# Propaganda data investigation

Vladyslav Bezborodov CS-3
Id name 24
Computational Social Science
Teacher: Andrew Kurochkin

5.12.2023

# Agenda

1. Introduction
2. Data acquisition
3. Exploratory data analysis
4. Further work
5. Sources

# Introduction

- A few seconds about dataset
- What I tried investigating
- Presenting my results

# Data acquisition

- Used a propaganda dataset by Kate Burovova
- ~300 of russian telegram propaganda channels
  (just needed to join them in one file, used ipynb file for merging from previous homework)
- Obtaining time : 30-60 min
- Possible problems : had to adapt ipynb file for current task

- **Data statistics :**
  - Posts quantity : 8108693
  - Dataset size : 6670 Megabytes

Brief look into the data

```
In [4]: df.head(10)
```

| | id | date | views | reactions | to_id | fwd_from | message | type | duration | channel_id |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 12602.0 | 2022-12-19 13:05:23+00:00 | 3645.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🏴‍ Исламабад сделал ставку на афганских тали... | photo | NaN | Abbasdjuma |
| 1 | 12601.0 | 2022-12-19 09:52:21+00:00 | 5831.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🕯 6 лет назад, 19 декабря 2016 года, в резуль... | photo | NaN | Abbasdjuma |
| 2 | 12600.0 | 2022-12-19 09:18:53+00:00 | 3944.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🇮🇷 Глава МИД Ирана Хосейн Амир Абдоллахиян с... | photo | NaN | Abbasdjuma |
| 3 | 12599.0 | 2022-12-19 08:32:39+00:00 | 2970.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | MessageFwdHeader(date=datetime.datetime(2022, ... | Наши Друзья открыли \nсбор  для одного из Доне... | photo | NaN | Abbasdjuma |
| 4 | 12598.0 | 2022-12-18 21:41:25+00:00 | 4993.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🇷🇺 Сегодня, 19 декабря в России празднуют День... | photo | NaN | Abbasdjuma |
| 5 | 12597.0 | 2022-12-18 13:39:35+00:00 | 5713.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🇷🇺❤️\n\n#Дайджест_СМИ \n\n 🇻🇦 ABC: Папа римский... | photo | NaN | Abbasdjuma |
| 6 | 12596.0 | 2022-12-18 08:38:12+00:00 | 6186.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🇷🇸❤️ Российский посол в Сербии Александр Боц... | photo | NaN | Abbasdjuma |
| 7 | 12595.0 | 2022-12-17 14:37:20+00:00 | 6181.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | MessageFwdHeader(date=datetime.datetime(2022, ... | 🇷🇺 Директор Центрального разведывательного у... | photo | NaN | Abbasdjuma |
| 8 | 12594.0 | 2022-12-17 08:31:08+00:00 | 80764.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | NaN | 🇷🇺⚡️⚡️⚡️⚡️\n\nНедавно я писал об освобожд... | video | 103.0 | Abbasdjuma |
| 9 | 12593.0 | 2022-12-17 08:00:35+00:00 | 6889.0 | MessageReactions(results=[ReactionCount(reacti... | PeerChannel(channel_id=1261603870) | MessageFwdHeader(date=datetime.datetime(2022, ... | 🇮🇳 В Индии уверены, что Запад не откажется о... | photo | NaN | Abbasdjuma |

```
In [5]: df.shape
```

```
Out[5]: (8108693, 10)
```

# Data Analysis

Investigated questions

- How has the overall activity (number of posts) in these propaganda channels was evolving before and after the beginning of the war?

```
In [13]: # For our storytelling start we would like to begin with presenting the number of all posts for every channel
         # and then we will gradually deepen into the analysis.

         # Here we have a sorted df of posts over all years. In the next cell we would like to show the posts quantity
         # before and after the war

         pd.crosstab(index = df['channel_id'], columns = 'overall_num_of_posts').sort_values(by='overall_num_of_posts',
                                                                                             ascending= False).head(30)
```
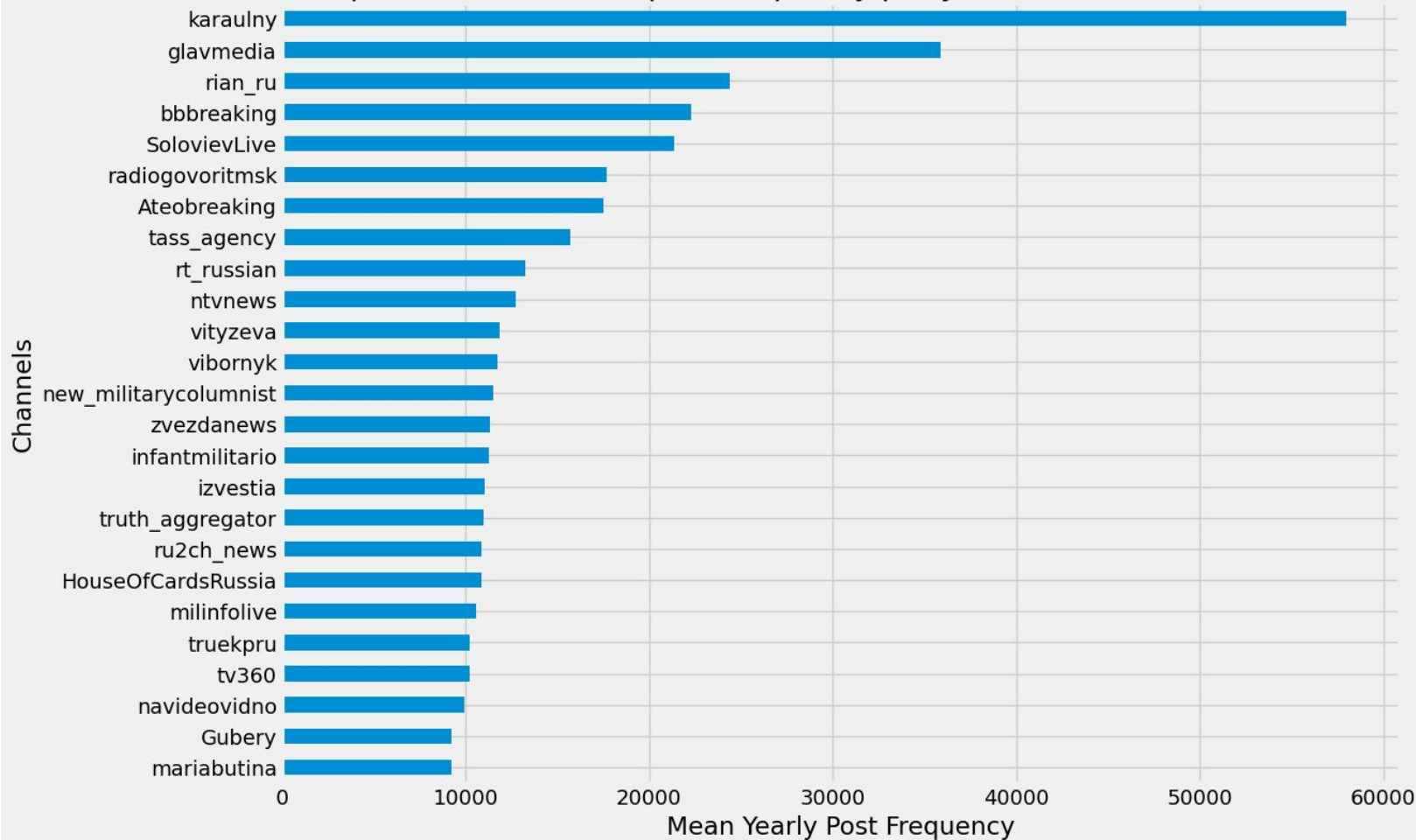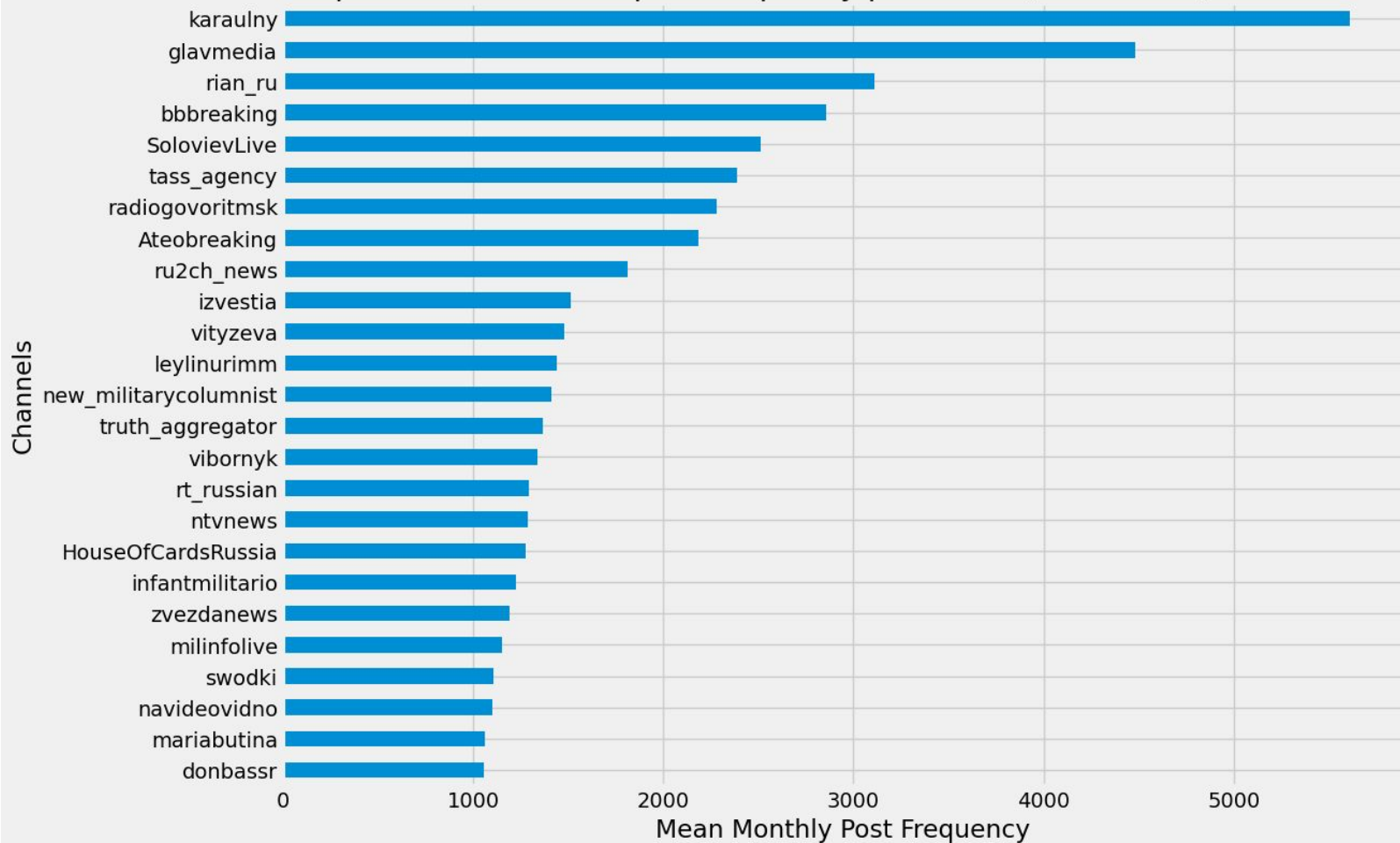
Out[13]:

| col_0 | overall_num_of_posts |
|---|---|
| channel_id | |
| karaulny | 435487 |
| glavmedia | 218510 |
| swodki | 195882 |
| rian_ru | 187755 |
| tass_agency | 170746 |
| SolovievLive | 140116 |
| bbbreaking | 139745 |
| rt_russian | 138568 |
| radiogovoritmsk | 117498 |
| izvestia | 113400 |
| ntvnews | 105117 |
| zvezdanews | 102190 |
| tv360 | 100284 |
| truekpru | 100131 |

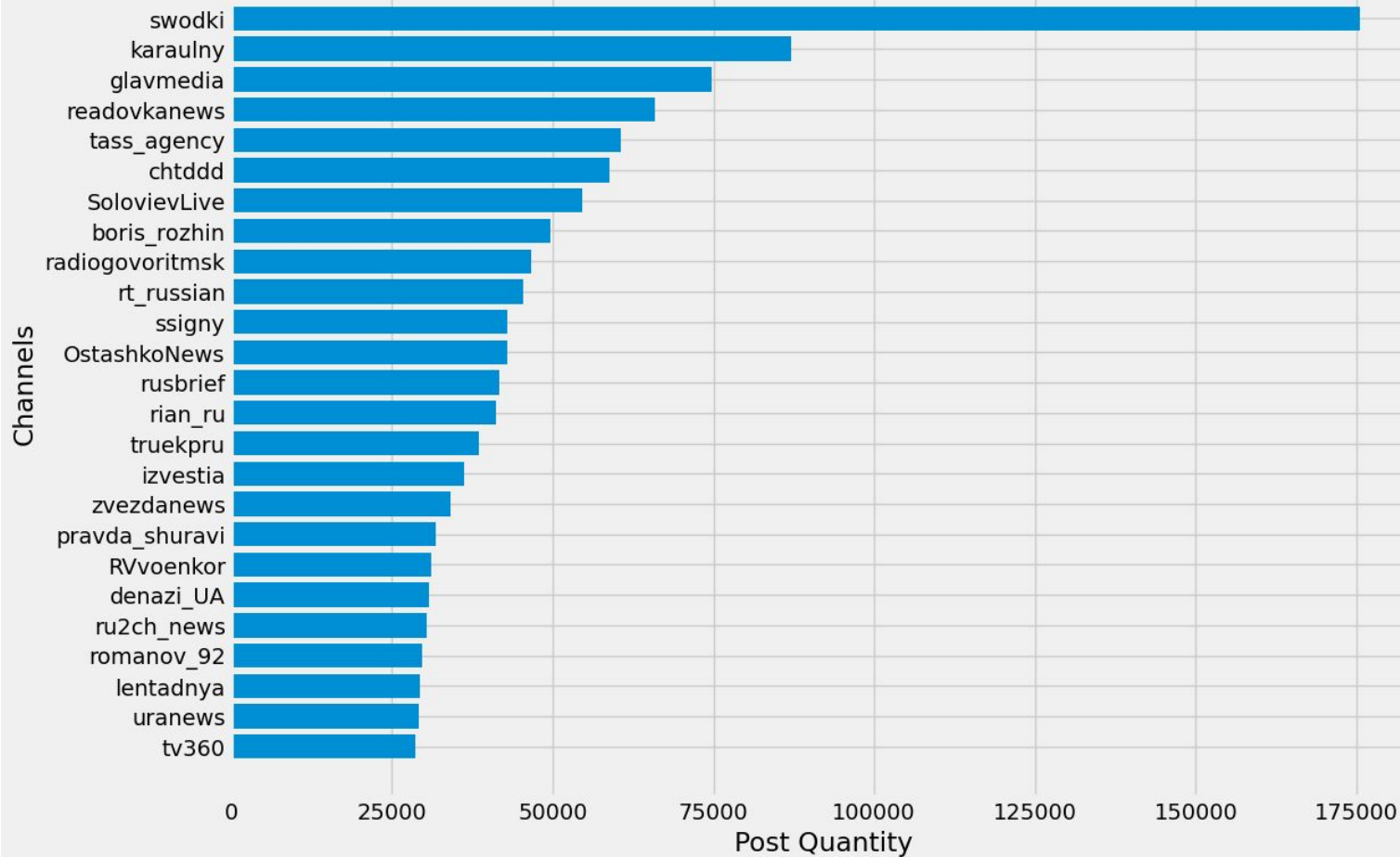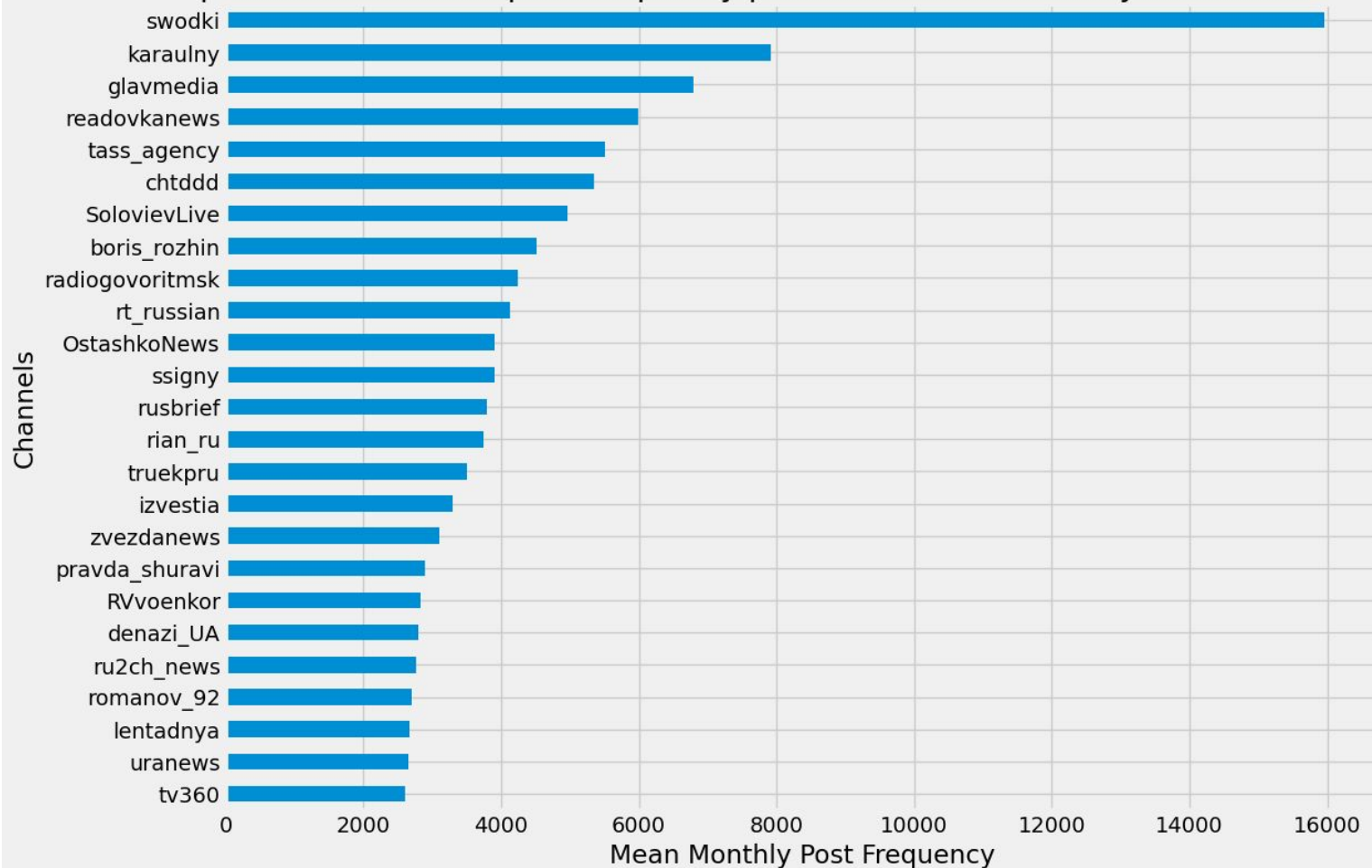Top 25 channel mean post frequency per year(2017-2022) before war

Top 25 channel mean post frequency per month(2017-2022) before war

Channels (top to bottom): karaulny, glavmedia, rian_ru, bbbreaking, SolovievLive, tass_agency, radiogovoritmsk, Ateobreaking, ru2ch_news, izvestia, vityzeva, leylinurimm, new_militarycolumnist, truth_aggregator, vibornyk, rt_russian, ntvnews, HouseOfCardsRussia, infantmilitario, zvezdanews, milinfolive, swodki, navideovidno, mariabutina, donbassr
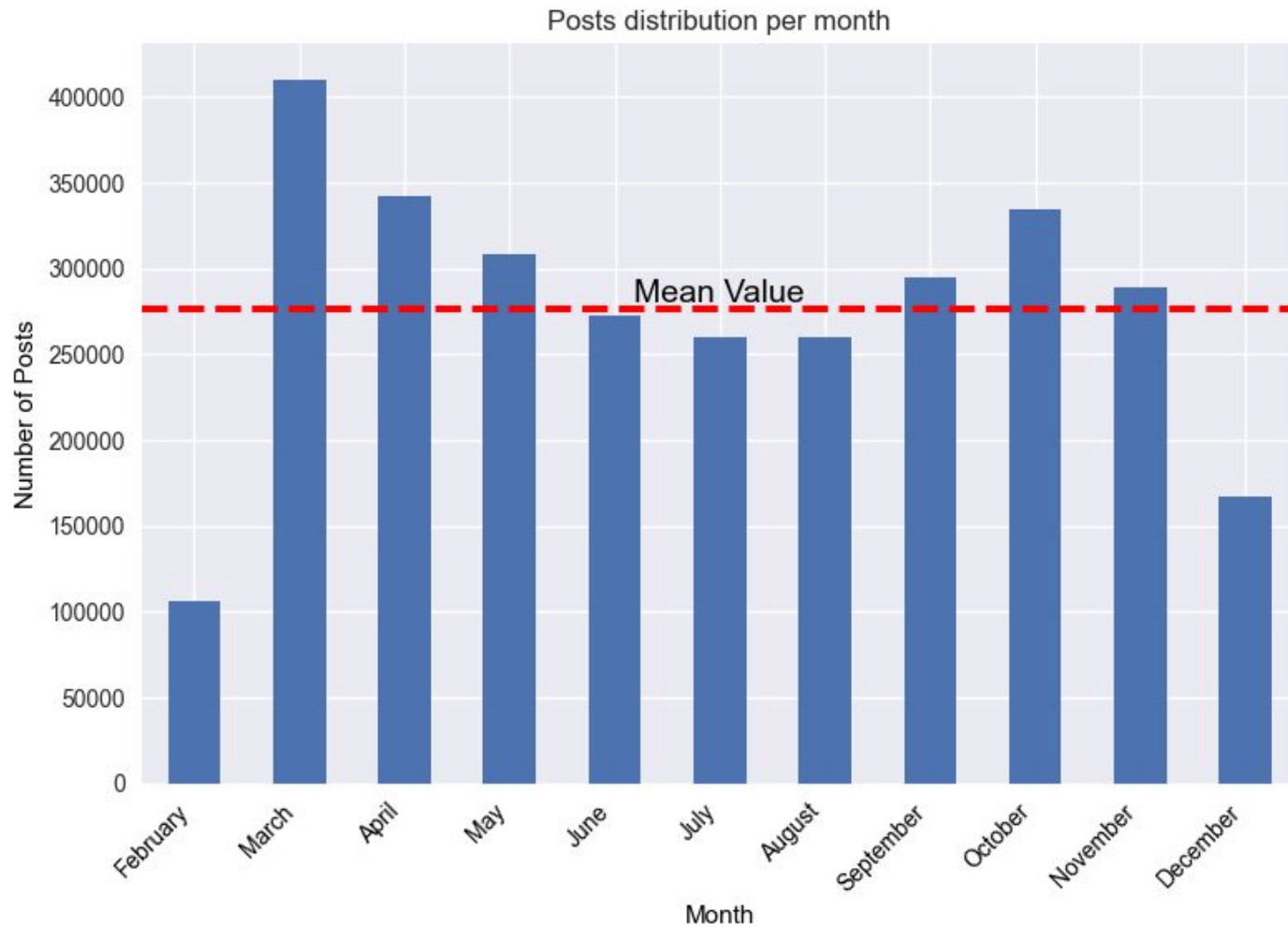
X-axis: Mean Monthly Post Frequency

Top 25 Channel Post Quantity from February to December 2022

Top 25 channel mean post frequency per month from February to December 2022

- What is the distribution of posts across months?

Posts distribution per month

As observed, there are 6 months that surpass the mean threshold value for post activity, with the highest post amounts recorded in March, April, and October, respectively. Notably, these peaks coincide with significant events that likely influenced the surge in post activity.
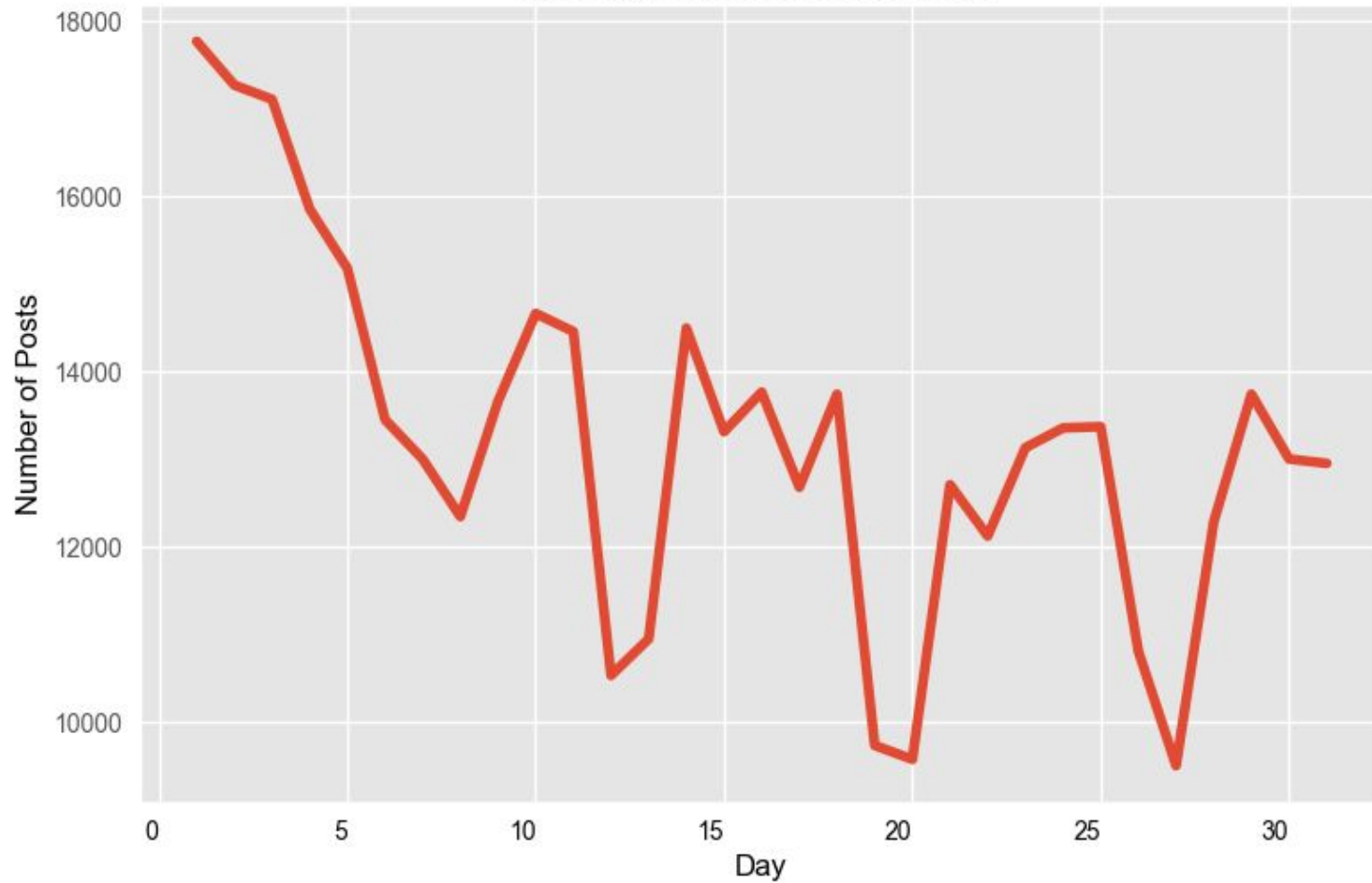
In March, there was heightened activity due to the commencement of warfare, particularly in the Kyiv region, with notable fights in Bucha, Irpin, and Hostomel, which may explain the increased posts during this period.

Continuing into April and May, conflicts persisted in the East and South of Ukraine, culminating in the extended battle in Mariupol, which garnered substantial attention and continued through May. The conclusion of this conflict, particularly impactful for russia, likely contributed to a surge in posts during this timeframe.

Moreover, in September were counter-offensive operation in the Kharkiv region and russia announced mobilization, a significant event that could have provoked the high post flow observed during that month. In October was first explosion on the Kerch bridge and massive air missiles atacks from russia side. And in November Ukraine managed to get the right-bank occupied districts of the Kherson region back.

- Let's take a closer look at March as it stands out for high news flow.
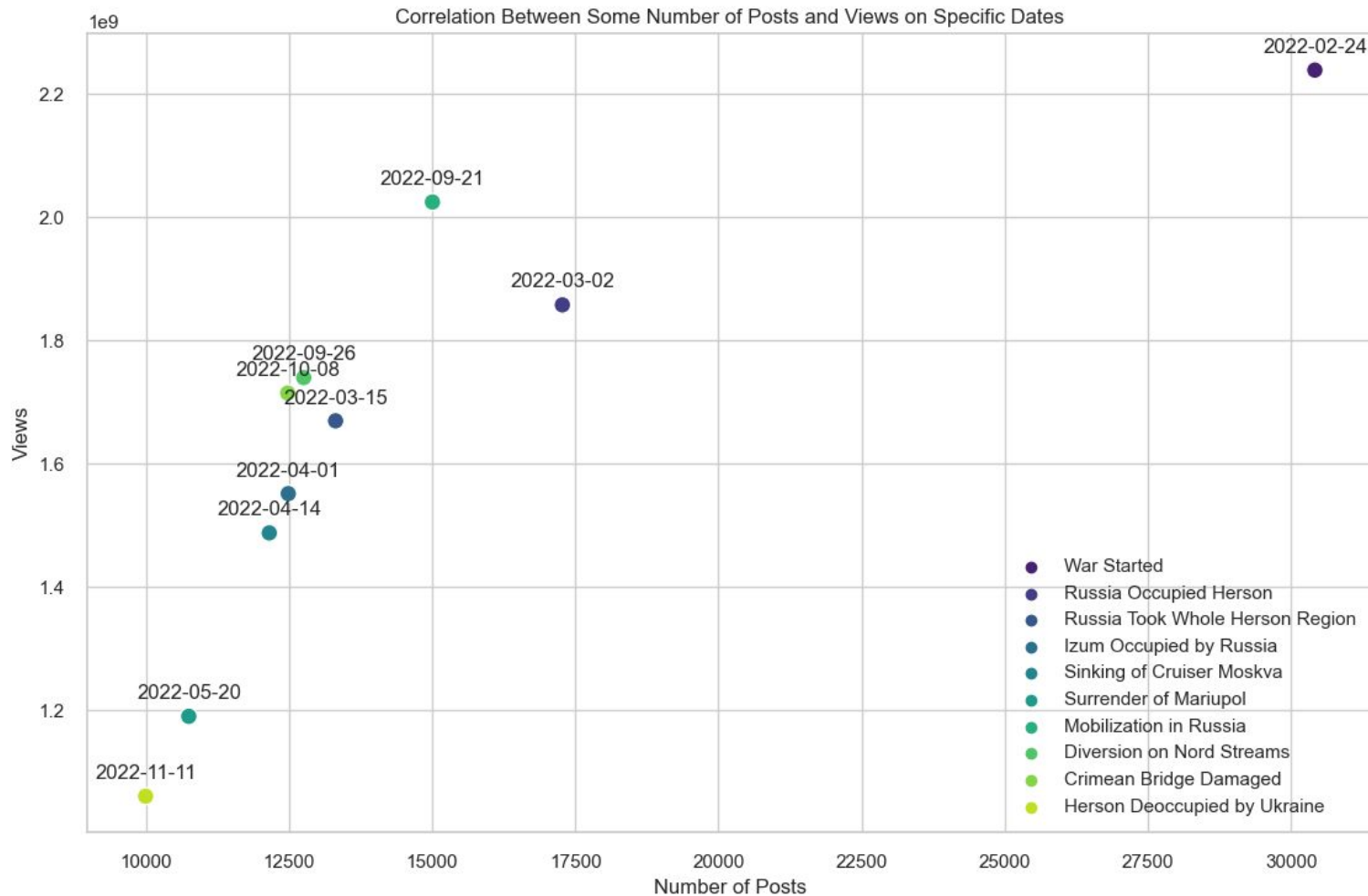
Message distribution in March

The high news flow was only at the beginning of the month as in this days might have happened most essential events for russian army.
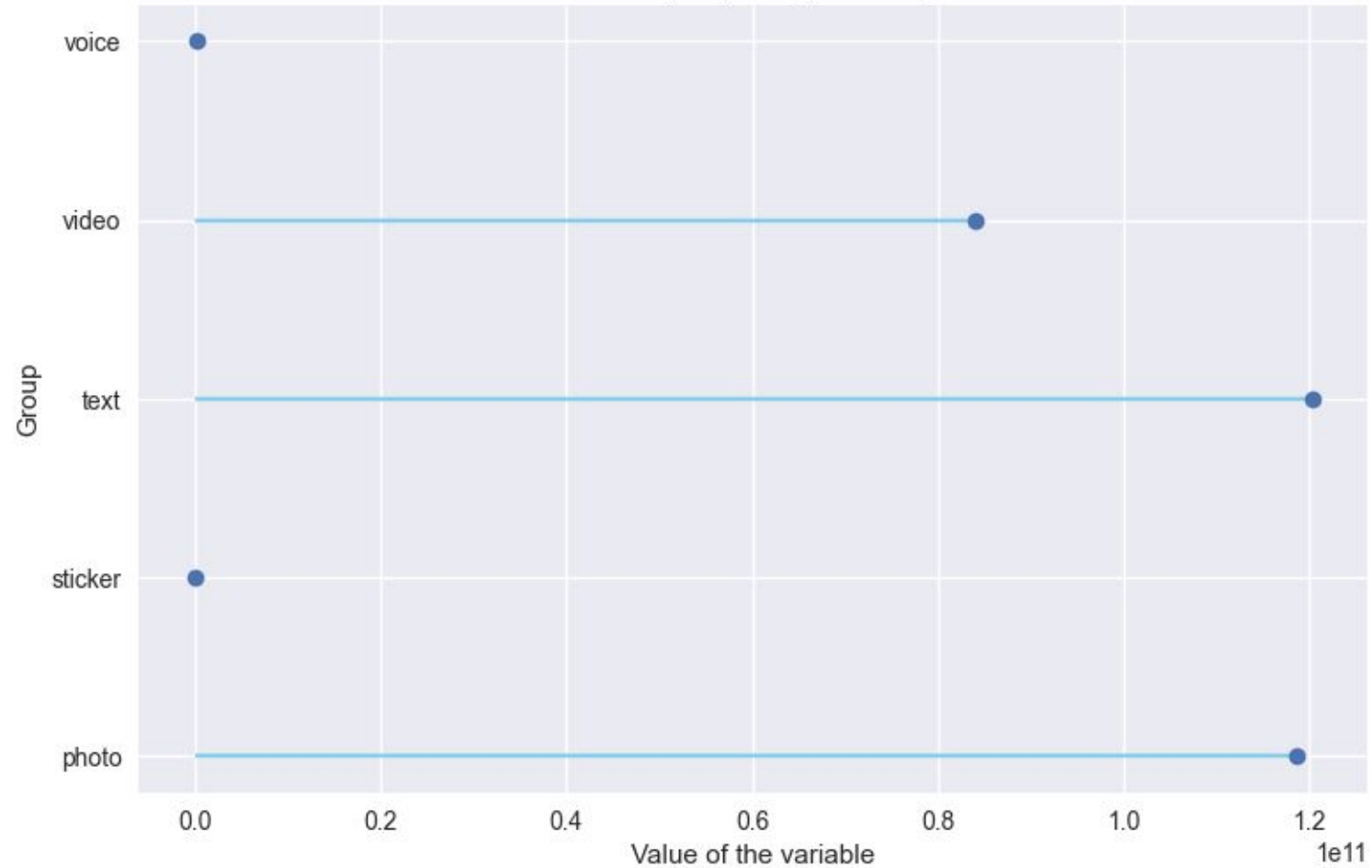
According to the official sources, by the first days of March Russia has taken vast territories and was fighting in main directions toward Kyiv and south region.

- Can we identify any significant events or spikes in post activity based on specific dates?
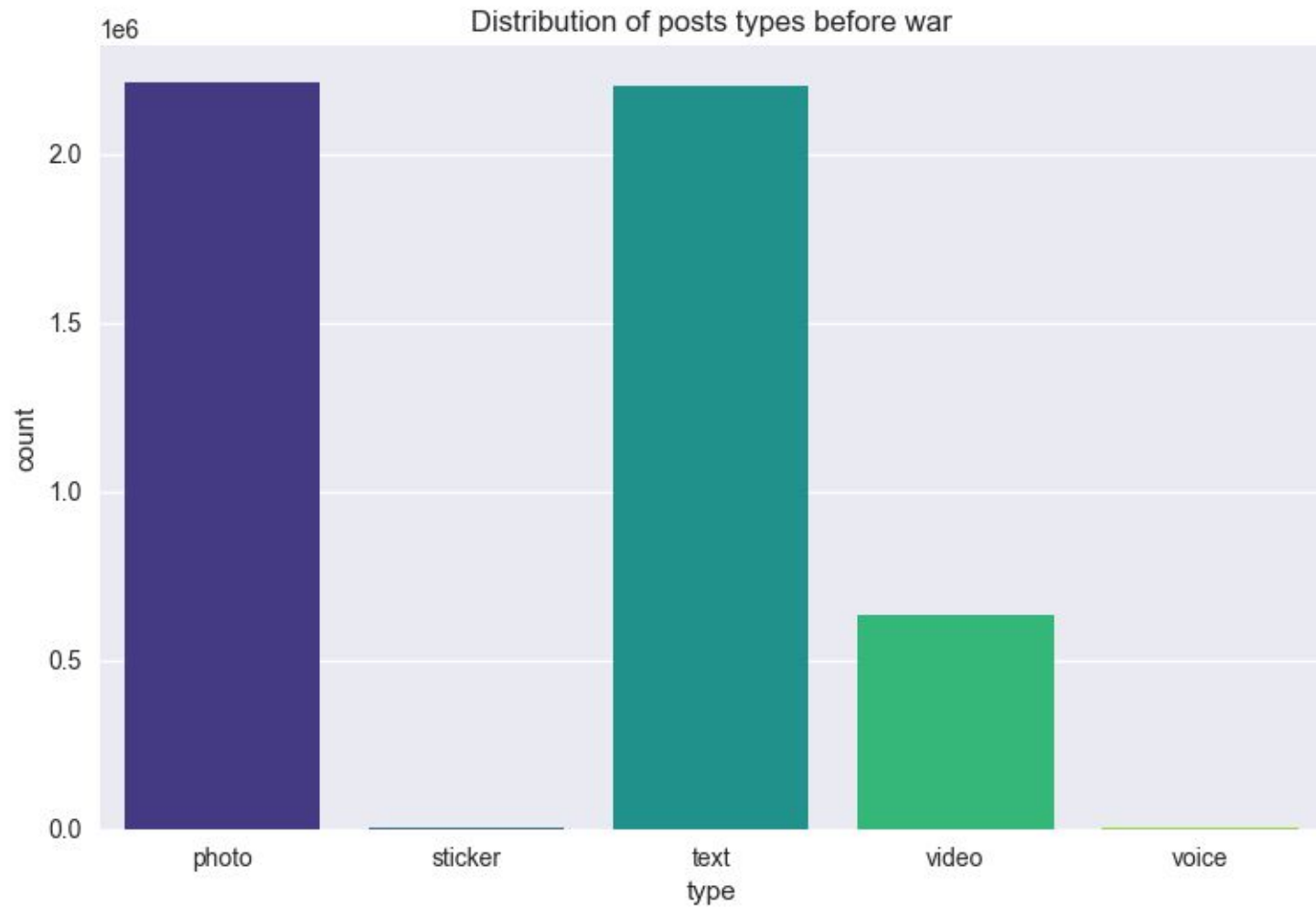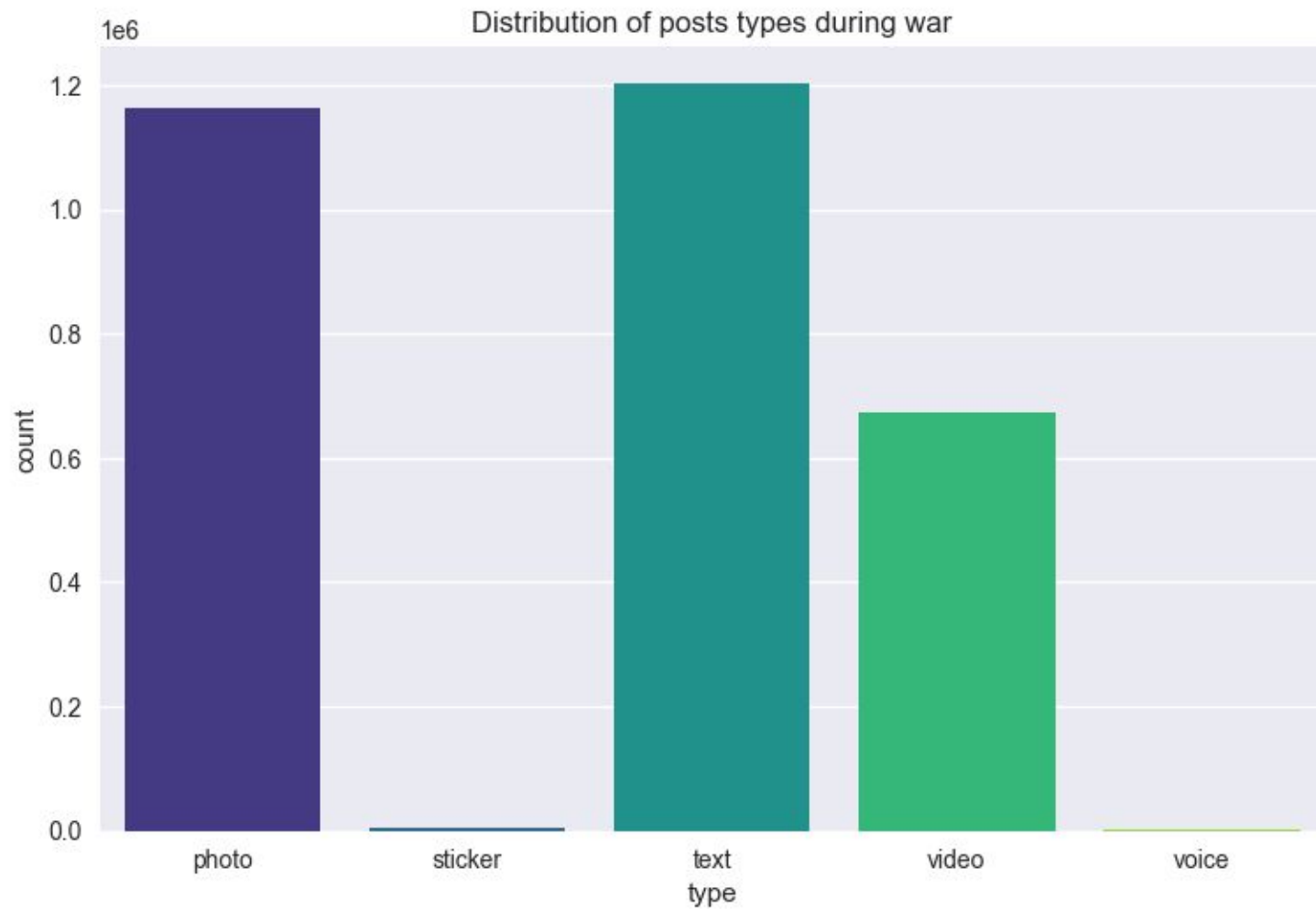
Correlation Between Some Number of Posts and Views on Specific Dates

- What is the distribution of views per post types?

Total views per post type during war

- What is the distribution of posts types?

Distribution of posts types before war
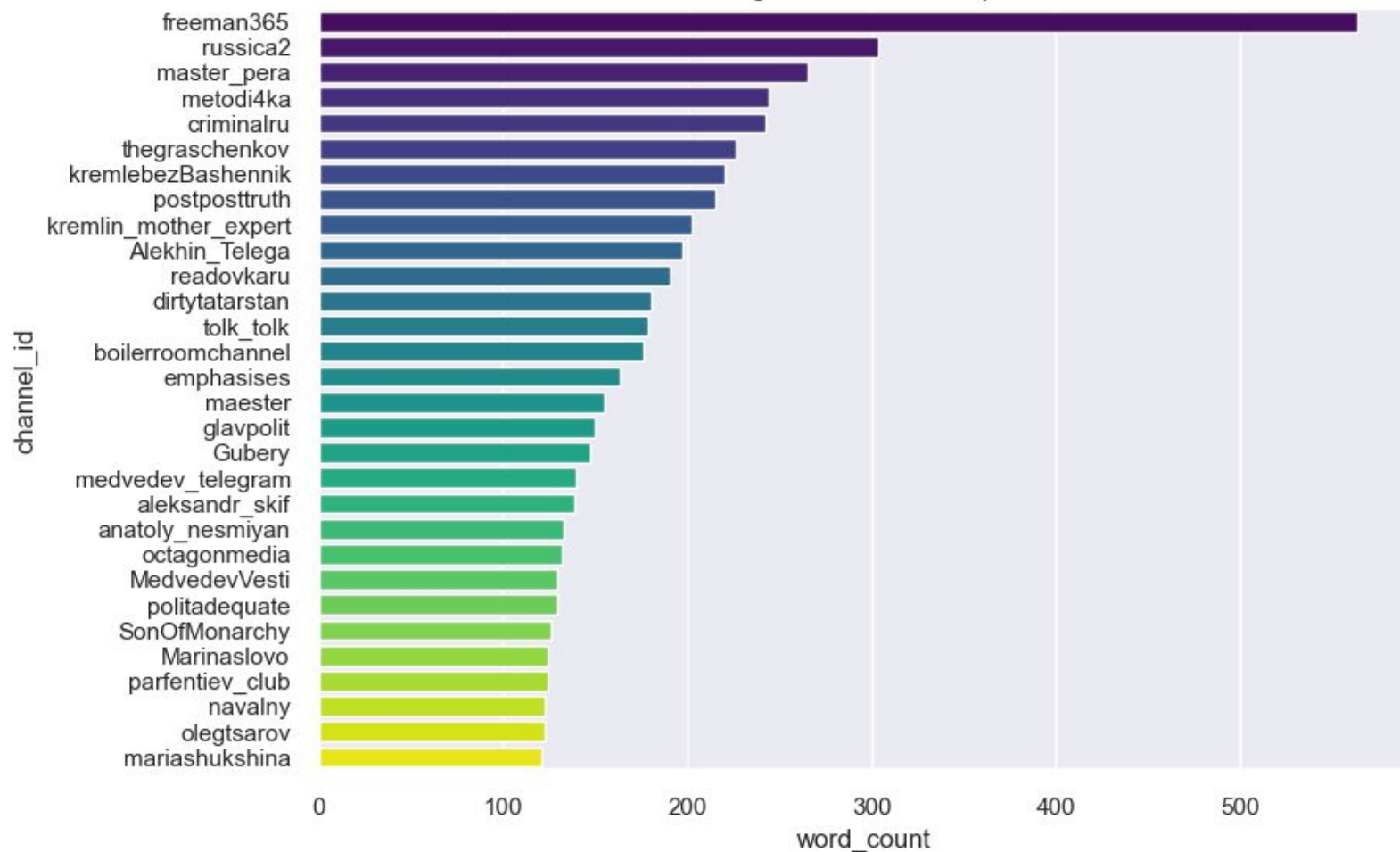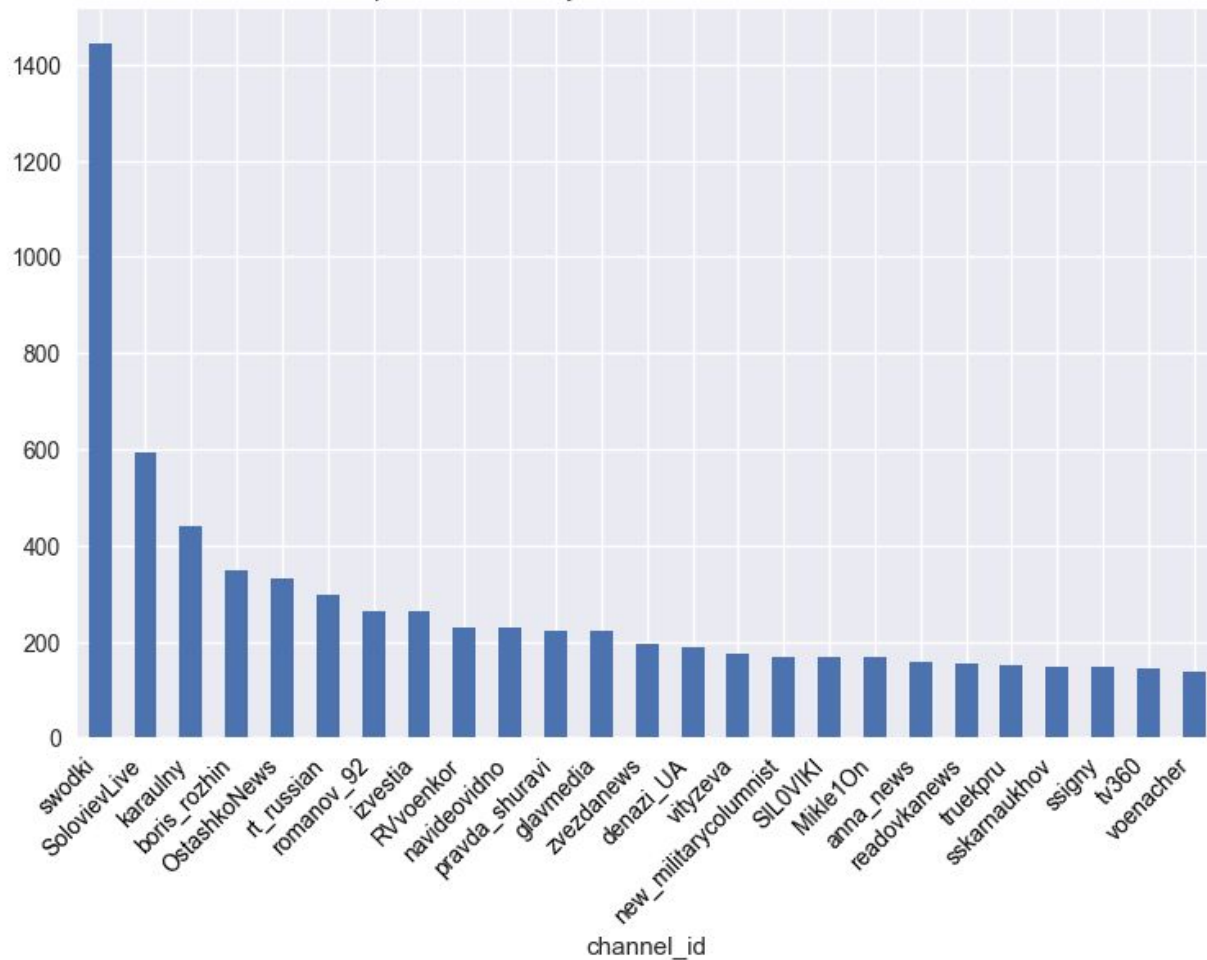
Distribution of posts types during war

- What is the average post size of channels posts?

Average size of channels posts

- What is total duration of all videos per channel?

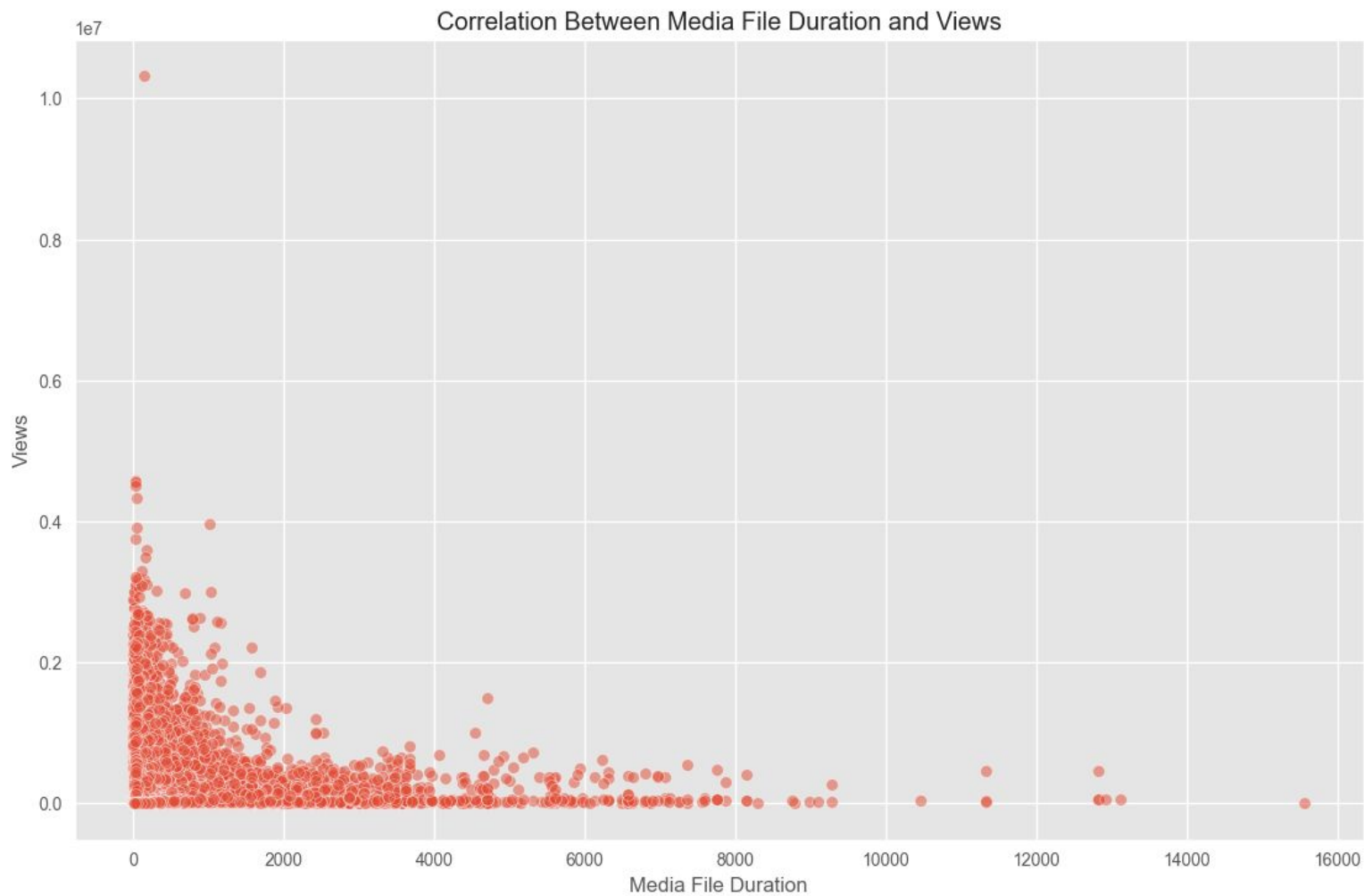top 25 channels by duration of all videos in hours

- How the file duration(audio video) might influence engagement(views)?
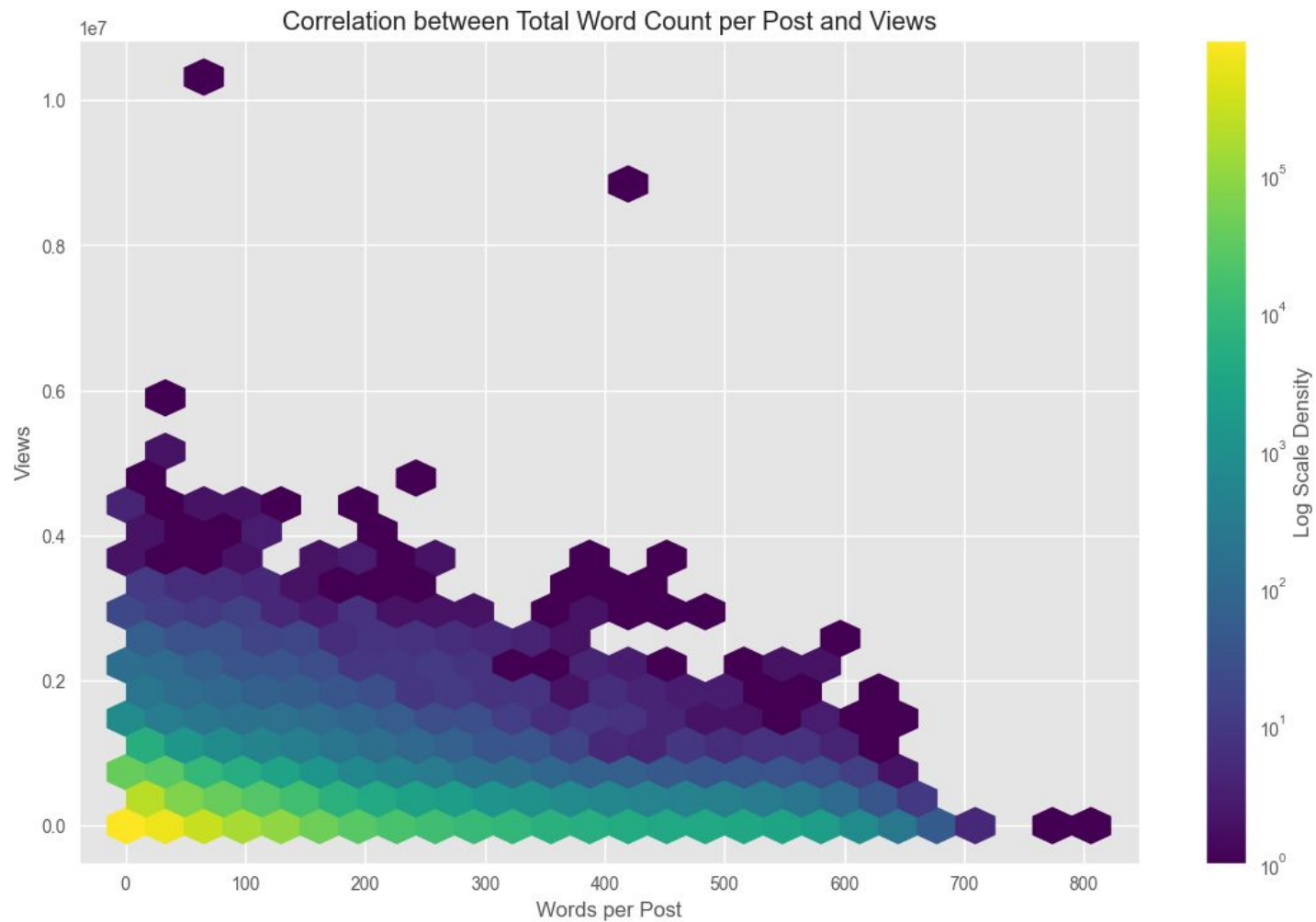
  reactions under posts are not included because 66% of all reactions are lost(NaN values)

  ```
  1  df.isna().sum()
  id                0
  date              0
  views        204614
  reactions   6207754
  to_id             0
  fwd_from    6088700
  message     1092428
  type              0
  duration    6791376
  channel_id        0
  dtype: int64
  ```

  here is a screenshot that shows distribution of lost data
  in columns of dataframe

Correlation Between Media File Duration and Views

- How post size might influence engagement?

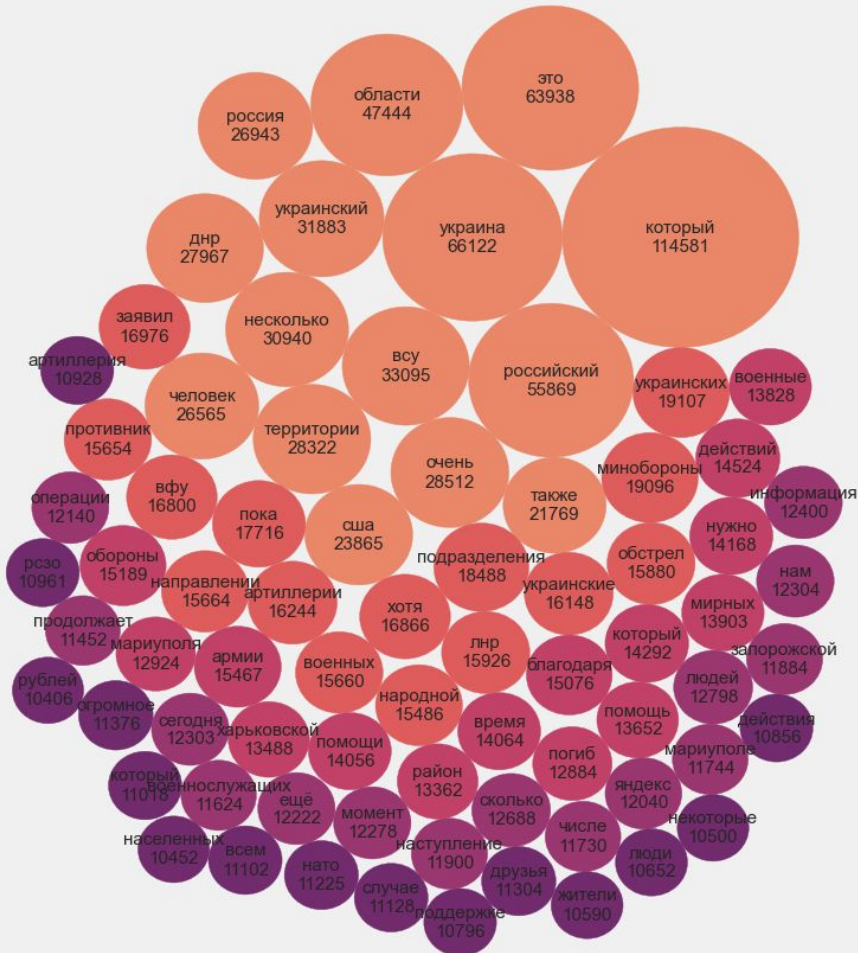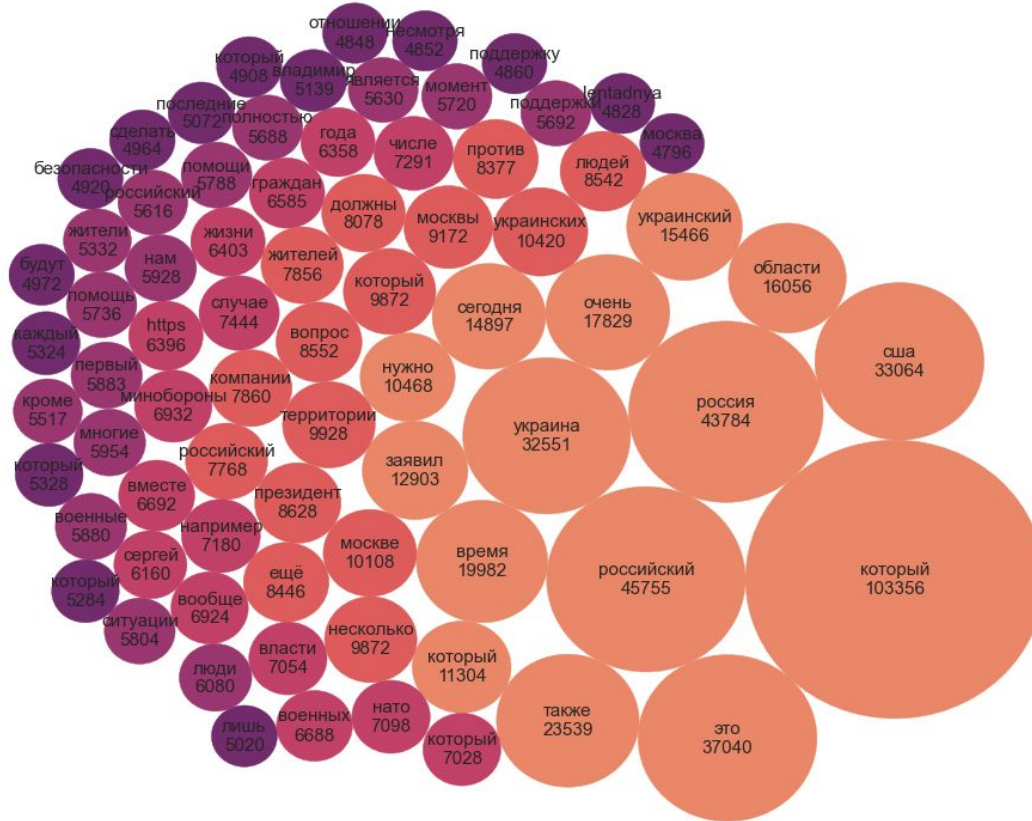Correlation between Total Word Count per Post and Views

- What are the most frequently used words in the posts throughout the warfare?

  Let's take the most active channels, swodki, karaulny and glavmedia, and based on their posts investigate the question.

Swodki channel top 75 frequent words
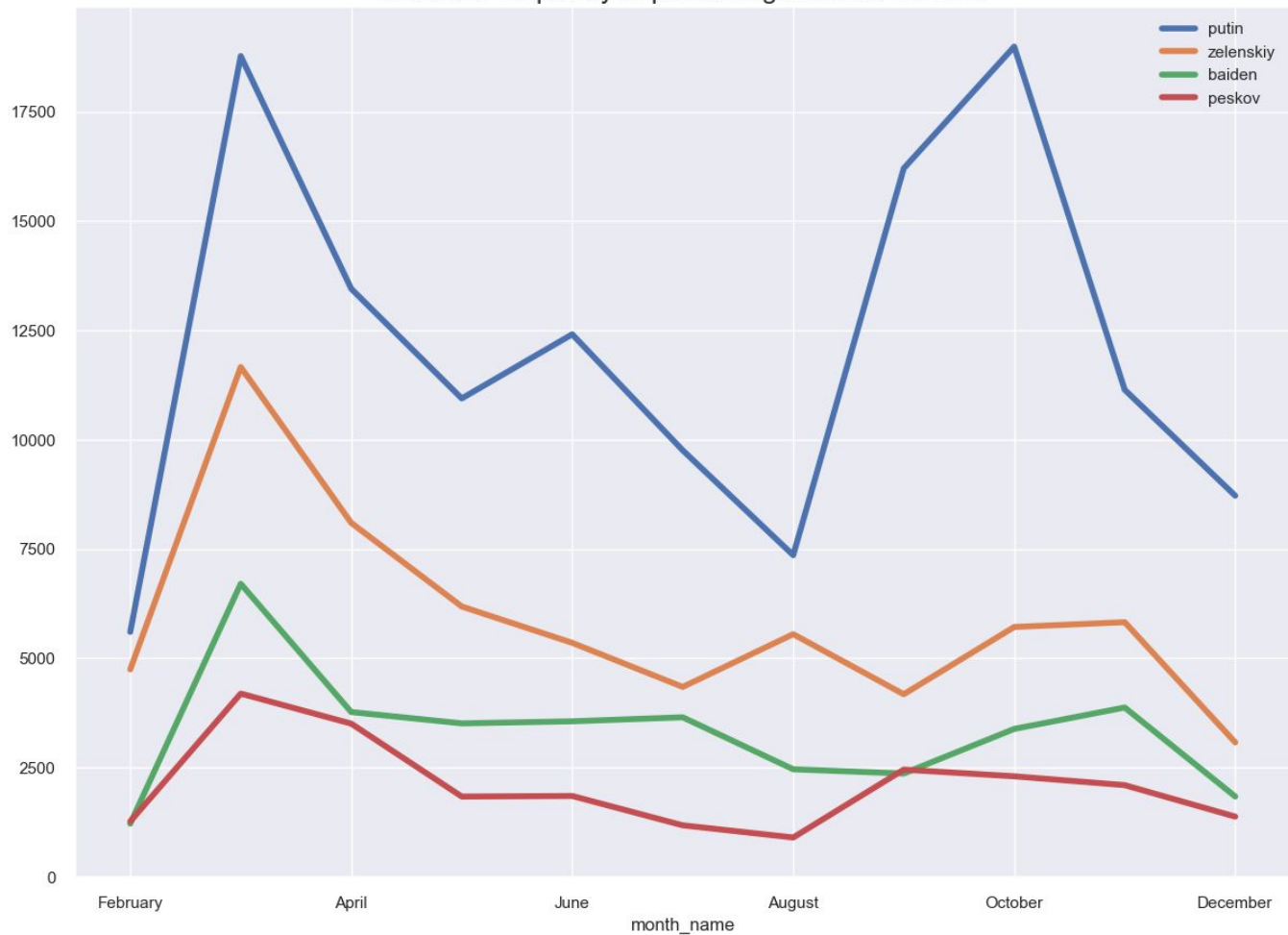
Karaulny channel top 75 frequent words
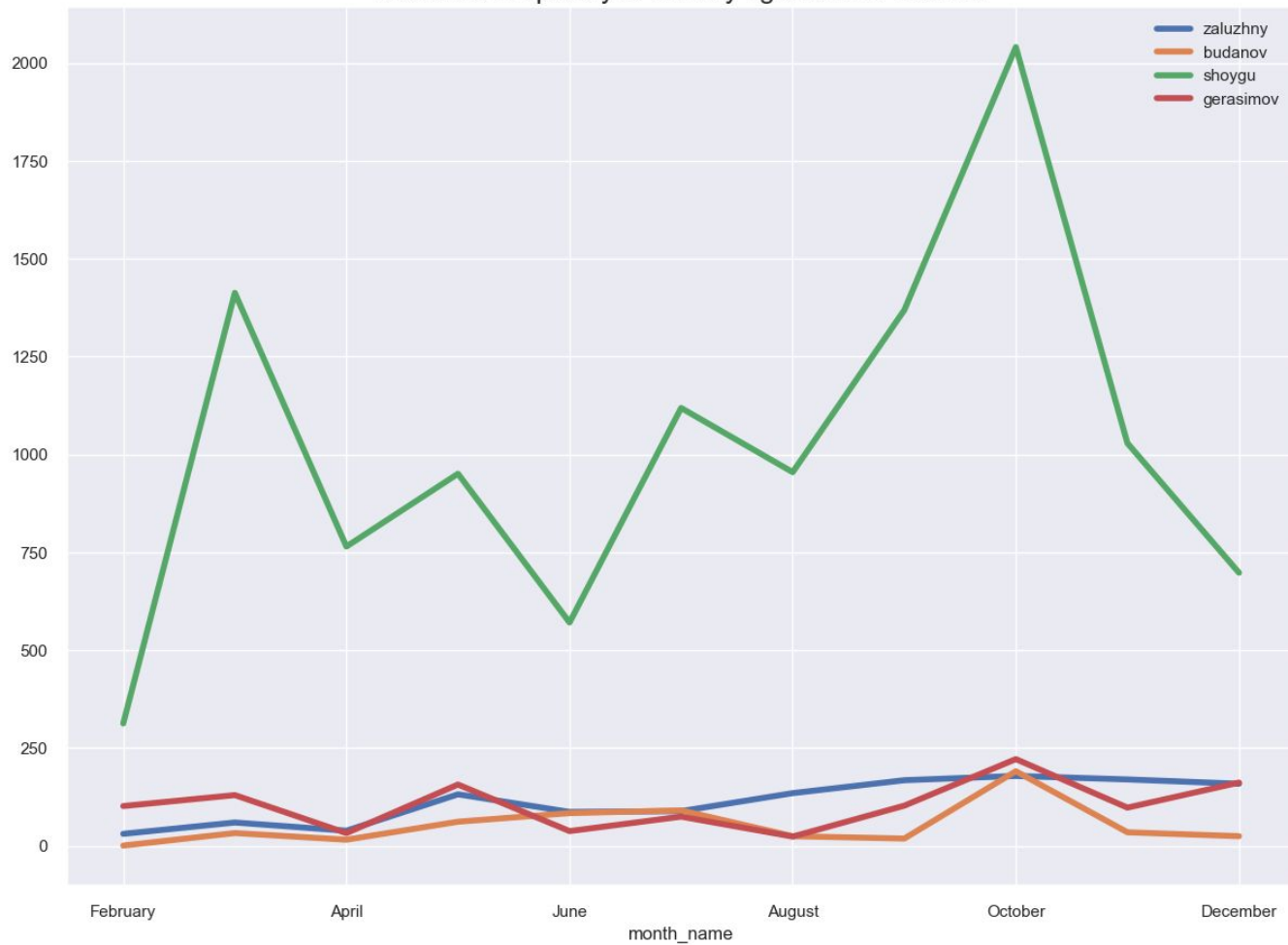
Glavmedia channel top 75 frequent words

- What is the dynamic change of some words, such as Ukrainian cities, political or military figures in whole dataset?
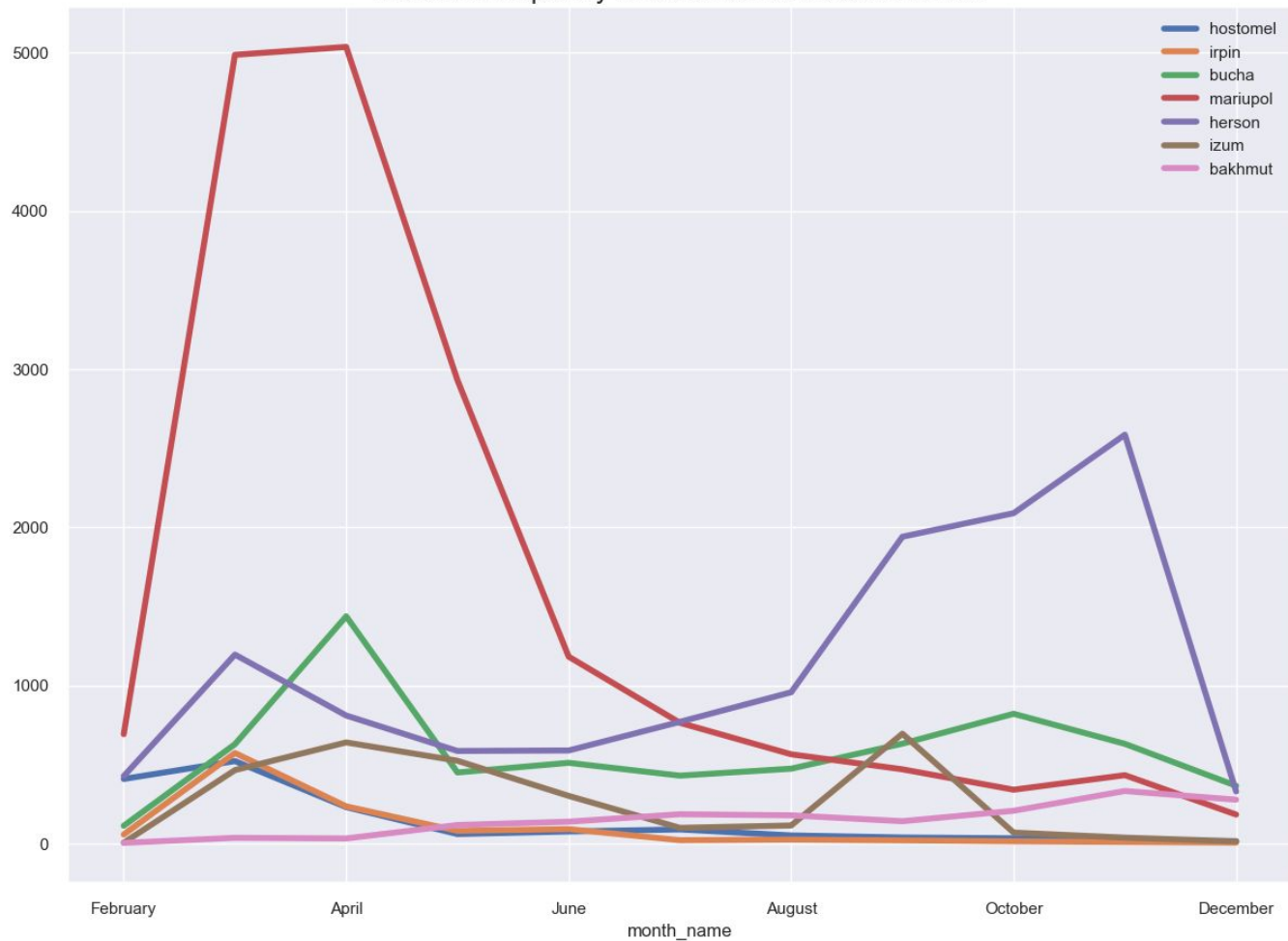
Mentions frequency of political figures over months

Mentions frequency of military figures over months

Mentions frequency of Ukrainian cities over months

# Further work

I have minimum three ideas what to explore in the dataset:

- Do successes or failures of russians provoke more reactions in posts?
- Co-Occurrence Plot Analysis
- Find a distribution of posts types per channel

# Sources

- Link to github repository :

  https://github.com/vladyslavBrothervinn/russian-propaganda-data-investigation

- Other links :

  https://python-graph-gallery.com/

  https://www.data-to-viz.com/

  https://seaborn.pydata.org/index.html

Thank you for your attention!