 'commenter_url': 'https://facebook.com/nhicuti15vn2eavvAfY',
  'comment_id': '1380205486256733', 'comment_url': 'https://facebook.com/1380205486256733', 'commenter_id': '100011249204666', 'commenter_url': 'https://facebook.com/tramss.nguyenv2eavv',
  'comment_id': '760717819217496', 'comment_url': 'https://facebook.com/760717819217496', 'commenter_id': '100062854610936', 'commenter_url': 'https://facebook.com/julii.julii.94842eavvA',
  'comment_id': '1543654649788402', 'comment_url': 'https://facebook.com/1543654649788402', 'commenter_id': '100014953336998', 'commenter_url': 'https://facebook.com/phuongtruc.nguyen.1',
  'comment_id': '249410978141387', 'comment_url': 'https://facebook.com/249410978141387', 'commenter_id': '100026963393521', 'commenter_url': 'https://facebook.com/athu.50702762eavvAfBc',
  'comment_id': '1283512088981307', 'comment_url': 'https://facebook.com/1283512088981307', 'commenter_id': '100007625986479', 'commenter_url': 'https://facebook.com/KhanhNgoc32502eavvA',
  'comment_id': '889984165811383', 'comment_url': 'https://facebook.com/889984165811383', 'commenter_id': '100044229679401', 'commenter_url': 'https://facebook.com/zchieu2eavvAfBjVrCO3H',
  'comment_id': '301219232882997', 'comment_url': 'https://facebook.com/301219232882997', 'commenter_id': '100072819962559', 'commenter_url': 'https://facebook.com/profile.php?id=100072',
  ...
  'comment_id': '548096843974438', 'comment_url': 'https://facebook.com/548096843974438', 'commenter_id': '100006064170185', 'commenter_url': 'https://facebook.com/minhhang.10032eavvAfY',
  'comment_id': '695766505661064', 'comment_url': 'https://facebook.com/695766505661064', 'commenter_id': '100052954426741', 'commenter_url': 'https://facebook.com/hanhuong.ho.73157207e',
  'comment_id': '228414213035658', 'comment_url': 'https://facebook.com/228414213035658', 'commenter_id': '100048765469085', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=100048',
  'comment_id': '237942805319139', 'comment_url': 'https://facebook.com/237942805319139', 'commenter_id': '100023845962353', 'commenter_url': 'https://facebook.com/nrofile.rhn?id=100023',
  ... }
```

Sau xử lí

	commenter_id	commenter_name	comment_text	comment_time
0	100043997325521	Thắng Anh Con Em	Nghe riết cái thấy cũng thích thích hay hay 🤔	2023-11-23
1	100082965051760	Trang Lê	Bài gì z:))	2023-11-23
2	100066800510946	Bích Trâm	Nghe xong ăn mì 3 miễn	2023-11-23
3	100022189274901	Trần Diệu Hiền	Này thiệt nha, kiểu bị cuốn cuốn sao í 🤔	2023-11-23
4	100091124224730	Hân Lê	Hay	2023-11-23
...	...	...	...	...
7466	100029720314950	Kiều Hồ Trung Dũng	Nếu nói theo hướng tích cực thì những cái cây ...	2023-03-30
7467	100022488562817	Đức Đức	Đào Lê Tâm An	2023-03-30
7468	100012760427859	Huỳnh Trang Hoàng Mỹ	Minh Duy Trịnh Phan để cho dễ nhận diện thôi...	2023-03-30
7469	100002686756259	Minh Duy Trịnh Phan	Huỳnh Trang Hoàng Mỹ mình thích ý nghĩ của bạn 🤔	2023-03-30
7470	100081237549027	Ngọc Như Ý	Hà Phương sẽ thật tuyệt zôiiii	2023-03-30

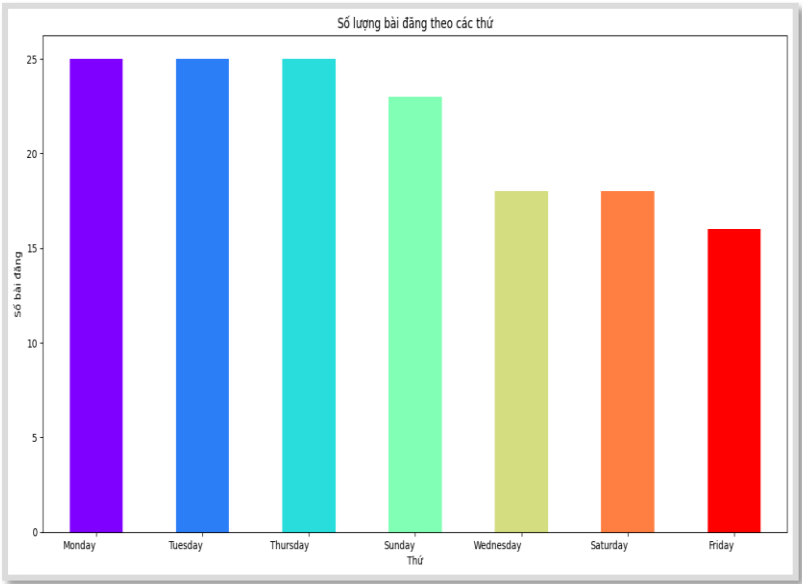
7471 rows × 4 columns

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7471 entries, 0 to 7470
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   commenter_id    7471 non-null   object
1   commenter_name  7471 non-null   object
2   comment_text    7471 non-null   object
3   comment_time    7471 non-null   datetime64[ns]
dtypes: datetime64[ns](1), object(3)
memory usage: 233.6+ KB
```

Phần 3: Phân tích

I. Bài viết

1. Bài đăng theo các thứ



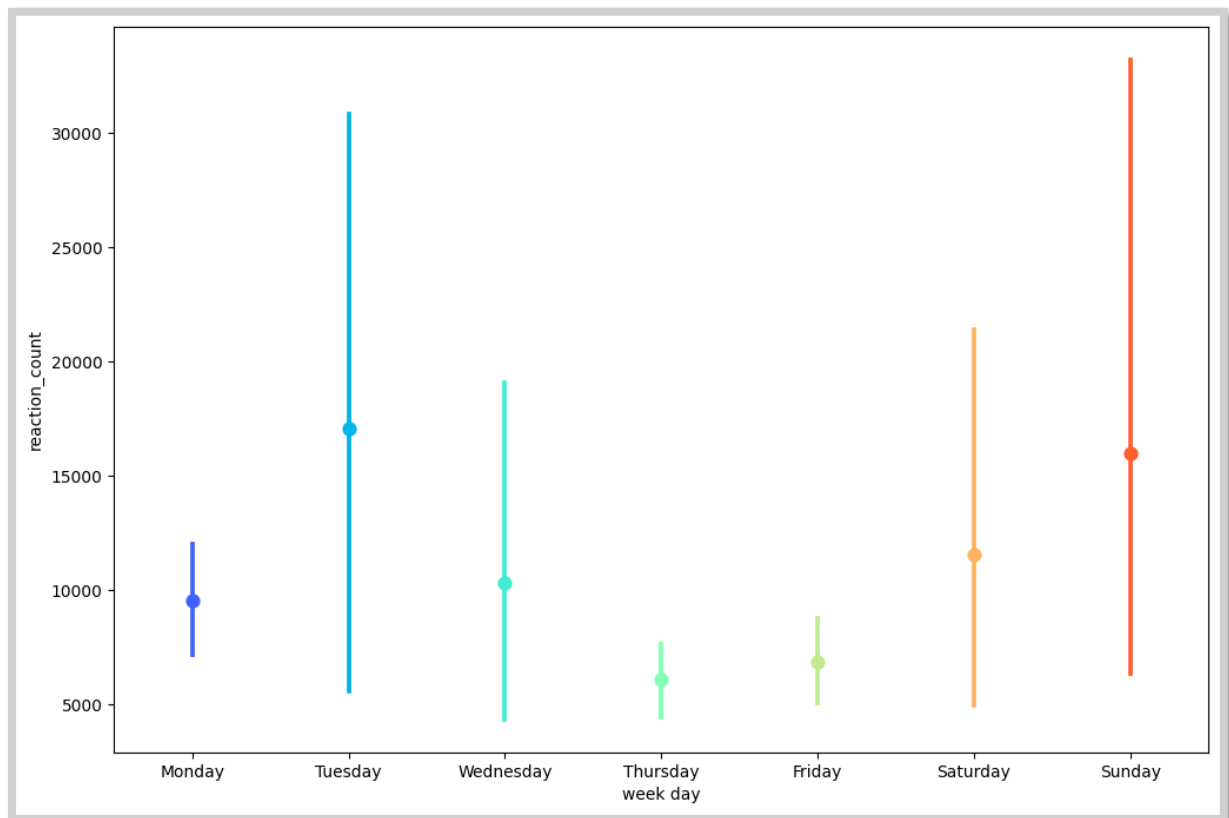
Ta thấy thứ **hai, ba và năm** cùng có số bài đăng là **nhiều nhất**, 25 bài ; ngày có số bài đăng **ít nhất** là thứ **sáu**

⇒ Như vậy fanpage này chủ yếu hoạt động vào đầu tuần và kém hoạt động hơn vào cuối tuần

Các bài đăng vào hôm thứ 2

	post_id	text	post_text	time	image	video	comments	shares	reaction_count	week_day
1	944225883728392	elm mới hai tuổi, mà suốt ngày dính phờ-ớt, tó...	elm mới hai tuổi, mà suốt ngày dính phờ-ớt, tó...	2023-11-27 11:59:35	1	0	446	148	15157	Monday
20	919536056197375	Ước gì có cái công tắc để bật tắt mọi suy nghĩ...	Ước gì có cái công tắc để bật tắt mọi suy nghĩ...	2023-10-16 12:00:48	1	0	296	603	7312	Monday
33	902544574563190	Có những tổn thương chỉ lộ ra khi người ta nhậ...	Có những tổn thương chỉ lộ ra khi người ta nhậ...	2023-09-18 12:10:38	1	0	304	1210	15843	Monday
43	890615962422718	Khi bạn còn trẻ, họ mặc định bạn không biết gì...	Khi bạn còn trẻ, họ mặc định bạn không biết gì...	2023-08-28 12:46:14	1	0	2558	3174	6015	Monday
47	886134539537527	Có biết vì sao bạn hay yêu nhầm người không? V...	Có biết vì sao bạn hay yêu nhầm người không? V...	2023-08-21 01:33:54	1	0	747	239	14899	Monday

## Đăng bài vào thứ mấy thì nhận được nhiều tương tác ?



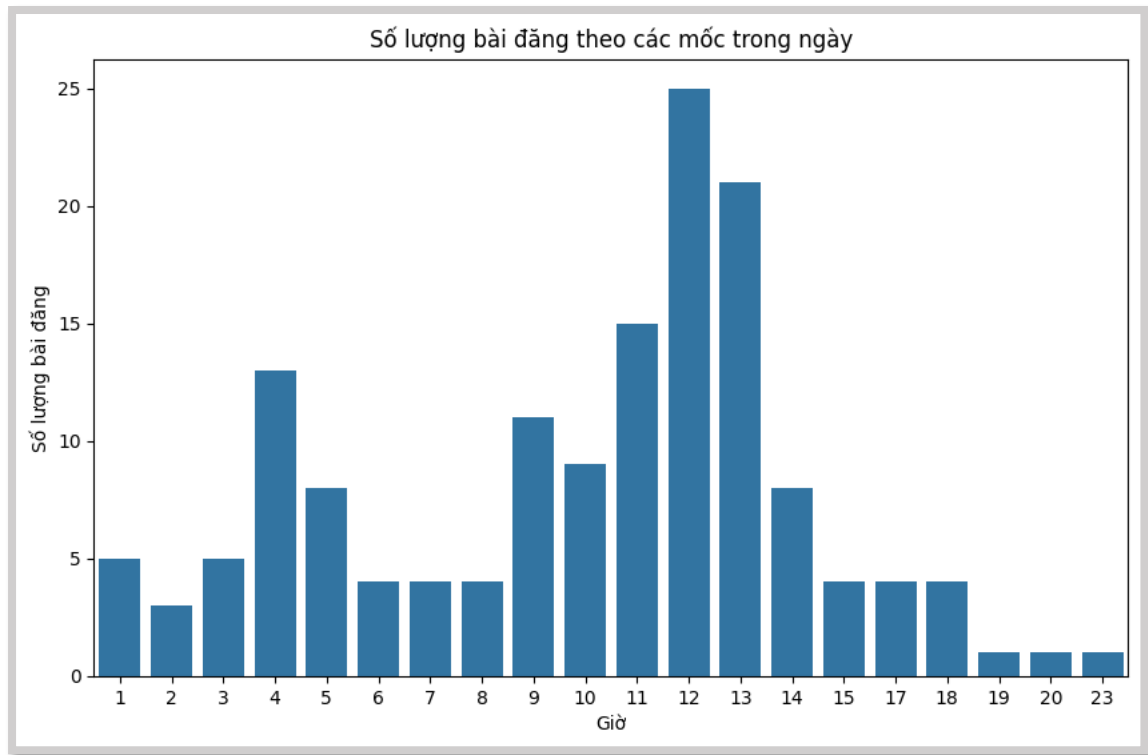
- Trung bình thứ ba, chủ nhật sẽ nhận được nhiều tương tác nhất, thứ năm, sáu nhận được ít nhất
- Thứ năm và sáu đăng tương đối nhiều (41/150 bài) nhưng nhận lại không được nhiều tương tác

Bốn ngày có nhiều bài viết nhất

```
Monday: 25  
Sunday: 23  
Friday: 16  
Thursday: 25
```

- ⇒ Đăng bài vào thứ ba hoặc chủ nhật để nhận được nhiều tương tác

## Các mốc thời gian trong ngày hay đăng bài

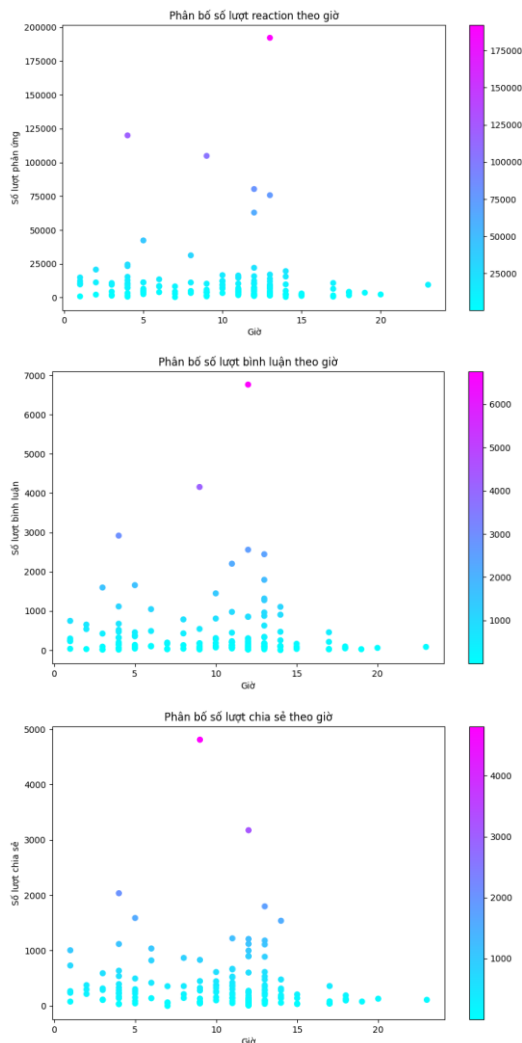


- Mốc 11-13h đăng nhiều bài nhất là 61 bài
- Mốc 19-2h đăng ít bài nhất là 11 bài

```
... Mốc 1 giờ: 5 bài viết
Mốc 2 giờ: 3 bài viết
Mốc 3 giờ: 5 bài viết
Mốc 4 giờ: 13 bài viết
Mốc 5 giờ: 8 bài viết
Mốc 6 giờ: 4 bài viết
Mốc 7 giờ: 4 bài viết
Mốc 8 giờ: 4 bài viết
Mốc 9 giờ: 11 bài viết
Mốc 10 giờ: 9 bài viết
Mốc 11 giờ: 15 bài viết
Mốc 12 giờ: 25 bài viết
Mốc 13 giờ: 21 bài viết
Mốc 14 giờ: 8 bài viết
Mốc 15 giờ: 4 bài viết
Mốc 16 giờ: 0 bài viết
Mốc 17 giờ: 4 bài viết
Mốc 18 giờ: 4 bài viết
Mốc 19 giờ: 1 bài viết
Mốc 20 giờ: 1 bài viết
Mốc 21 giờ: 0 bài viết
Mốc 22 giờ: 0 bài viết
Mốc 23 giờ: 1 bài viết
```

⇒ Fanpage chủ yếu đăng vào lúc trưa, đăng rất ít vào buổi tối

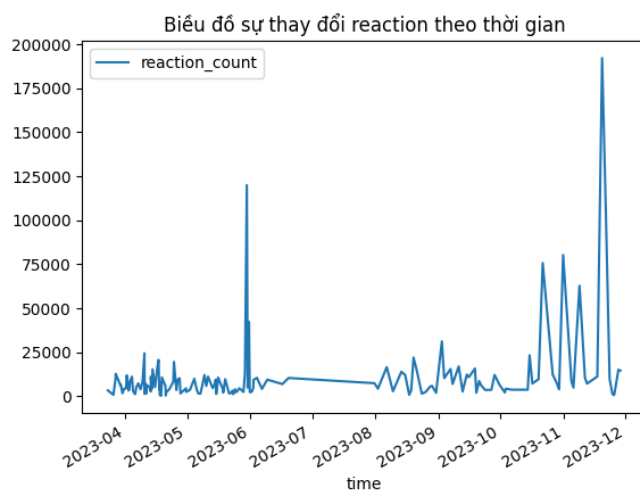
## Phân bố các lượt tương tác theo giờ



- Trong khung giờ từ 11-13h, lượt tương tác nhận lại được rất nhiều ngược lại buổi tối giao đoạn 19-22h rất ít tương tác. Sự khác biệt này phản ánh mẫu quen thuộc hành vi trực tuyến của giới trẻ. Buổi trưa là thời gian nghỉ trưa, mọi người có xu hướng tìm đến những thứ giải trí sau giờ học hoặc làm việc căng thẳng. Trong khi đó, buổi tối có thể là thời gian họ dành cho các hoạt động khác như học tập, gia đình, bạn bè, ...

⇒ Tóm lại, thời gian thích hợp nhất để Fanpage Hai Chiều đăng bài là trưa của thứ ba hoặc chủ nhật, hạn chế đăng vào các tối của thứ năm và sáu vì khi ấy có đăng nhiều cũng không nhận được nhiều phản hồi

## Phân tích về reaction



Độ giao động của reaction\_count là: 22014.571394267226

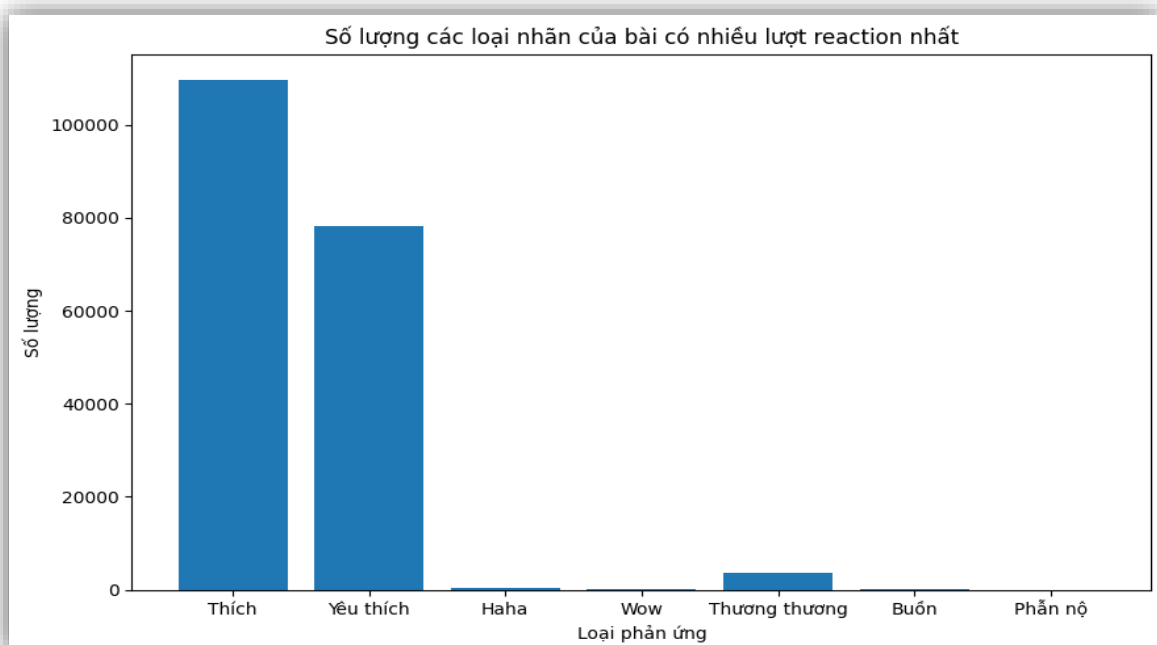
- Số lượng react tăng đột ngột vào cuối tháng 5 và cuối tháng 11
- Có nhiều biến động, độ giao động lớn ~ 22000
- Giai đoạn tháng 6-8 hầu như không có sự phát triển

## Bài viết có nhiều lượt reaction nhất

## Bài viết có ít lượt reaction nhất

Thông tin bài viết có nhiều reaction nhất	
post_id	939784010839246
text	Ở trường, những bài học không chỉ nằm trên bảng.
post_text	Ở trường, những bài học không chỉ nằm trên bảng.
time	2023-11-19 13:01:39
image	1
video	0
comments	1315
shares	1799
comments_full	[{'comment_id': '897054048437725', 'comment_ur...
reaction_count	192186
week_day	Sunday
hour	13
thích	109684.0
yêu thích	78073.0
haha	464.0
wow	103.0
thương thương	3731.0
buồn	125.0
phản nộ	6.0

Thông tin bài viết có ít lượt reaction nhất	
post_id	809295750554740
text	Poster của team, ad dễ cho page công ty =)) lự...
post_text	Poster của team, ad dễ cho page công ty =)) lự...
time	2023-04-18 07:43:49
image	1
video	0
comments	24
shares	1
comments_full	[{'comment_id': '53757775215269', 'comment_ur...
reaction_count	345
week_day	Tuesday
hour	7
thích	226.0
yêu thích	108.0
haha	5.0
wow	1.0
thương thương	5.0
buồn	0.0
phản nộ	0.0

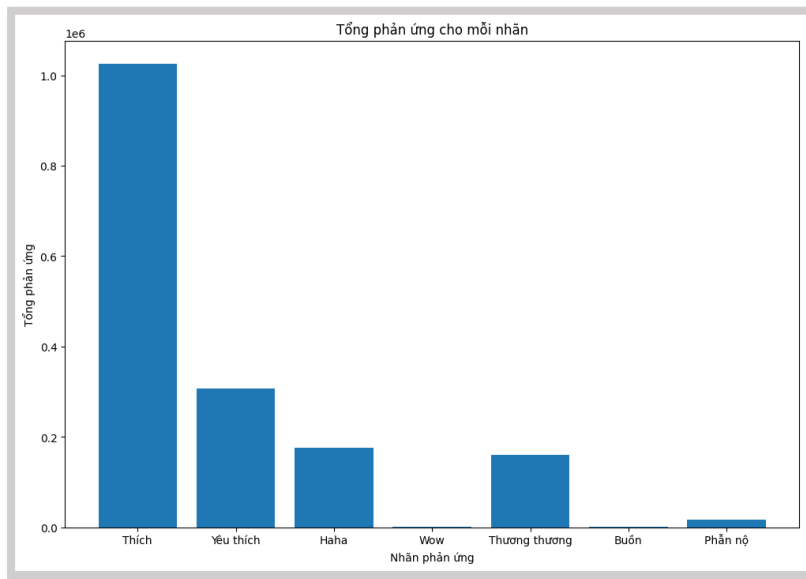


Số lượng reaction bài cao nhất này chủ yếu đến từ nhãn Thích và Yêu thích

⇒ Nhận xét về các ảnh trên:

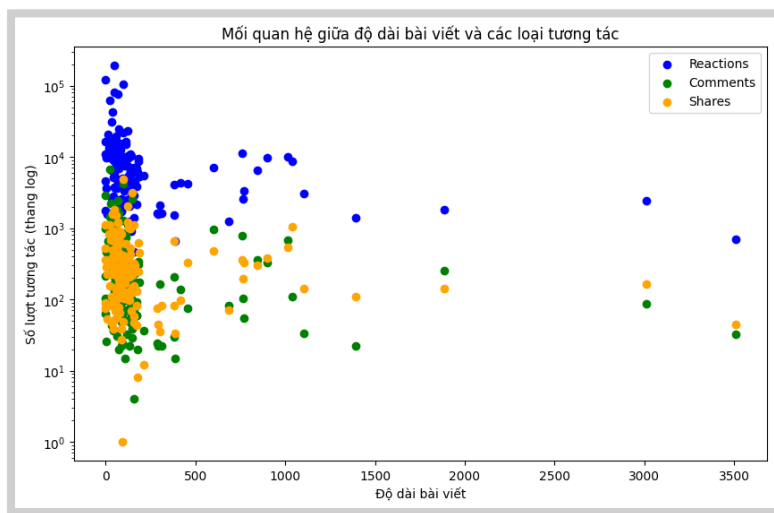
Số lượng reaction tăng đột ngột vào giai đoạn cuối tháng 5 và cuối tháng 11, đạt đỉnh là bài đăng vào lúc 13h ngày 19-11 với 192186 reaction, chủ yếu đến từ các nhãn Thích và Yêu thích; đạt đáy vào lúc 7h ngày 18-4 với 345 reaction. Bài cao nhất gấp 557 lần bài thấp nhất

## Tổng react cho mỗi nhãn



⇒ Lượt react của các bài viết chủ yếu đến từ nhãn Thích, điều này cho thấy người dùng có xu hướng thể hiện sự ủng hộ hoặc sự thích với nội dung một cách đơn giản, thay vì biểu lộ ra cảm xúc thật của mình thông qua các nhãn khác

## Mối quan hệ giữa các loại tương tác và độ dài bài viết



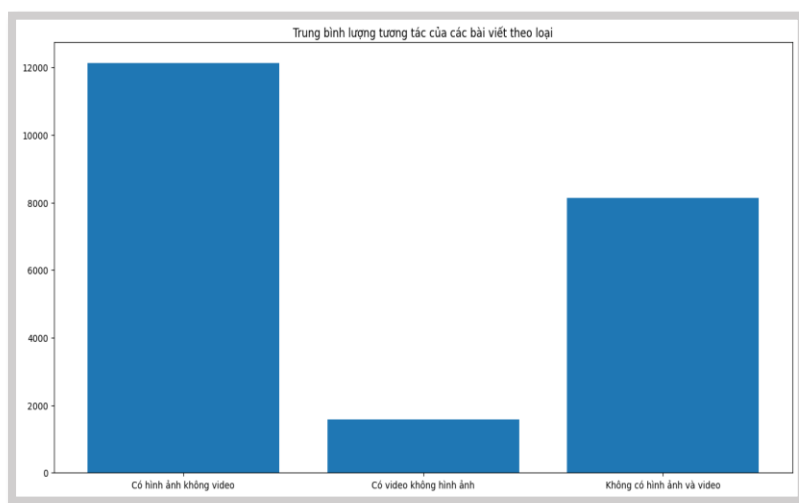
- Văn bản càng ngắn lượt tương tác càng nhiều  
⇒ Người dùng facebook có xu hướng ưa chuộng những nội dung ngắn gọn, dễ tiếp cận, và nhanh chóng



## Việc có hình ảnh hoặc video trong bài viết liệu thu hút người dùng hơn không?

```
Số bài viết có hình ảnh không video: 122
Số bài viết có video không hình ảnh: 3
Số bài viết không có hình ảnh và video: 25
```

Fanpage chủ yếu đăng bài kèm theo hình ảnh



Lượt tương tác khi có hình ảnh là rất cao, trong khi có video thì lại rất thấp. Điều này có thể cho thấy rằng người dùng thích sự ngắn gọn, dễ tiếp cận. Ngoài ra, nội dung hình ảnh có thể cung cấp thông tin một cách trực quan và dễ tiếp nhận hơn video, làm tăng khả năng thu hút sự chú ý và tương tác từ phía người dùng

	comments	shares	reaction_count
0	2202	671	14595
1	446	148	15157
2	82	70	1234
3	15	33	661
4	30	83	1546
...	...	...	...
146	322	637	12811
147	254	1006	9757
148	36	79	891
149	39	145	3409
mean	414	397	11259

⇒ Trung bình mỗi bài viết sẽ có khoảng 400 comments, 400 lượt shares và 11300 lượt react. Dựa trên số liệu này ta có thể khẳng định fanpage có sức hút khá lớn

## Phân tích nội dung bài viết

## Các từ thông dụng



Các từ có thể dễ dàng nhìn thấy như là: có, bạn, không, người, ...

10 từ xuất hiện nhiều nhất trong bài viết

[('có', 141), ('là', 127), ('không', 115), ('bạn', 100), ('một', 96), ('người', 91), ('của', 69), ('những', 65), ('và', 63), ('sẽ', 58)]

Tìm thử một từ khóa bất kì

**Tổng số 150 bài viết**

3 bài viết chứa từ khóa 'xã hội'

[ 'Mạng xã hội đang dần biến con người thành một cỗ máy vô cảm khi nhiều người tự nhốt bản thân vào đó mà mất dần khả năng kết nối với thế giới thực.',  
'Là một sinh vật xã hội, chúng ta luôn khao khát một định hướng, một điểm tựa nhất định trong khi không ai thực sự hiểu bản thể của mình như thế nào.',  
'Mạng xã hội thực chất là nơi kết nối...những kẻ đơn luyến/Nghe có vẻ cực đoan ha, nhưng một nghiên cứu năm 2021 chỉ ra rằng 10 ngày sử dụng các nền tảng mạng :

## Trong bài viết có emoji không?

Có 72 emoji được sử dụng trong 150 bài viết

Trong đó có 46 bài viết sử dụng emoji

⇒ Có tận 46 bài viết chứa emoji, 72 emoji đã được sử dụng. Chủ fanpage có thể thuộc thế hệ Gen Z hoặc có xu hướng sử dụng các phong cách giao tiếp phổ biến trong thế hệ này. Việc sử dụng rộng rãi emoji của tác giả phản ánh sự sáng tạo và mong muốn tạo ra nội dung thú vị, đáng yêu và gần gũi với cộng đồng facebook

## 1. Phân tích comment

## Độ dài trung bình của comment?

Độ dài trung bình của bình luận là: 82 kí tự

## Các từ thông dụng



Các từ thông dụng có thể dễ dàng nhìn thấy là: là, thì, mình, có, mà, ...

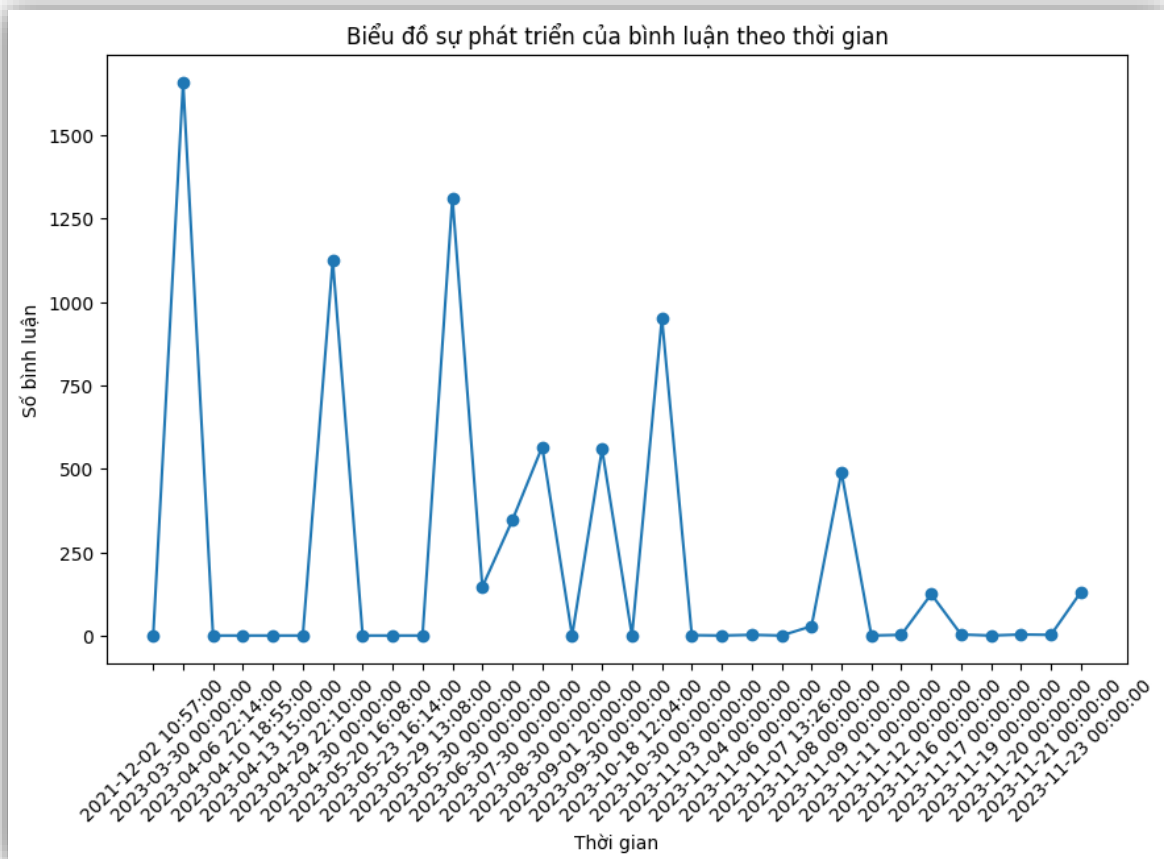
10 từ xuất hiện nhiều nhất trong nội luận  
[('là', 2483), ('có', 1951), ('thì', 1610), ('mình', 1565), ('không', 1238), ('ngươi', 1208), ('bạn', 1158), ('mà', 1124), ('cũng', 1041), ('con', 955)]

## Emoji trong bình luận

Ở phần phân tích nội dung bài viết, ta đã thấy xuất hiện rất nhiều emoji rồi, nên từ đó có thể chắc chắn rằng phần bình luận cũng có.

Có 3249 emoji được sử dụng trong 7471 bình luận  
Trong đó có 2024 bình luận sử dụng emoji

## Độ phát triển của comment theo thời gian



- Ở biểu đồ này, do đã chuẩn hóa một số giá trị nên trông nó không được liên tục đẹp mắt, nhưng nhìn chung có thể thấy được số lượng comment phát triển rất mạnh thời gian đầu tháng 2 – 4 năm 2023 và giảm mạnh về cuối năm
- Có lẽ Fanpage này nổi do các bài viết tạo trend ở thời gian đầu và hiện tại không còn nhiều nữa nên mới gây ra sự suy giảm như này

### Những ai là fan cứng của Fanpage?

Ở đây em dựa theo tiêu chí số lượng comment của mỗi cá nhân để đánh giá

Top 5 người bình luận nhiều nhất là

	commenter_id	commenter_name	count
0	100044229679401	Hai Chiều	321
1	100043997325521	Thằng Anh Con Em	61
2	100065791324776	Xuân Nghi Vương	44
3	100028374075249	Reanzares Limuel	30
4	1829242625	Phan Hồng Đức	28

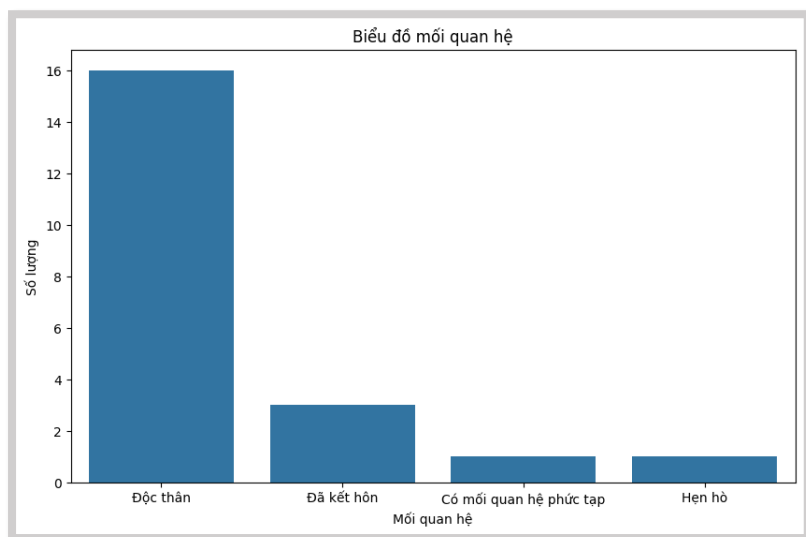
- Xuất hiện tên fanpage “Hai Chiều ” với số lượng nhiều nhất , điều này cho thấy họ có sự tôn trọng đối với cộng đồng của mình, hoặc có thể là họ đang giải đáp thắc của người hâm mộ
- Như vậy, ‘Thăng Anh Con Em’, ‘Xuân Nghi Vương’, ‘Reanzares Limuel’, ‘Phan Hồng Đức’ là các fan cứng.

## 2. Phân tích commenter

Do không đủ kiến thức để crawl lại dữ liệu bị thiếu nên phần này em chỉ phân tích 66 profile nhận được từ dữ liệu cũ

	Followler_count	id	Name	Thông tin cơ bản	Học vấn	Mối quan hệ	Nơi từng sống
0	None	100043997325521	Thăng Anh Con Em	Ngày sinh	NaN	NaN	NaN
1	None	100082965051760	Trang Lê	NaN	NaN	NaN	NaN
2	None	100066800510946	Bích Trâm	NaN	Đại học Văn Hiến\nCao đẳng/Đại học\n1 tháng 10...	Độc thân	NaN
3	None	100022189274901	Trần Diệu Hiền	Nữ\nGiới tính	Trường Đại học Mở TP. HCM - Tư vấn tuyển sinh\...	Độc thân	NaN
4	None	100091124224730	Hân Lê	NaN	NaN	Lucie Tranová\nĐã kết hôn	Thành phố Hồ Chí Minh\nTỉnh/Thành phố hiện tại
...	...	...	...	...	...	...	...
61	None	61550611575810	The Lyricalogy	Ngày sinh	NaN	NaN	NaN
62	None	100090007038592	Quoc Anh Ng	NaN	NaN	Độc thân	NaN
63	None	100041235171670	L'une	Ngày sinh	NaN	NaN	NaN
64	None	100024095117283	Đặng Bạch Phát	Nam\nGiới tính	Đại học Bách khoa Hà Nội - Hanoi University of...	Độc thân	Hà Nội\nTỉnh/Thành phố hiện tại\nHà Nội\nQuốc gia
65	None	100071394758382	Xóm Xi Xăm	Ngày sinh	NaN	NaN	NaN

## Những người độc thân

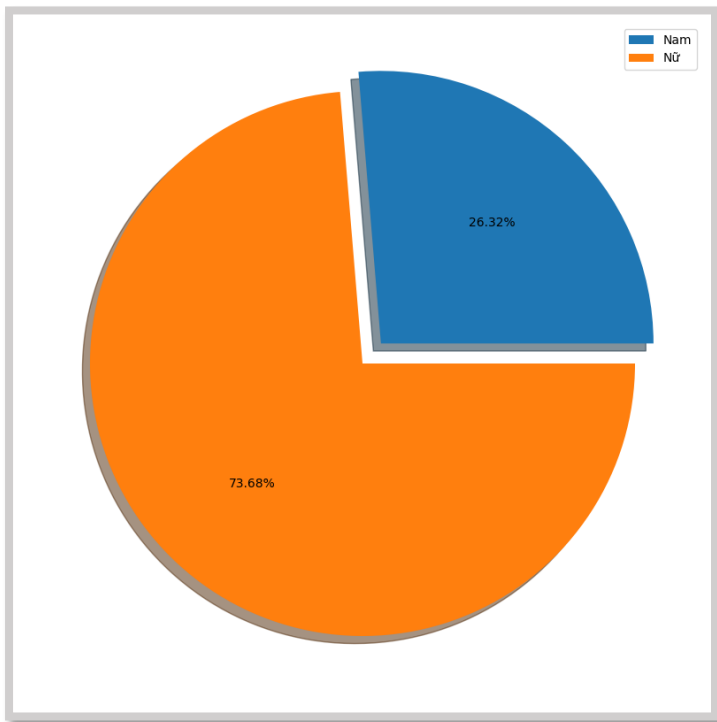


Phần lớn những commenter đều là người độc thân, khoảng 25%

Tỷ lệ số người độc thân chiếm 0.2424242424242424

- ⇒ Bài viết có thể chứa nội dung hoặc chủ đề mà người độc thân thường quan tâm hoặc muốn thảo luận
- ⇒ Fanpage có thể đã triển khai một chiến lược tương tác đặc biệt để kích thích sự tham gia từ nhóm đối tượng này

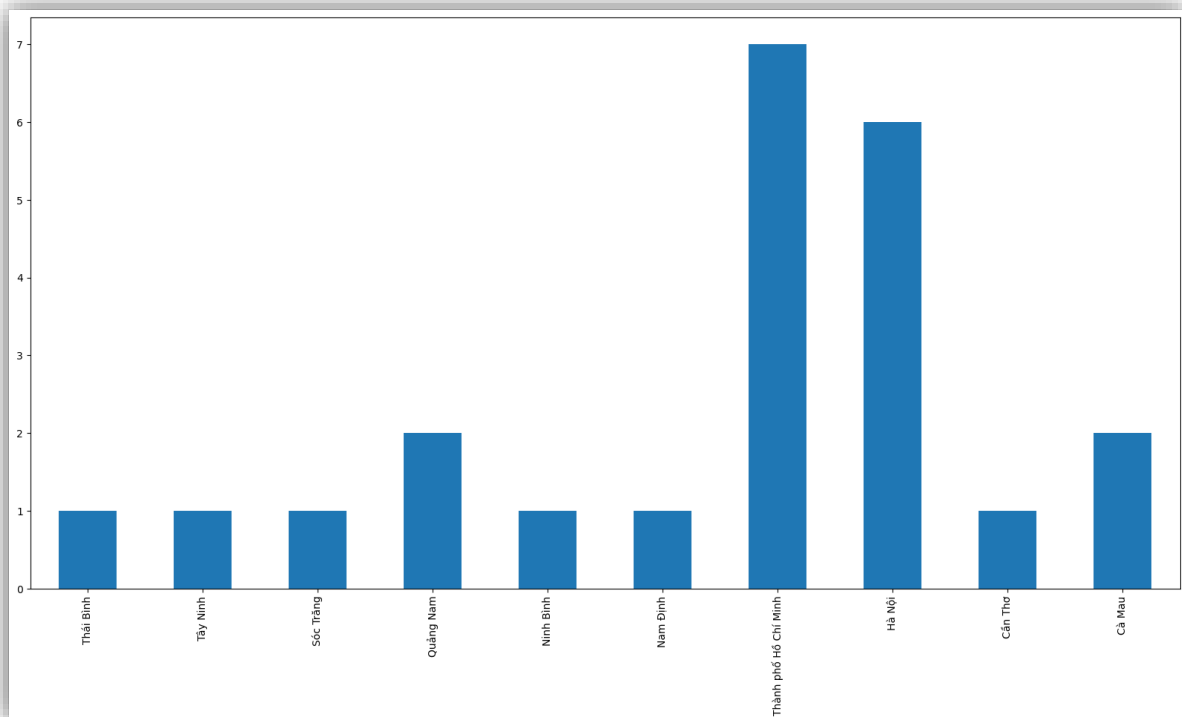
## Tỷ lệ giới tính



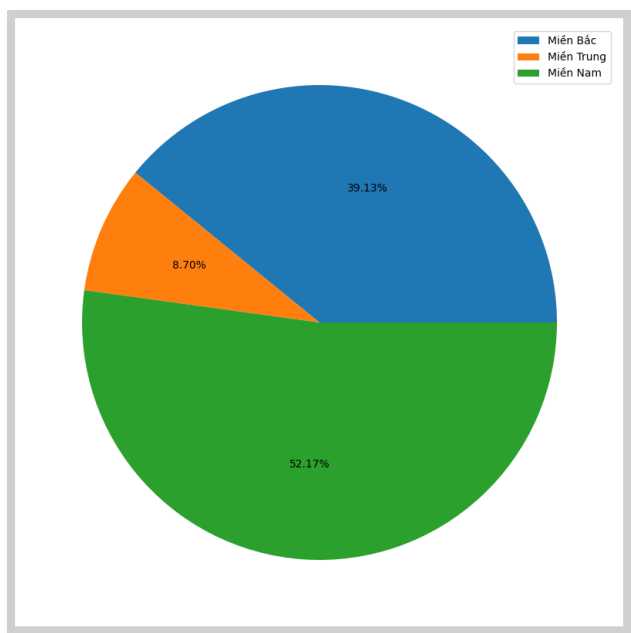
Những người comment chủ yếu là nam, số lượng gấp 3 lần nữ giới

⇒ Nội dung bài viết phù hợp với đối tượng nam

## Nơi ở



- Commenter đến từ chủ yếu 2 thành phố lớn của Việt Nam là TP.HCM và TP. Hà Nội



- Miền Nam chiếm chủ yếu ~ 2 miền Bắc và Trung gộp lại