

```
#Preprocessing, EDA and Evaluation

!pip install transformers
!pip install datasets

Looking in indexes: https://pypi.org/simple, https://us-
python.pkg.dev/colab-wheels/public/simple/
Collecting transformers
  Downloading transformers-4.28.1-py3-none-any.whl (7.0 MB)
  7.0/7.0 MB 75.0 MB/s eta
0:00:00
    requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.9/dist-
    packages (from transformers) (4.65.0)
    Requirement already satisfied: requests in
    /usr/local/lib/python3.9/dist-packages (from transformers) (2.27.1)
Collecting huggingface-hub<1.0,>=0.11.0
  Downloading huggingface_hub-0.13.4-py3-none-any.whl (200 kB)
  200.1/200.1 kB 22.6 MB/s eta
0:00:00
    requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.9/dist-
    packages (from transformers) (6.0)
    Requirement already satisfied: numpy>=1.17 in
    /usr/local/lib/python3.9/dist-packages (from transformers) (1.22.4)
    Requirement already satisfied: regex!=2019.12.17 in
    /usr/local/lib/python3.9/dist-packages (from transformers)
    (2022.10.31)
    Requirement already satisfied: filelock in
    /usr/local/lib/python3.9/dist-packages (from transformers) (3.11.0)
    Requirement already satisfied: packaging>=20.0 in
    /usr/local/lib/python3.9/dist-packages (from transformers) (23.0)
Collecting tokenizers!=0.11.3,<0.14,>=0.11.1
  Downloading tokenizers-0.13.3-cp39-cp39-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl (7.8 MB)
  7.8/7.8 MB 91.6 MB/s eta
0:00:00
    requirement already satisfied: typing-extensions>=3.7.4.3 in
    /usr/local/lib/python3.9/dist-packages (from huggingface-
    hub<1.0,>=0.11.0->transformers) (4.5.0)
    Requirement already satisfied: certifi>=2017.4.17 in
    /usr/local/lib/python3.9/dist-packages (from requests->transformers)
    (2022.12.7)
    Requirement already satisfied: idna<4,>=2.5 in
    /usr/local/lib/python3.9/dist-packages (from requests->transformers)
    (3.4)
    Requirement already satisfied: urllib3<1.27,>=1.21.1 in
    /usr/local/lib/python3.9/dist-packages (from requests->transformers)
    (1.26.15)
    Requirement already satisfied: charset-normalizer~=2.0.0 in
    /usr/local/lib/python3.9/dist-packages (from requests->transformers)
    (2.0.12)
```

```
Installing collected packages: tokenizers, huggingface-hub,
transformers
Successfully installed huggingface-hub-0.13.4 tokenizers-0.13.3
transformers-4.28.1
Looking in indexes: https://pypi.org/simple, https://us-
python.pkg.dev/colab-wheels/public/simple/
Collecting datasets
    Downloading datasets-2.11.0-py3-none-any.whl (468 kB)
    ━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 468.7/468.7 kB 30.6 MB/s eta
0:00:00
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.9/dist-
packages (from datasets) (1.22.4)
Collecting dill<0.3.7,>=0.3.0
    Downloading dill-0.3.6-py3-none-any.whl (110 kB)
    ━━━━━━━━━━━━━━━━━━━━━━━━━━ 110.5/110.5 kB 14.8 MB/s eta
0:00:00
Requirement already satisfied: tqdm>=4.62.1 in /usr/local/lib/python3.9/dist-
packages (from datasets) (4.65.0)
Requirement already satisfied: fsspec[http]>=2021.11.1 in
/usr/local/lib/python3.9/dist-packages (from datasets) (2023.4.0)
Collecting multiprocess
    Downloading multiprocess-0.70.14-py39-none-any.whl (132 kB)
    ━━━━━━━━━━━━━━━━━━━━━━━━ 132.9/132.9 kB 17.9 MB/s eta
0:00:00
Requirement already satisfied: pyarrow>=8.0.0 in
/usr/local/lib/python3.9/dist-packages (from datasets) (9.0.0)
Requirement already satisfied: huggingface-hub<1.0.0,>=0.11.0 in
/usr/local/lib/python3.9/dist-packages (from datasets) (0.13.4)
Requirement already satisfied: packaging in
/usr/local/lib/python3.9/dist-packages (from datasets) (23.0)
Requirement already satisfied: pandas in
/usr/local/lib/python3.9/dist-packages (from datasets) (1.5.3)
Collecting responses<0.19
    Downloading responses-0.18.0-py3-none-any.whl (38 kB)
Requirement already satisfied: pyyaml>=5.1 in
/usr/local/lib/python3.9/dist-packages (from datasets) (6.0)
Collecting aiohttp
    Downloading aiohttp-3.8.4-cp39-cp39-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl (1.0 MB)
    ━━━━━━━━━━━━━━━━━━━━━━━━ 1.0/1.0 MB 75.0 MB/s eta
0:00:00
Requirement already satisfied: requests>=2.19.0 in
/usr/local/lib/python3.9/dist-packages (from datasets) (2.27.1)
Collecting xxhash
    Downloading xxhash-3.2.0-cp39-cp39-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl (212 kB)
    ━━━━━━━━━━━━━━━━━━━━━━━━ 212.2/212.2 kB 25.2 MB/s eta
0:00:00
Requirement already satisfied: charset-normalizer<4.0,>=2.0 in
/usr/local/lib/python3.9/dist-packages (from aiohttp->datasets)
```

```
(2.0.12)
Collecting frozenlist>=1.1.1
    Downloading frozenlist-1.3.3-cp39-cp39-
manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux
2014_x86_64.whl (158 kB) ━━━━━━━━━━━━━━━━ 158.8/158.8 kB 21.3 MB/s eta
0:00:00
anylinux_2_17_x86_64.manylinux2014_x86_64.whl (264 kB) ━━━━━━━━━━━━ 264.6/264.6 kB 24.5 MB/s eta
0:00:00
ultidict<7.0,>=4.5
    Downloading multidict-6.0.4-cp39-cp39-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl (114 kB) ━━━━━━━━ 114.2/114.2 kB 15.3 MB/s eta
0:00:00
ent already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.9/dist-
packages (from aiohttp->datasets) (22.2.0)
Collecting async-timeout<5.0,>=4.0.0a3
    Downloading async_timeout-4.0.2-py3-none-any.whl (5.8 kB)
Requirement already satisfied: filelock in
/usr/local/lib/python3.9/dist-packages (from huggingface-
hub<1.0.0,>=0.11.0->datasets) (3.11.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in
/usr/local/lib/python3.9/dist-packages (from huggingface-
hub<1.0.0,>=0.11.0->datasets) (4.5.0)
Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.9/dist-packages (from requests>=2.19.0-
>datasets) (2022.12.7)
Requirement already satisfied: idna<4,>=2.5 in
/usr/local/lib/python3.9/dist-packages (from requests>=2.19.0-
>datasets) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in
/usr/local/lib/python3.9/dist-packages (from requests>=2.19.0-
>datasets) (1.26.15)
Requirement already satisfied: python-dateutil>=2.8.1 in
/usr/local/lib/python3.9/dist-packages (from pandas->datasets) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in
/usr/local/lib/python3.9/dist-packages (from pandas->datasets)
(2022.7.1)
Requirement already satisfied: six>=1.5 in
/usr/local/lib/python3.9/dist-packages (from python-dateutil>=2.8.1-
>pandas->datasets) (1.16.0)
Installing collected packages: xxhash, multidict, frozenlist, dill,
async-timeout, yarl, responses, multiprocess, aiosignal, aiohttp,
datasets
Successfully installed aiohttp-3.8.4 aiosignal-1.3.1 async-timeout-
4.0.2 datasets-2.11.0 dill-0.3.6 frozenlist-1.3.3 multidict-6.0.4
multiprocess-0.70.14 responses-0.18.0 xxhash-3.2.0 yarl-1.8.2
```

```
from datasets import load_dataset

#Importing the dataset from hugging face
dataset = load_dataset('jeremyf/fanfiction_z')

{"model_id":"f3e4e73d73324b86b270c723286bc84b","version_major":2,"vers
ion_minor":0}

Downloading and preparing dataset json/jeremyf--fanfiction_z to
/root/.cache/huggingface/datasets/jeremyf__json/jeremyf--
fanfiction_z-eela1447c757b7ff/0.0.0/
fe5dd6ea2639a6df622901539cb550cf8797e5a6b2dd7af1cf934bed8e233e6e...

{"model_id":"15051b5c076c42fc93d12322b9794ebd","version_major":2,"vers
ion_minor":0}

{"model_id":"b339064bdac041ef878c8633a8cb7caf","version_major":2,"vers
ion_minor":0}

 {"model_id":"1ee6b5b479d9415f9a564e3269d316c1","version_major":2,"vers
ion_minor":0}

 {"model_id":"2cd37f3097e54fd8balcceaa69bf60ddb","version_major":2,"vers
ion_minor":0}

Dataset json downloaded and prepared to
/root/.cache/huggingface/datasets/jeremyf__json/jeremyf--
fanfiction_z-eela1447c757b7ff/0.0.0/
fe5dd6ea2639a6df622901539cb550cf8797e5a6b2dd7af1cf934bed8e233e6e.
Subsequent calls will reuse this data.

 {"model_id":"dd54e990ede14e08af71b3abaaff6468","version_major":2,"vers
ion_minor":0}

dataset

DatasetDict({
    train: Dataset({
        features: ['story', 'title', 'category'],
        num_rows: 943
    })
})

import pandas as pd
import re

#Fetch the store list
story_list = [[i] for i in list(dataset['train']['story'])]
df = pd.DataFrame(story_list,columns=['story'])

#Extract the titles and the categories of the stories
df['title'] = dataset['train']['title']
df['category'] = dataset['train']['category']
```

```

#Clean up line breaks, extra punctuations etc
df['story_cleaned'] = df['story'].apply(lambda x: x.replace('**',''))
df['story_cleaned'] = df['story_cleaned'].apply(lambda x: x.replace("\n\n", '\n'))
df['story_cleaned'] = df['story_cleaned'].apply(lambda x:
re.sub(r'^([a-z])\1*$', 'WORD', x))
# for i in range(df.shape[0]):
#     title = df.iloc[i]['title'].partition('by')[0]
#     df.iloc[i]['story_cleaned'].replace(title, '')
#     print(df.iloc[i]['story_cleaned'])

df

                           story \
0      just a question\n\nWhat I thought when I first...
1          pain\n\nIt may be sad but it can be better\n\n...
2      Zombie Mayhem\n\nZombie Mayhem\n\nIn 2050, the...
3      Zero Wing 2: AYBASBTU Game Script\n\nZERO WING...
4      Not prosaic\n\nDisclaimer, I do not claim to o...
..                                 ...
938    More\n\n**I've watched this movie 3 times this...
939    Change Comes From Within\n\n**As usual, I own ...
940    Changes\n\nYou know that feeling where people ...
941    Undying Hero\n\nZero percent.\n\nVictory canno...
942    Unlosing Haunted Mansion\n\n**Disclaimer**: I ...

                           title \
0      just a question by multifics123
1          pain by Isaiah Thomas
2      Zombie Mayhem by kelbey342
3      Zero Wing 2: AYBASBTU Game Script by Rodrigo Shin
4          Not prosaic by Fault
..                                 ...
938          More by Lizwontcry
939    Change Comes From Within by Dominican Girl
940          Changes by michieexx3
941          Undying Hero by Arukoir
942    Unlosing Haunted Mansion by Smak64

                           category \
0          Zom-B
1          Zenda
2      Zoombie Blondes
3          Zero Wing
4          Zero Wing
..                                 ...
938    Zack and Miri Make a Porno
939    Zack and Miri Make a Porno
940    Zack and Miri Make a Porno
941          Z.H.P.

```

```
story_cleaned
0 just a question\nWhat I thought when I first r...
1 pain\nIt may be sad but it can be better\nThe ...
2 Zombie Mayhem\nZombie Mayhem\nIn 2050, the mil...
3 Zero Wing 2: AYBASBTU Game Script\nZERO WING 2...
4 Not prosaic\nDisclaimer, I do not claim to own...
...
938 More\nI've watched this movie 3 times this wee...
939 Change Comes From Within\nAs usual, I own noth...
940 Changes\nYou know that feeling where people be...
941 Undying Hero\nZero percent.\nVictory cannot be...
942 Unlosing Haunted Mansion\nDisclaimer: I do not...
```

[943 rows x 4 columns]

EDA

```
from matplotlib import pyplot as plt
import matplotlib

#Find and print different types of categories
df.describe()
df['category'].unique()

array(['Zom-B', 'Zenda', 'Zoombie Blondes', 'Zero Wing', 'Zoom',
       'Zettai Karen Children', 'Zero Day', 'Zombie Prom',
       'Zombie Survival Guide', 'Zeke and Luther', 'Zen', 'Zoop',
       "Zanna, Don't!", 'Zenonia', 'Zetman', 'Z Nation', 'Zatch Bell',
       'Zombieland', 'Zone of The Enders', "Zeke's Pad", 'Zoey 101',
       'Zinda/ज़िंदा', 'Zoo Tycoon', 'Zenon', 'Zombie Fallout',
       'Zombies, Run!', 'Z for Zachariah', 'Zapped', 'Zoombinis',
       'Zombie Powder', 'Zorro', 'Zoey 101, iCarly', 'Zoolander',
       'Zoids',
       'Zevo-3', 'Zero Dark Thirty', 'Zombie-Loan', 'Zodiac P.I.',
       'Zeta Project', 'Zathura', 'Zack and Miri Make a Porno',
       'Z.H.P.'],
      dtype=object)
```

```
#Find out dominant categories
category_count = df['category'].value_counts()
category_count
```

Zoey 101	372
Zoids	150
Zatch Bell	113
Zorro	80
Zombieland	52
Zombie-Loan	38
Zombie Survival Guide	21

Zeta Project	13
Zombies, Run!	12
Zombie Powder	12
Zoom	9
Zombie Fallout	7
Zenon	6
Zeke and Luther	5
Zoey 101, iCarly	5
Zone of The Enders	5
Zapped	4
Zen	3
Zoo Tycoon	3
Zombie Prom	3
Zero Wing	3
Zack and Miri Make a Porno	3
Z.H.P.	2
Zetman	2
Zodiac P.I.	2
Zoop	2
Zevo-3	1
Zoolander	1
Zathura	1
Zero Dark Thirty	1
Zom-B	1
Zoombinis	1
Z for Zachariah	1
Zenda	1
Zeke's Pad	1
Z Nation	1
Zenonia	1
Zanna, Don't!	1
Zero Day	1
Zettai Karen Children	1
Zombie Blondes	1
Zinda/ज़िंदा	1

Name: category, dtype: int64

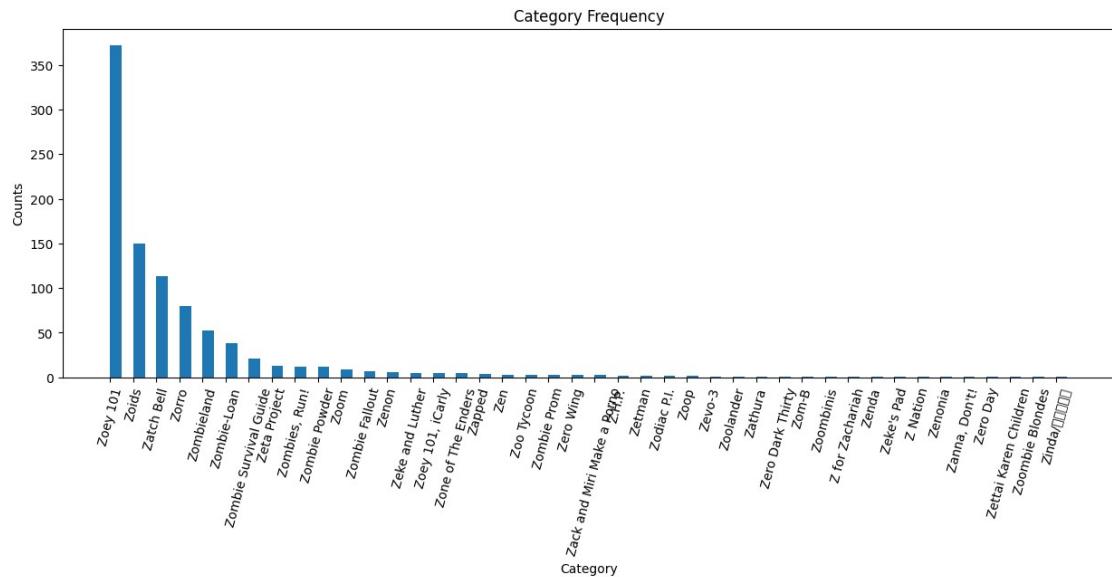
```
#Plot counts of stories per category
plt.figure(figsize=(15,5))
plt.xticks(rotation = 75)
plt.bar(category_count.index,category_count,align='edge', width= 0.5)
plt.title('Category Frequency')
plt.xlabel('Category')
plt.ylabel('Counts')
plt.show()

/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Glyph 2332 (\N{DEVANAGARI LETTER JA}) missing from
current font.
    fig.canvas.print_figure(bytes_io, **kw)
```

```

/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Matplotlib currently does not support Devanagari
natively.
    fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Glyph 2364 (\N{DEVANAGARI SIGN NUKTA}) missing from
current font.
    fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Glyph 2306 (\N{DEVANAGARI SIGN ANUSVARA}) missing from
current font.
    fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Glyph 2342 (\N{DEVANAGARI LETTER DA}) missing from
current font.
    fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
UserWarning: Glyph 2366 (\N{DEVANAGARI VOWEL SIGN AA}) missing from
current font.
    fig.canvas.print_figure(bytes_io, **kw)

```



### Expand Contractions

```

import re
contractions_dict = { "ain't": "are not", "'s": " is", "aren't": "are
not",
                      "can't": "cannot", "can't've": "cannot have",
                      "'cause": "because", "could've": "could
have", "couldn't": "could not",
                      "couldn't've": "could not have", "didn't": "did
not", "doesn't": "does not",
                      "don't": "do not", "hadn't": "had
not", "hadn't've": "had not have",

```

"hasn't": "has not", "haven't": "have not", "he'd": "he would",  
"he'd've": "he would have", "he'll": "he will",  
"he'll've": "he will have",  
"how'd": "how did", "how'd'y": "how do you", "how'll": "how will",  
"I'd": "I would", "I'd've": "I would have", "I'll": "I will",  
"I'll've": "I will have", "I'm": "I am", "I've": "I have", "isn't": "is not",  
"it'd": "it would", "it'd've": "it would have", "it'll": "it will",  
"it'll've": "it will have", "let's": "let us", "ma'am": "madam",  
"mayn't": "may not", "might've": "might have", "mightn't": "might not",  
"mightn't've": "might not have", "must've": "must have", "mustn't": "must not",  
"mustn't've": "must not have", "needn't": "need not",  
"needn't've": "need not have", "o'clock": "of the clock", "oughtn't": "ought not",  
"oughtn't've": "ought not have", "shan't": "shall not", "sha'n't": "shall not",  
"shan't've": "shall not have", "she'd": "she would", "she'd've": "she would have",  
"she'll": "she will", "she'll've": "she will have", "should've": "should have",  
"shouldn't": "should not", "shouldn't've": "should not have", "so've": "so have",  
"that'd": "that would", "that'd've": "that would have", "there'd": "there would",  
"there'd've": "there would have", "they'd": "they would",  
"they'd've": "they would have", "they'll": "they will",  
"they'll've": "they will have", "they're": "they are", "they've": "they have",  
"to've": "to have", "wasn't": "was not", "we'd": "we would",  
"we'd've": "we would have", "we'll": "we will", "we'll've": "we will have",  
"we're": "we are", "we've": "we have", "weren't": "were not", "what'll": "what will",  
"what'll've": "what will have", "what're": "what are", "what've": "what have",  
"when've": "when have", "where'd": "where did", "where've": "where have",  
"who'll": "who will", "who'll've": "who will have", "who've": "who have",

```

    "why've": "why have", "will've": "will
have", "won't": "will not",
    "won't've": "will not have", "would've": "would
have", "wouldn't": "would not",
    "wouldn't've": "would not have", "y'all": "you
all", "y'all'd": "you all would",
    "y'all'd've": "you all would have", "y'all're":
"you all are",
    "y'all've": "you all have", "you'd": "you
would", "you'd've": "you would have",
    "you'll": "you will", "you'll've": "you will
have", "you're": "you are",
    "you've": "you have"}

```

*# Regular expression for finding contractions*

```

contractions_re=re.compile('(%s)' %
'|'.join(contractions_dict.keys())))

```

*# Function for expanding contractions*

```

def expand_contractions(text,contractions_dict=contractions_dict):
    def replace(match):
        return contractions_dict[match.group(0)]
    return contractions_re.sub(replace, text)

```

*# Expanding Contractions in the reviews*

```

df['story_cleaned']=df['story_cleaned'].apply(lambda
x:expand_contractions(x))

```

### **Lower case**

```

df['story_cleaned']=df['story_cleaned'].apply(lambda x: x.lower())

```

df

```

                           story \
0   just a question\n\nWhat I thought when I first...
1   pain\n\nIt may be sad but it can be better\n\n...
2   Zombie Mayhem\n\nZombie Mayhem\n\nIn 2050, the...
3   Zero Wing 2: AYBASBTU Game Script\n\nZERO WING...
4   Not prosaic\n\nDisclaimer, I do not claim to o...
..
938  More\n\n**I've watched this movie 3 times this...
939  Change Comes From Within\n\n**As usual, I own ...
940  Changes\n\nYou know that feeling where people ...
941  Undying Hero\n\nZero percent.\n\nVictory canno...
942  Unlosing Haunted Mansion\n\n**Disclaimer**: I ...

```

```

                           title \
0   just a question by multifics123
1   pain by Isaiah Thomas

```

```

2           Zombie Mayhem by kelbey342
3   Zero Wing 2: AYBASBTU Game Script by Rodrigo Shin
4                   Not prosaic by Fault
..
938           ...
939           More by Lizwontcry
940       Change Comes From Within by Dominican Girl
941           Changes by michieexx3
942           Undying Hero by Arukoir
943       Unlosing Haunted Mansion by Smak64

          category \
0           Zom-B
1           Zenda
2       Zoombie Blondes
3           Zero Wing
4           Zero Wing
..
938   ...
939   Zack and Miri Make a Porno
940   Zack and Miri Make a Porno
941   Zack and Miri Make a Porno
942           Z.H.P.
943           Z.H.P.

          story_cleaned
0   just a question\nwhat i thought when i first r...
1   pain\nit may be sad but it can be better\nthe ...
2   zombie mayhem\nzombie mayhem\nin 2050, the mil...
3   zero wing 2: aybasbtu game script\nzero wing 2...
4   not prosaic\ndisclaimer, i do not claim to own...
..
938   ...
939   more\ni have watched this movie 3 times this w...
940   change comes from within\nas usual, i own noth...
941   changes\nyou know that feeling where people be...
942   undying hero\nzero percent.\nvictory cannot be...
943   unlosing haunted mansion\ndisclaimer: i do not...

```

[943 rows x 4 columns]

### ***Remove digits and words containing digits***

```
df['story_cleaned']=df['story_cleaned'].apply(lambda x: re.sub('\w*\d\w*', '', x))
```

### ***Remove Punctuations***

```
import string
df['story_cleaned']=df['story_cleaned'].apply(lambda x: re.sub('[\%]' % re.escape(string.punctuation), '', x))

df
```

0 just a question\n\nWhat I thought when I first... story  
1 pain\n\nIt may be sad but it can be better\n\n...  
2 Zombie Mayhem\n\nZombie Mayhem\n\nIn 2050, the...  
3 Zero Wing 2: AYBASBTU Game Script\n\nZERO WING...  
4 Not prosaic\n\nDisclaimer, I do not claim to o...  
..  
938 More\n\n\*\*I've watched this movie 3 times this...  
939 Change Comes From Within\n\n\*\*As usual, I own ...  
940 Changes\n\nYou know that feeling where people ...  
941 Undying Hero\n\nZero percent.\n\nVictory canno...  
942 Unlosing Haunted Mansion\n\n\*\*Disclaimer\*\*: I ...

		title \
0	just a question by multifics123	
1	pain by Isaiah Thomas	
2	Zombie Mayhem by kelbey342	
3	Zero Wing 2: AYBASBTU Game Script by Rodrigo Shin	
4	Not prosaic by Fault	
..		...
938	More by Lizwontcry	
939	Change Comes From Within by Dominican Girl	
940	Changes by michieexx3	
941	Undying Hero by Arukoir	
942	Unlosing Haunted Mansion by Smak64	

		category	\
0		Zom-B	
1		Zenda	
2		Zoombie Blondes	
3		Zero Wing	
4		Zero Wing	
..		..	..
938	Zack and Miri Make a Porno		
939	Zack and Miri Make a Porno		
940	Zack and Miri Make a Porno		
941		Z.H.P.	
942		Z.H.P.	

0 just a question\nwhat i thought when i first r...  
1 pain\nit may be sad but it can be better\nthe ...  
2 zombie mayhem\nzombie mayhem\nin the military...  
3 zero wing aybasbtu game script\nzero wing \...  
4 not prosaic\nndisclaimer i do not claim to own ...  
..  
938 more\ni have watched this movie times this we...  
939 change comes from within\nas usual i own nothi...  
940 changes\nyou know that feeling where people be...  
941 undying hero\nzero percent\nvictory cannot be ...

```
942 unlosing haunted mansion\nDisclaimer i do not ...
```

```
[943 rows x 4 columns]
```

### **Removing extra spaces**

```
df['story_cleaned']=df['story_cleaned'].apply(lambda x: re.sub(' +', ' ',x))

df['story_cleaned']

0    just a question\nwhat i thought when i first r...
1    pain\nit may be sad but it can be better\nthe ...
2    zombie mayhem\nzombie mayhem\nin the military ...
3    zero wing aybasbtu game script\nzero wing \nal...
4    not prosaic\nDisclaimer i do not claim to own ...

         ..
938   more\ni have watched this movie times this wee...
939   change comes from within\nas usual i own nothi...
940   changes\nyou know that feeling where people be...
941   undying hero\nzero percent\nvictory cannot be ...
942   unlosing haunted mansion\nDisclaimer i do not ...

Name: story_cleaned, Length: 943, dtype: object
```

```
import spacy
```

```
#Loading the model and then lemmatization
nlp = spacy.load('en_core_web_sm', disable = ['parser','ner'])
df['story_lemmatized'] = df['story_cleaned'].apply(lambda x:''.join([token.lemma_ for token in list(nlp(x)) if (token.is_stop == False)]))

df

                           story \
0    just a question\n\nWhat I thought when I first...
1    pain\n\nIt may be sad but it can be better\n\n...
2    Zombie Mayhem\n\nZombie Mayhem\n\nIn 2050, the...
3    Zero Wing 2: AYBASBTU Game Script\n\nZERO WING...
4    Not prosaic\n\nDisclaimer, I do not claim to o...

         ..
938   More\n\n**I've watched this movie 3 times this...
939   Change Comes From Within\n\n**As usual, I own ...
940   Changes\n\nYou know that feeling where people ...
941   Undying Hero\n\nZero percent.\n\nVictory canno...
942   Unlosing Haunted Mansion\n\n**Disclaimer**: I ...

                           title \
0    just a question by multifics123
1                  pain by Isaiah Thomas
2      Zombie Mayhem by kelbey342
```

3 Zero Wing 2: AYBASBTU Game Script by Rodrigo Shin  
4 Not prosaic by Fault

..  
938 More by Lizwontcry  
939 Change Comes From Within by Dominican Girl  
940 Changes by michieexx3  
941 Undying Hero by Arukoir  
942 Unlosing Haunted Mansion by Smak64

category \  
0 Zom-B  
1 Zenda  
2 Zoombie Blondes  
3 Zero Wing  
4 Zero Wing

..  
938 Zack and Miri Make a Porno  
939 Zack and Miri Make a Porno  
940 Zack and Miri Make a Porno  
941 Z.H.P.  
942 Z.H.P.

story\_cleaned \  
0 just a question\nwhat i thought when i first r...  
1 pain\nit may be sad but it can be better\nthe ...  
2 zombie mayhem\nzombie mayhem\nin the military ...  
3 zero wing aybasbtu game script\nzero wing \nal...  
4 not prosaic\nDisclaimer i do not claim to own ...  
..  
938 more\ni have watched this movie times this wee...  
939 change comes from within\nas usual i own nothi...  
940 changes\nyou know that feeling where people be...  
941 undying hero\nzero percent\nvictory cannot be ...  
942 unlosing haunted mansion\nDisclaimer i do not ...

story\_lemmatized  
0 question \n think read b guy think plz leave r...  
1 pain \n sad well \n yellow sun glisten face \n...  
2 zombie mayhem \n zombie mayhem \n military tes...  
3 zero wing aybasbtu game script \n zero wing \n...  
4 prosaic \n disclaimer claim zero wing make mon...  
..  
938 \n watch movie time week love urge write littl...  
939 change come \n usual character belong \n \n ex...  
940 change \n know feel people believe impacting e...  
941 undying hero \n zero percent \n victory gain \n...  
942 unlose haunt mansion \n disclaimer zettai hero...

[943 rows x 5 columns]

```
#Grouping data according to the category
df_grouped =
df[['category','story_lemmatized']].groupby(by='category').agg(lambda
x: ' '.join(x))
df_grouped

story_lemmatized
category

Z Nation          violent delight \n violent delight \n
desire c...
Z for Zachariah   get save \n hey guy hope like god bless
check ...
Z.H.P.            undying hero \n zero percent \n victory
gain \...
Zack and Miri Make a Porno \n watch movie time week love urge write
littl...
Zanna, Don't!    story \n p styletextalign startnot story
try b...
Zapped           trust \n hello people \n see barely fan
fic th...
Zatch Bell        healing tear \n sherryou fight hardplease
for...
Zathura          dominique jakowec \n test type \n\n end
file \n
Zeke and Luther   wanna skateboard \n helllo know upload new
cha...
Zeke's Pad         meaning friendship \n write month ago see
fina...
Zen             informer \n informer \n zen arianna \n
second ...
Zenda           pain \n sad well \n yellow sun glisten
face \n...
Zenon           chapter \n flame red soil \n disclaimer
zeno...
Zenonia          chapter \n yeah yeah update schedule hm
anyw...
Zero Dark Thirty zero dark thirty alternate footage \n
member s...
Zero Day          mad world \n late christmas present good
frien...
Zero Wing         zero wing aybasbtu game script \n zero
wing \n...
Zeta Project      old thing \n old thing \n pass trough door
hea...
Zetman           poem \n test \n eye \n test story \n party
\n ...
Zettai Karen Children tea time \n tea time \n set postunlimite
andy ...
```

Zevo-3 memory \n memory \n \n unfortunately  
zevo \n j... \n planet p \n disclaimer zinda character \n \n ...  
Zinda/ज़िंदा Zodiac P.I. christmas shoe \n disclaimer zodiac pi  
Zodiac P.I. base st...  
base st... Zoey 101 song \n author note story want song  
think go... Zoey 101, iCarly cotton swab christmas tree \n \n cotton  
swab c... Zoids christmas proposal \n author note wow know  
day... Zom-B question \n think read b guy think plz  
leave r... Zombie Fallout end world end apolla \n end end apolla \n  
song... Zombie Powder wait \n title wait \n characterspairing  
elwood... Zombie Prom scream \n scream \n cry room \n scream bad  
doo... Zombie Survival Guide everstone \n monday run right scream kevin  
run... Zombie-Loan natural phenomenon \n person death natural  
phe... Zombieland zombie hunter \n zombie hunter \n season  
final... Zombies, Run! home \n procrastinate update fanfiction  
write ... Zone of The Enders leo play hotline miami \n leo play hotline  
mia... Zoo Tycoon fanfiction \n late night late usual dark  
wet... Zoolander hansel \n derek pout lip hansel think  
hanstupi... Zoom nightmare \n anjust quick oneshot \n  
disclaime... Zoombie Blondes zombie mayhem \n zombie mayhem \n military  
tes... Zoombinis put bubble machine inside mountain \n  
wonder b... Zoop jacksepticeye \n second day america  
birthday d... Zorro unconsciously conscious \n author note  
word dr...

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
#Create a Document term matrix for plotting word clouds
```

```
cv=CountVectorizer(analyzer='word')  
data=cv.fit_transform(df_grouped['story_lemmatized'])  
df_dtm = pd.DataFrame(data.toarray()),
```

```
columns=cv.get_feature_names_out()
df_dtm.index=df_grouped.index
df_dtm.head(10)

aa    aaa   aaaaaaaaaaaaaa
aaaaaaaaaaaaaaaaaaaaaa \
category

Z Nation          0    0    0
0
Z for Zachariah 0    0    0
0
Z.H.P.           0    0    0
0
Zack and Miri Make a Porno 0    0    0
0
Zanna, Don't!    0    0    0
0
Zapped            0    0    0
0
Zatch Bell        1    0    0
0
Zathura           0    0    0
0
Zeke and Luther  1    0    0
0
Zeke's Pad        0    0    0
0
```

```
aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa
aaaaaaaaaaaaahhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
h \
category
```

```
Z Nation
0

Z for Zachariah
0

Z.H.P.
0

Zack and Miri Make a Porno
0

Zanna, Don't!
0
```

Zapped  
0

Zatch Bell  
1

Zathura  
0

Zeke and Luther  
0

Zeke's Pad  
0

aaa  
aaaaahhh  
hh  
category

Z Nation  
0

Z for Zachariah  
0

Z.H.P.  
0

Zack and Miri Make a Porno  
0

Zanna, Don't!  
0

Zapped  
0

Zatch Bell  
1

Zathura  
0

Zeke and Luther

0

Zeke's Pad

0

aaaaaaaaaaaaaaaaaaaaahhhhhhhhhh \

category

Z Nation

0

Z for Zachariah

0

Z.H.P.

0

Zack and Miri Make a Porno

0

Zanna, Don't!

0

Zapped

0

Zatch Bell

0

Zathura

0

Zeke and Luther

0

Zeke's Pad

0

aaaaaaaaaaaaanyway aaaaaaaaaahhh

aaaaaaaaahhhhhh \

category

Z Nation

0

0

Z for Zachariah

0

0

Z.H.P.

0

0

Zack and Miri Make a Porno

0

0

0

Zanna, Don't!

0

0

Zapped

0

0

Zatch Bell

0

0

Zathura

0

0

0

Zeke and Luther

0

0

0

Zeke's Pad

0

0

0

... zzz zzzaker zzzoia zzzoid

zzzzzzzzz \

...

category

Z Nation

...

0

0

0

0

Z for Zachariah	...	0	0	0	0
Z.H.P.	...	0	0	0	0
Zack and Miri Make a Porno	...	0	0	0	0
Zanna, Don't!	...	0	0	0	0
Zapped	...	0	0	0	0
Zatch Bell	...	0	1	0	0
Zathura	...	0	0	0	0
Zeke and Luther	...	0	0	0	0
Zeke's Pad	...	0	0	0	0

zzzzzzzzzzzzzzzz zzzzzzzzzzzzzzzz \

## category

Z Nation	0	0
Z for Zachariah	0	0
Z.H.P.	0	0
Zack and Miri Make a Porno	0	0
Zanna, Don't!	0	0
Zapped	0	0
Zatch Bell	0	0
Zathura	0	0
Zeke and Luther	0	0
Zeke's Pad	0	0

ZZZZZZZZZZZZZZZZZZZZZ ZZZZZZZZZZZZZZZZZZZZ \

## category

Z Nation	0	0
Z for Zachariah	0	0
Z.H.P.	0	0
Zack and Miri Make a Porno	0	0
Zanna, Don't!	0	0
Zapped	0	0
Zatch Bell	0	0
Zathura	0	0
Zeke and Luther	0	0
Zeke's Pad	0	0

```
Z Nation
0
Z for Zachariah
0
Z.H.P.
0
Zack and Miri Make a Porno
0
Zanna, Don't!
0
Zapped
0
Zatch Bell
0
Zathura
0
Zeke and Luther
0
Zeke's Pad
0
```

[10 rows x 34485 columns]

```
# Importing wordcloud for plotting word clouds and textwrap for
wrapping longer text
from wordcloud import WordCloud
from textwrap import wrap

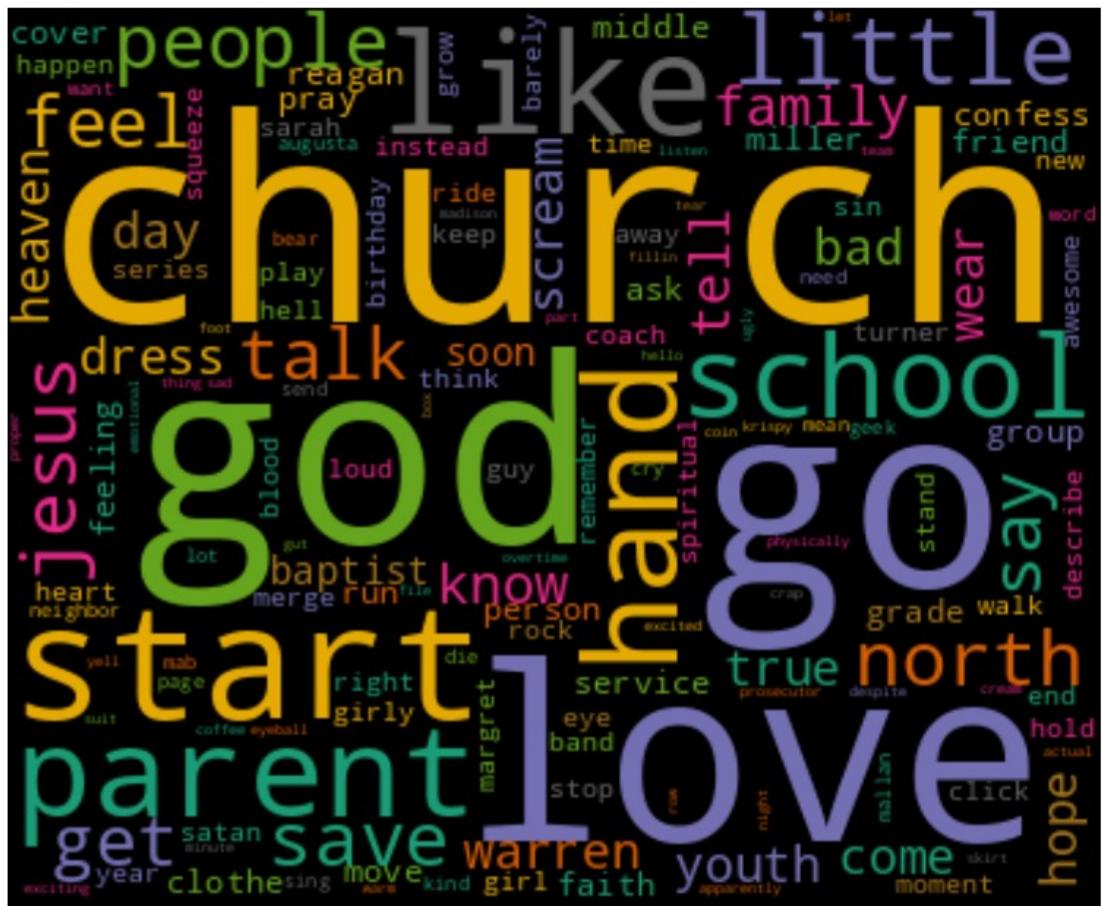
# Main method for generating word clouds (freq of different words in a
doc)
def generate_wordcloud(data,title):
    wc = WordCloud(width=400, height=330,
max_words=150,colormap="Dark2").generate_from_frequencies(data)
    plt.figure(figsize=(10,8))
    plt.imshow(wc, interpolation='bilinear')
    plt.axis("off")
    plt.title('\n'.join(wrap(title,60)),fontsize=13)
    plt.show()

# Transposing document term matrix
df_dtm=df_dtm.transpose()
# Plotting word cloud for each product
for index,product in enumerate(df_dtm.columns):

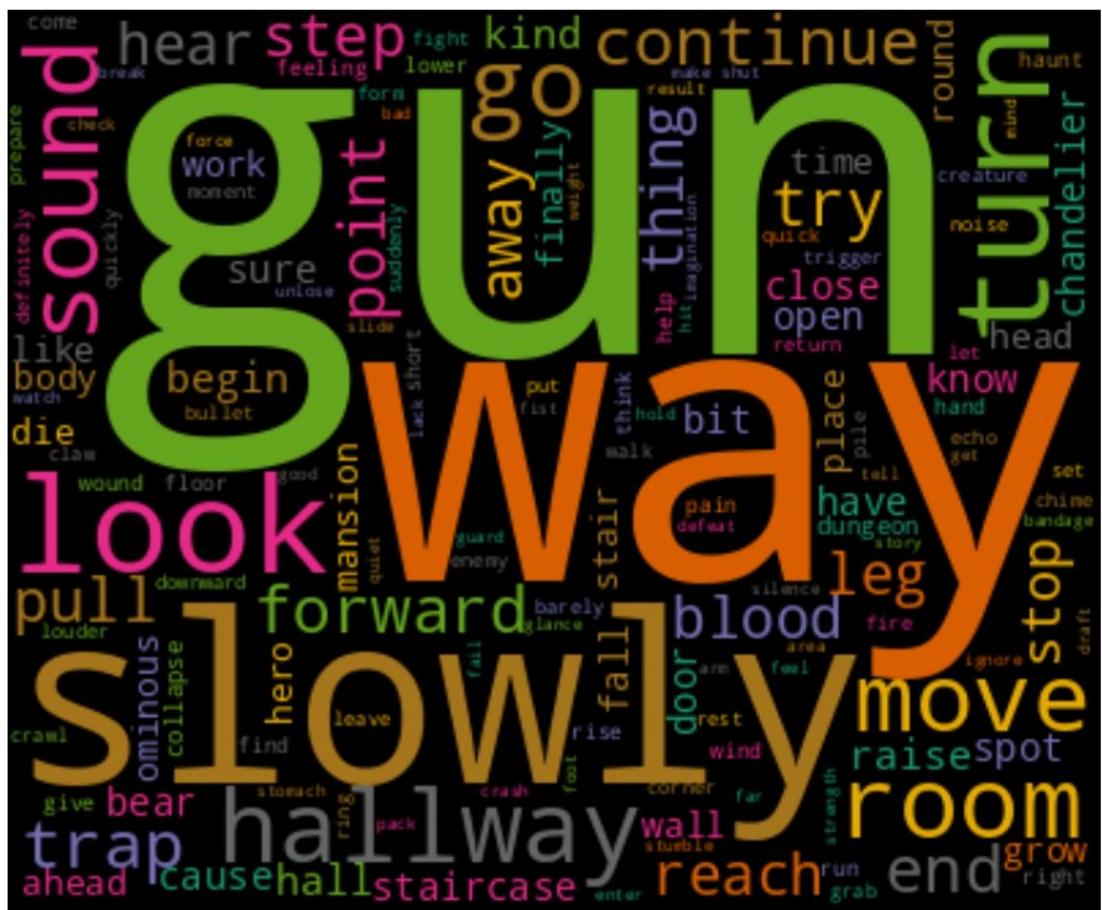
    generate_wordcloud(df_dtm[product].sort_values(ascending=False),produ
t)
```

Z Nation

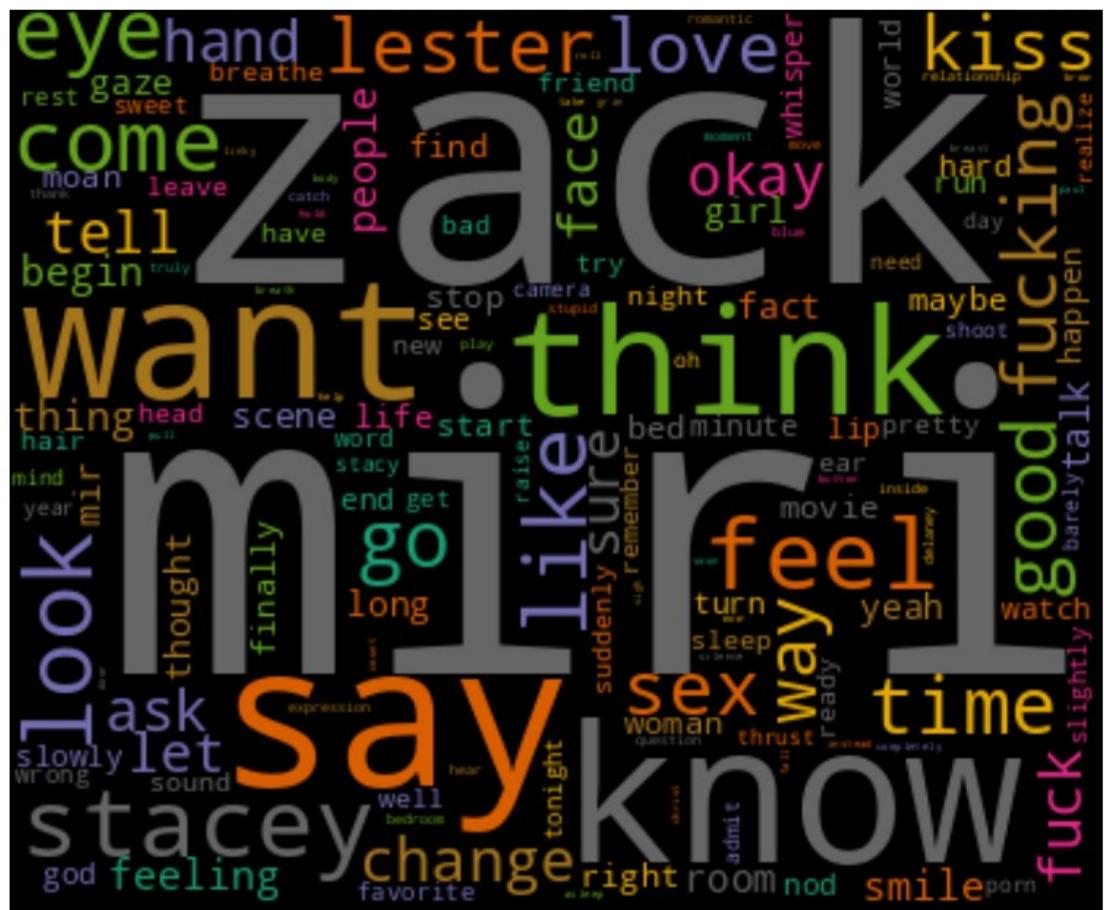
## Z for Zachariah



Z.H.P.



## Zack and Miri Make a Porno



Zanna, Don't!

startnot  
file end  
styletextalign

beta story  
trySorry

## Zapped



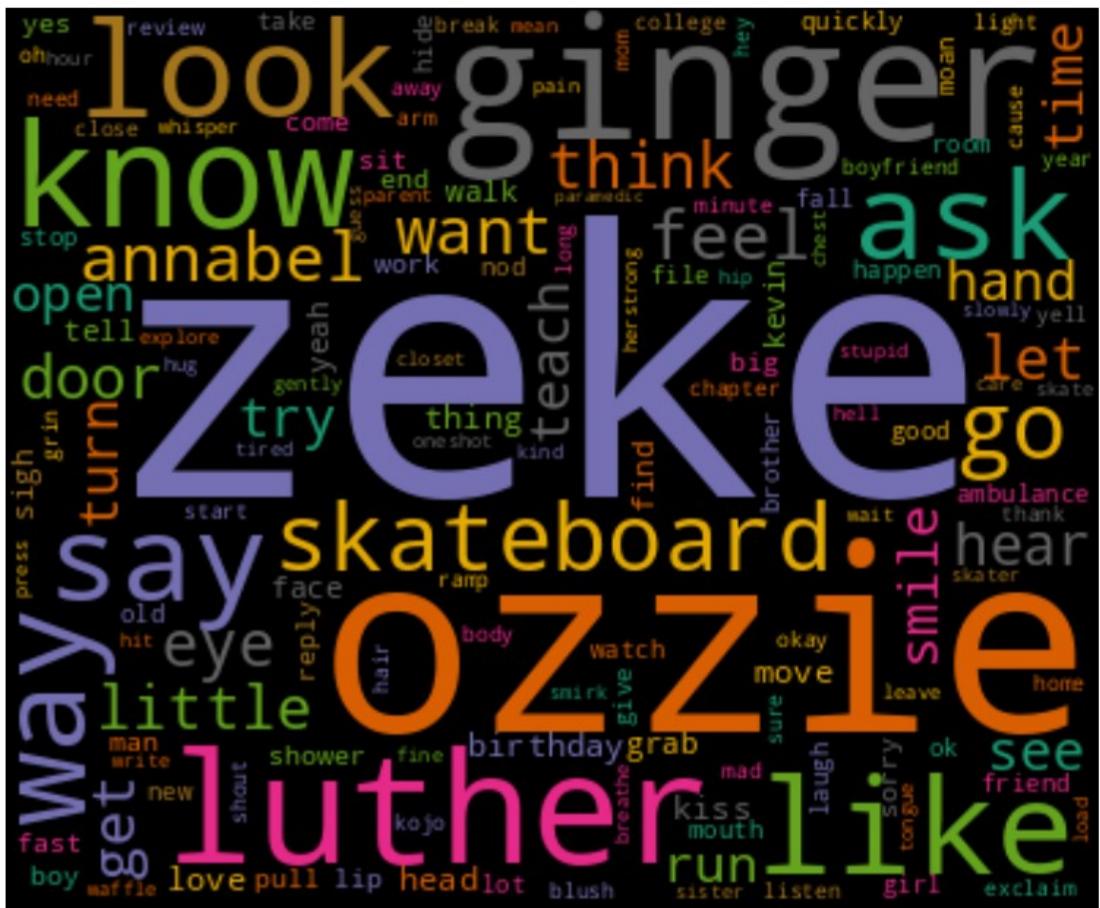
Zatch Bell



Zathura

jakowec  
dominique  
endtype  
file  
test

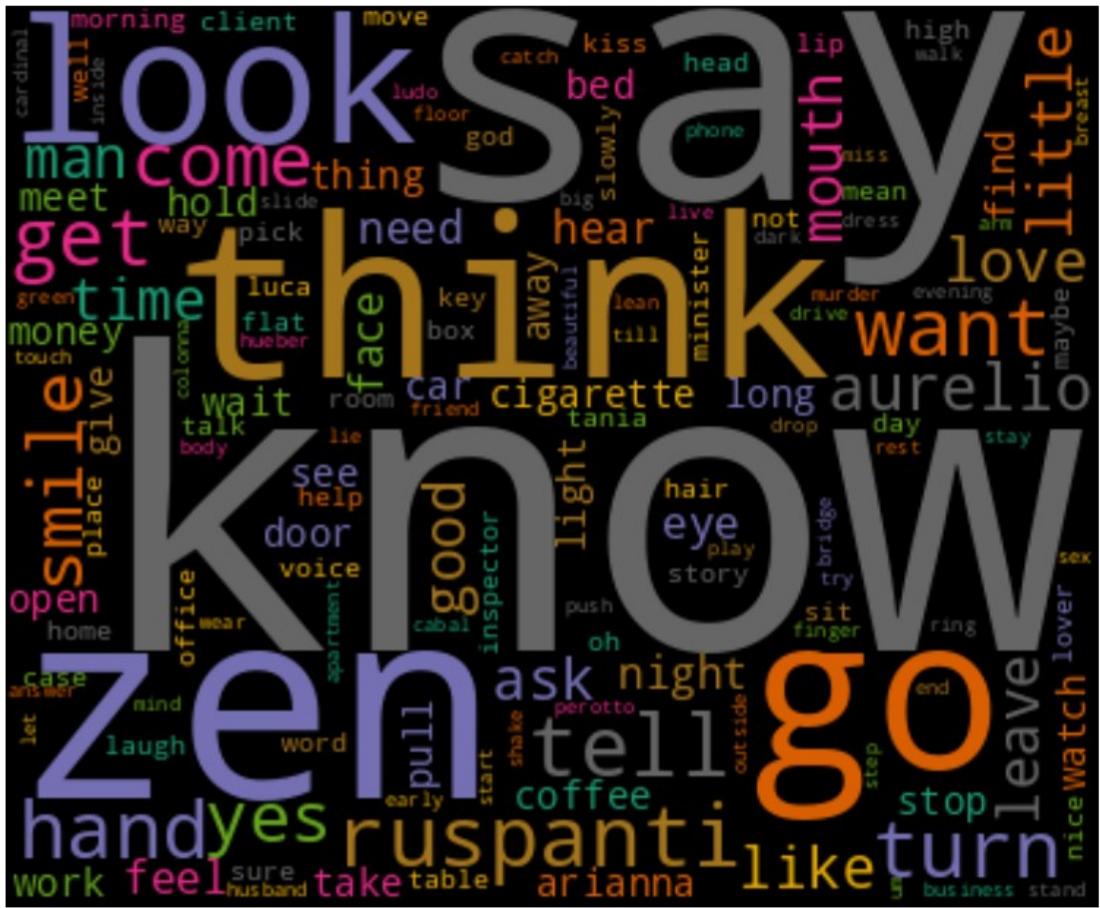
## Zeke and Luther



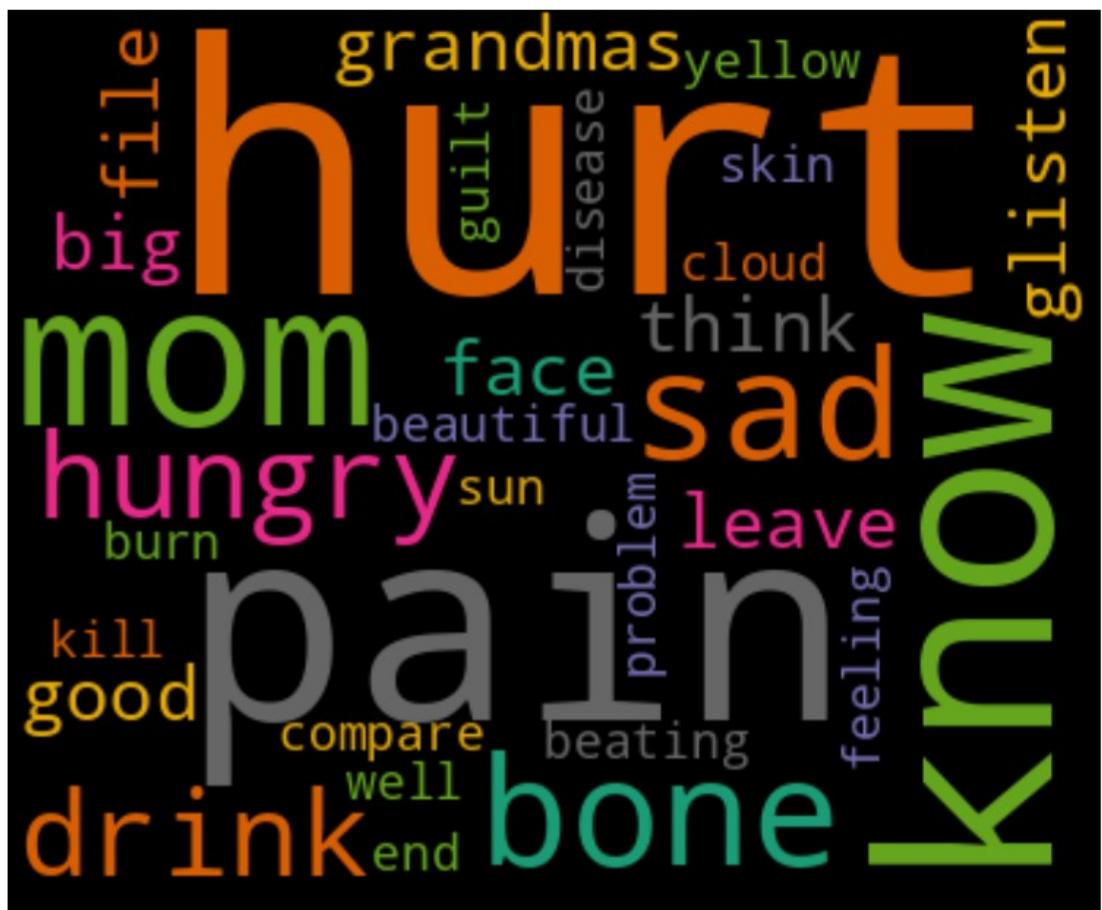
## Zeke's Pad

life rely friendless see matter type friendless need rush away month month fix sure lately mistake section file sleep far intuitively meet simply anymore time blast ear upload problem unlikely suddenly loneliness pillow tell book friend carelessness ago Zekere night help have collide want somebody finally get middle actually difficulty sleeping

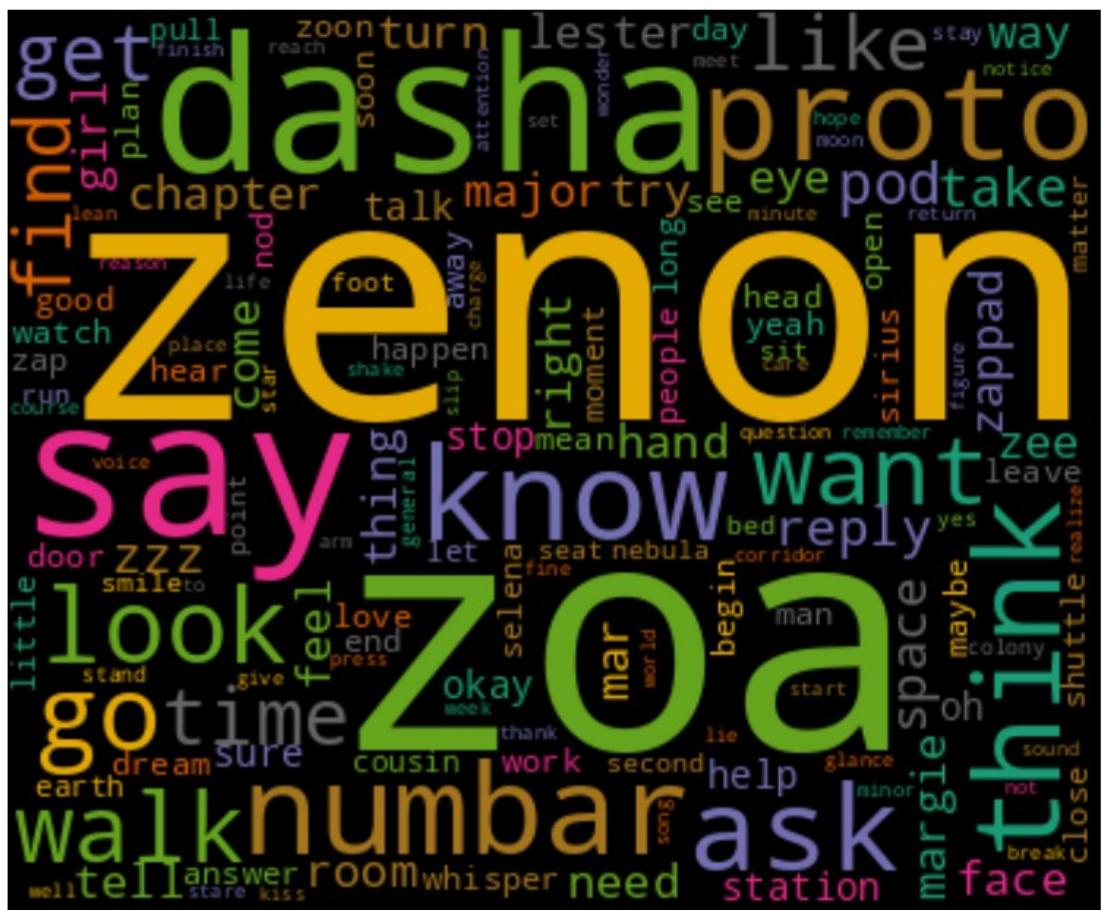
Zen



Zenda



Zenon



Zenonia

young long seat seat  
mutant chapter immature  
treat cause talk general  
stare hand space  
try write  
tell start  
eye finally believe  
bar book nice  
thing message  
read stop  
look yell  
open brunette monster  
near week  
fact give way  
chest hero take  
step  
body glass people  
begin understand  
small old  
anomaly step  
perpetually get  
small body glass people  
begin understand  
age turn suddenly  
know spread  
door pure away  
chase halfway  
smile blade drink  
say go room end  
aria go room end  
right friend true case  
entry day like  
leave guard live  
mean leave  
force shoot  
table mean  
utter like  
blood live  
think let past  
slowly story actually  
man little  
call hear single  
sigh single  
boy wall  
zenonia kiss

## Zero Dark Thirty

Zero Day



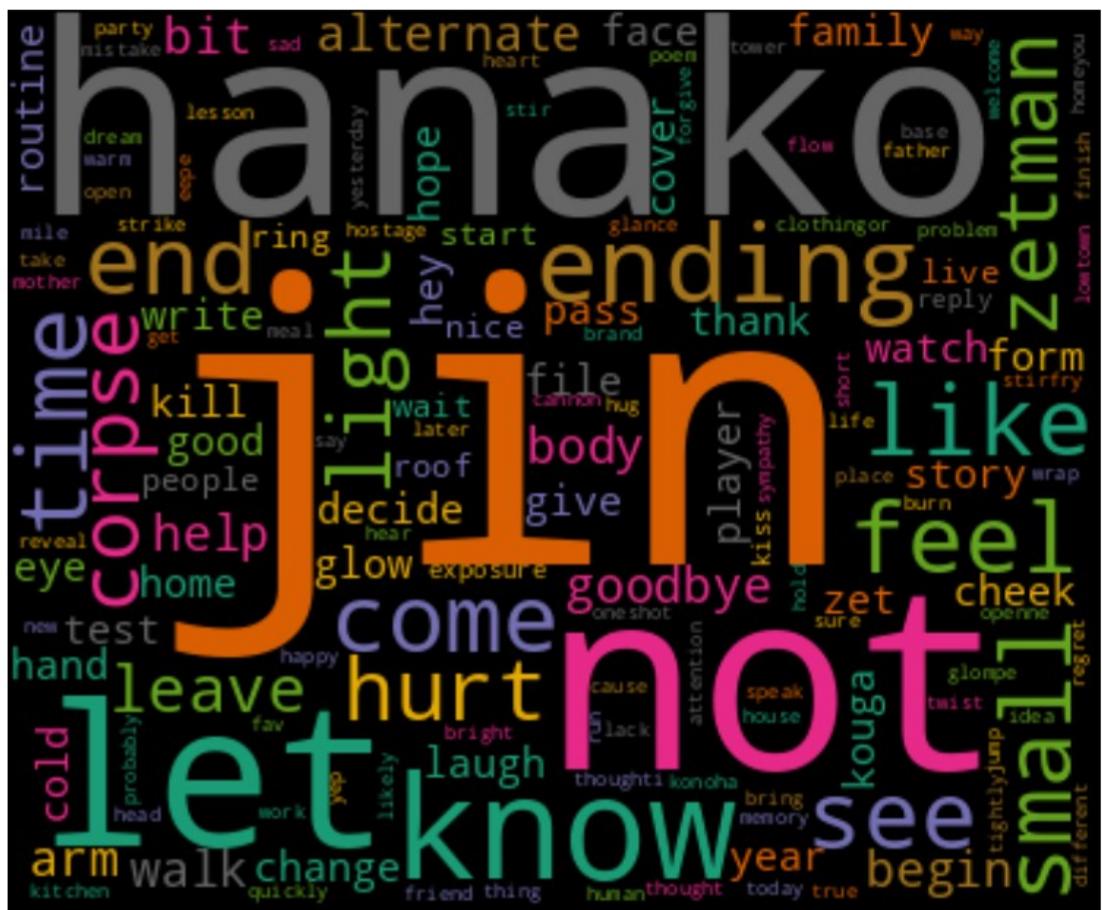
## Zero Wing



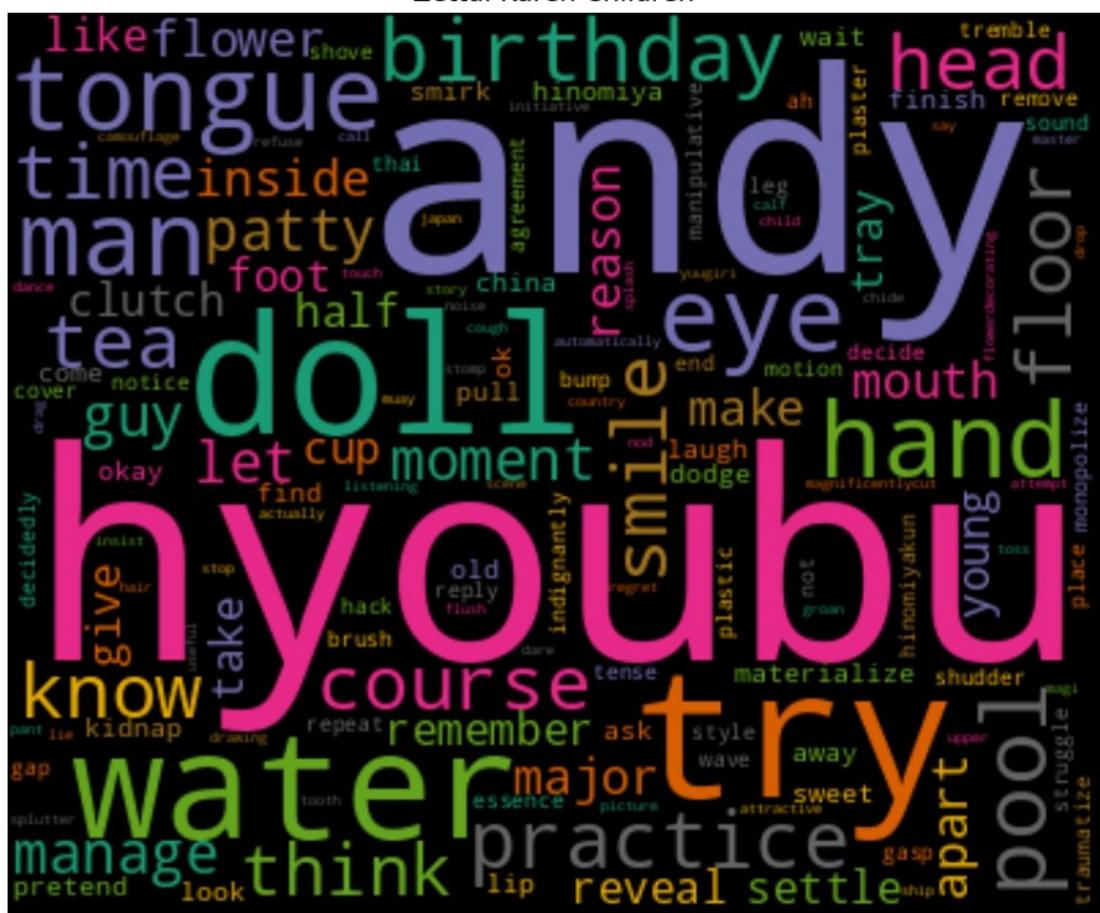
Zeta Project



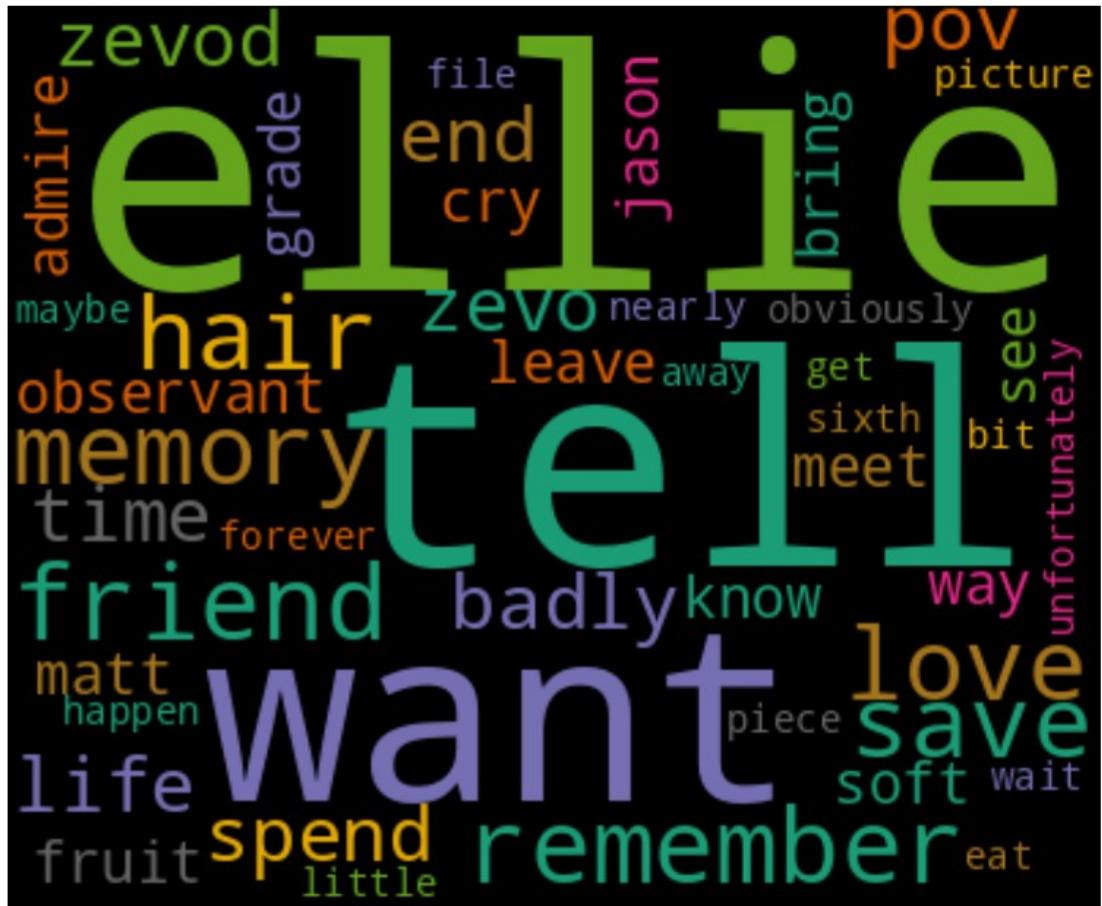
Zetman



## Zettai Karen Children



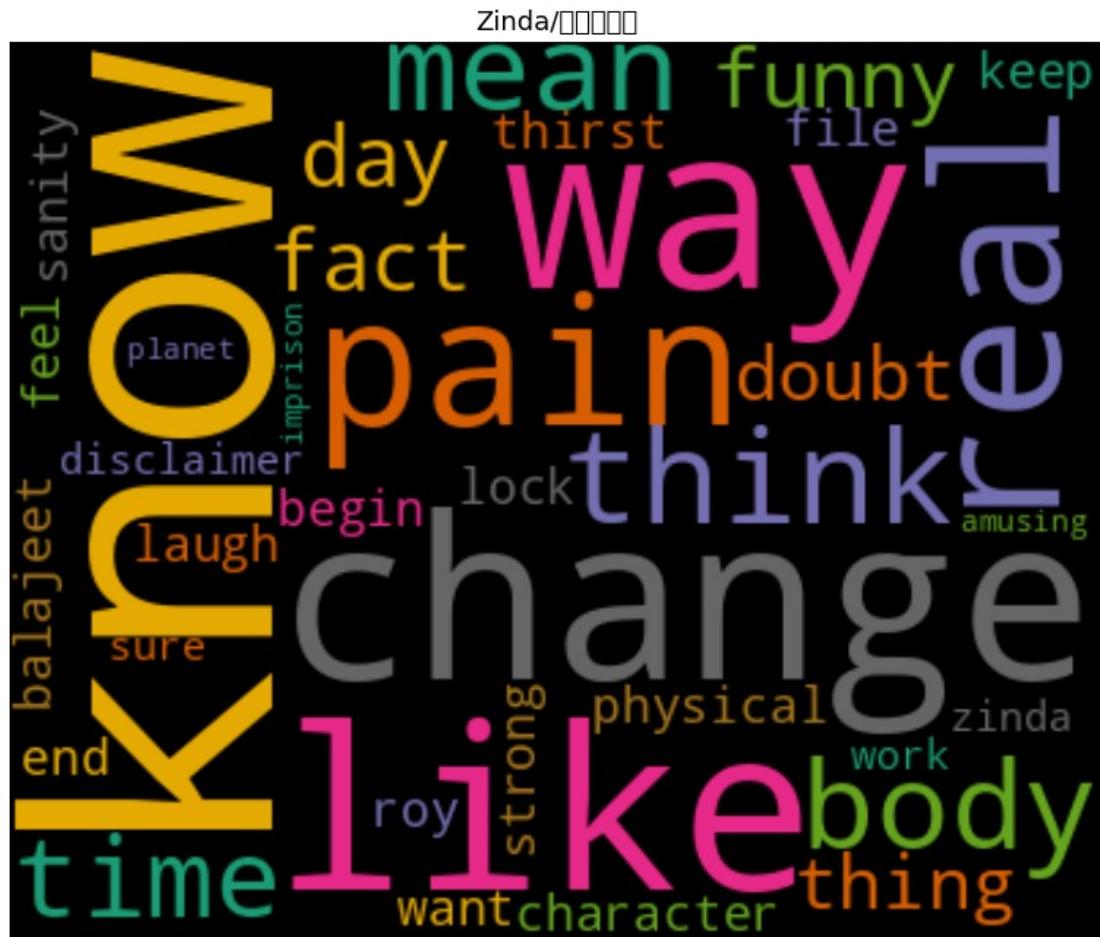
Zevo-3



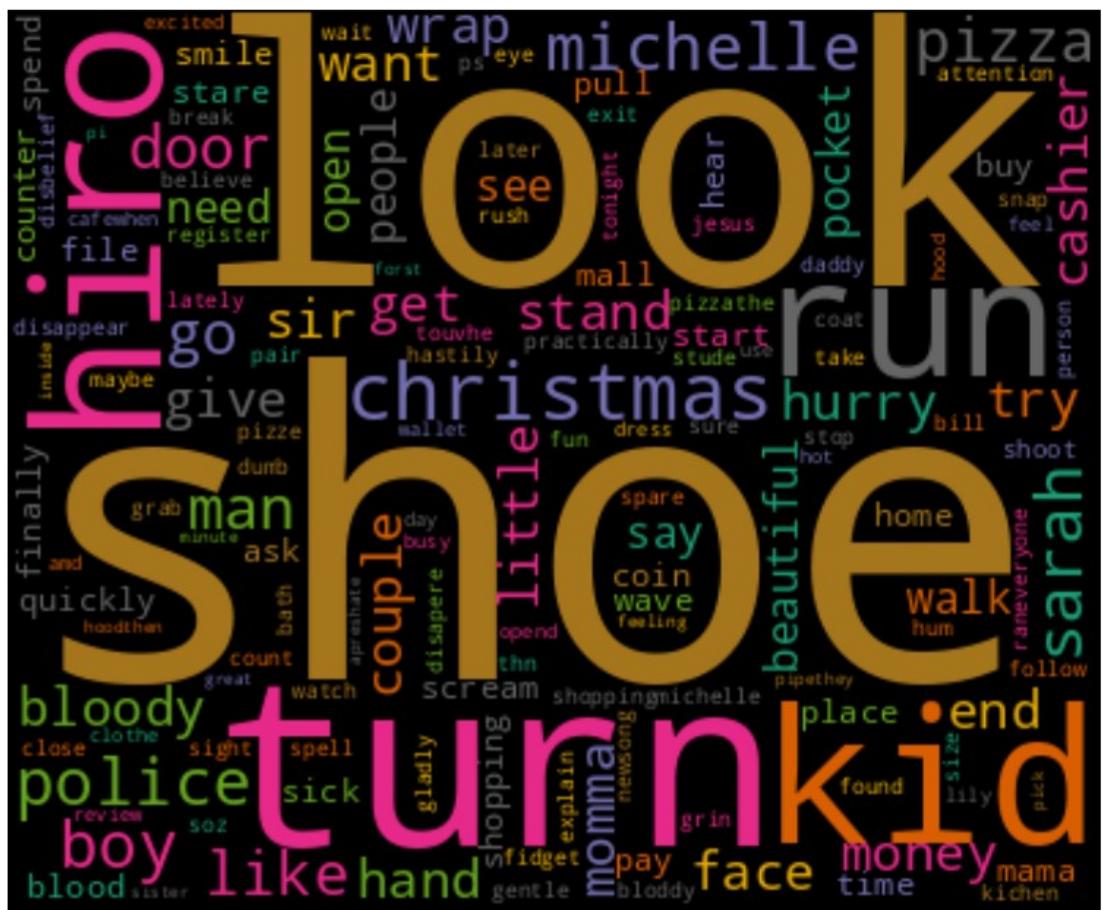
```
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:  
UserWarning: Glyph 2332 (\N{DEVANAGARI LETTER JA}) missing from  
current font.  
    fig.canvas.print_figure(bytes_io, **kw)  
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:  
UserWarning: Matplotlib currently does not support Devanagari  
natively.  
    fig.canvas.print_figure(bytes_io, **kw)  
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:  
UserWarning: Glyph 2364 (\N{DEVANAGARI SIGN NUKTA}) missing from  
current font.  
    fig.canvas.print_figure(bytes_io, **kw)  
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:  
UserWarning: Glyph 2306 (\N{DEVANAGARI SIGN ANUSVARA}) missing from  
current font.  
    fig.canvas.print_figure(bytes_io, **kw)  
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:  
UserWarning: Glyph 2342 (\N{DEVANAGARI LETTER DA}) missing from  
current font.  
    fig.canvas.print_figure(bytes_io, **kw)  
/usr/local/lib/python3.9/dist-packages/IPython/core/pylabtools.py:151:
```

UserWarning: Glyph 2366 (\N{DEVANAGARI VOWEL SIGN AA}) missing from current font.

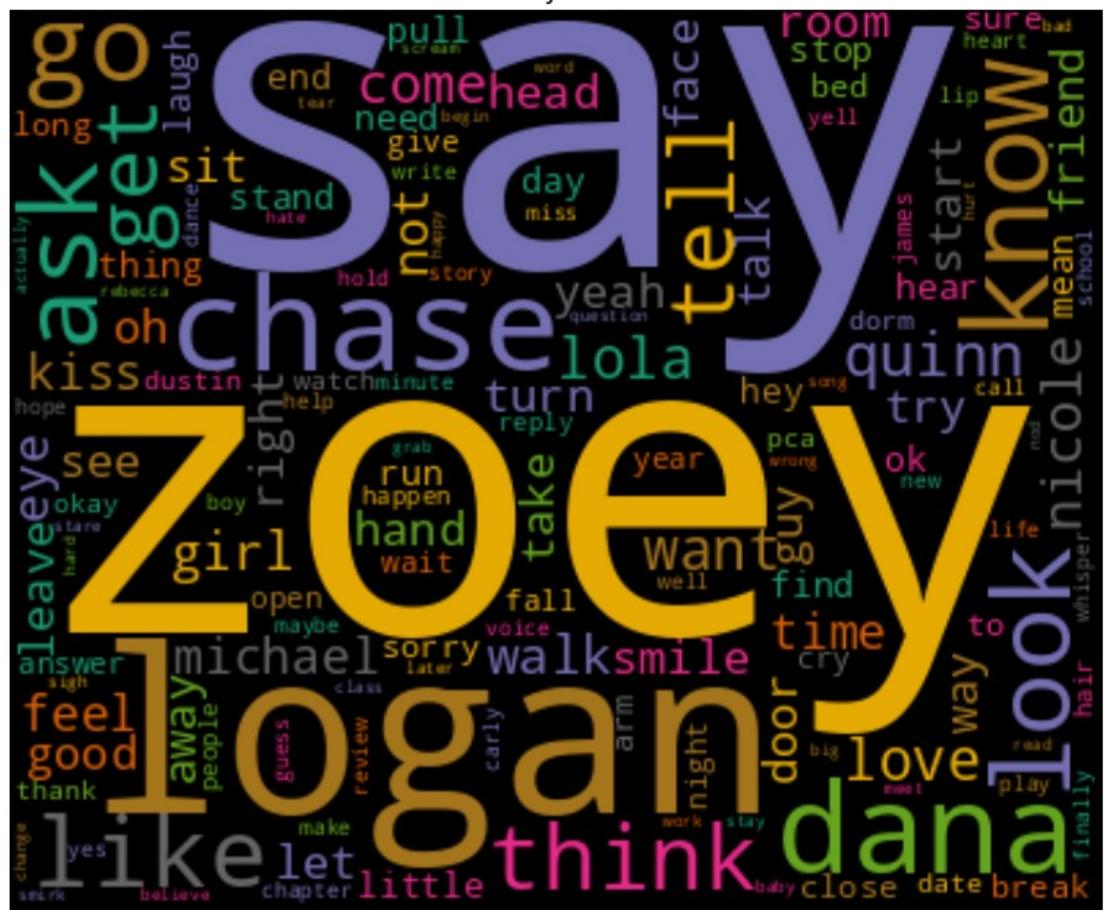
```
fig.canvas.print_figure(bytes_io, **kw)
```



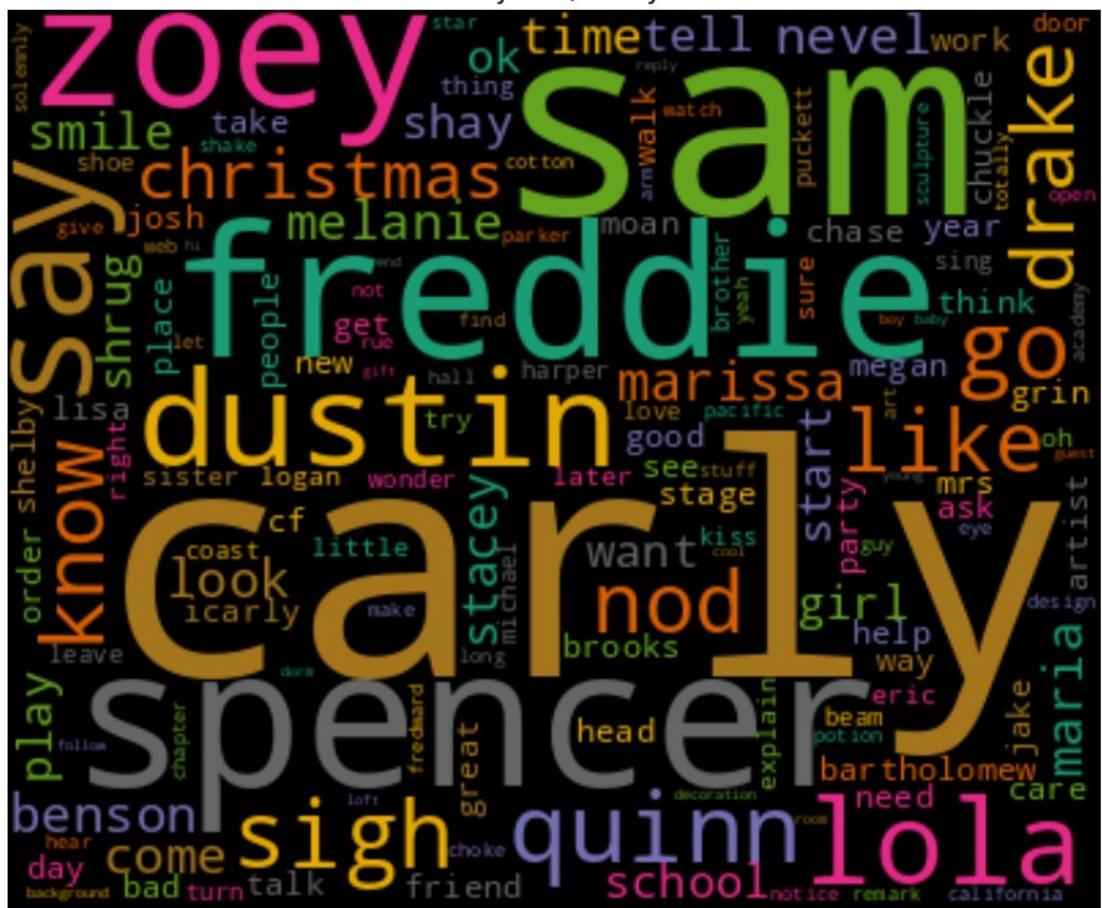
Zodiac P.I.



Zoey 101



Zoey 101, iCarly



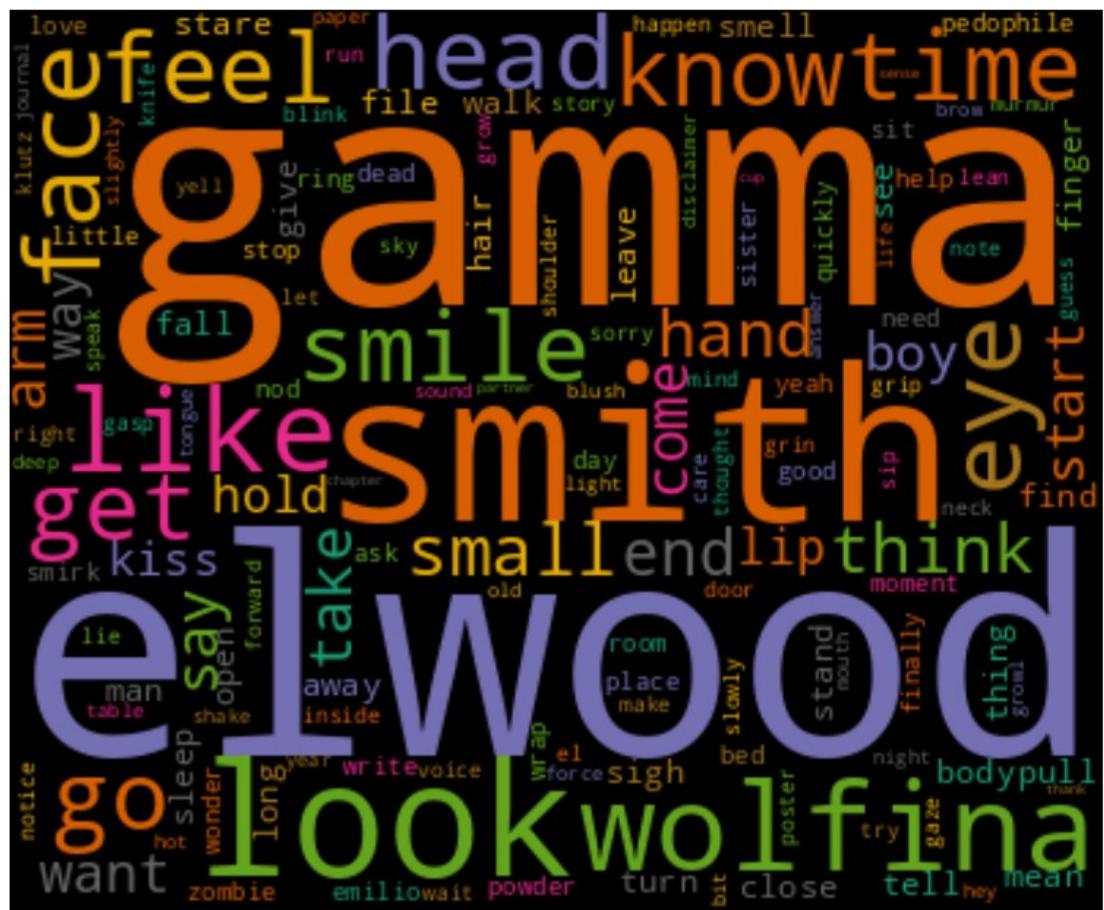
## Zoids

Zom-B

book love file end  
fanfiction  
thank  
series question guy  
plz add actual  
start read leave fanfic  
review think people

## Zombie Fallout

## Zombie Powder

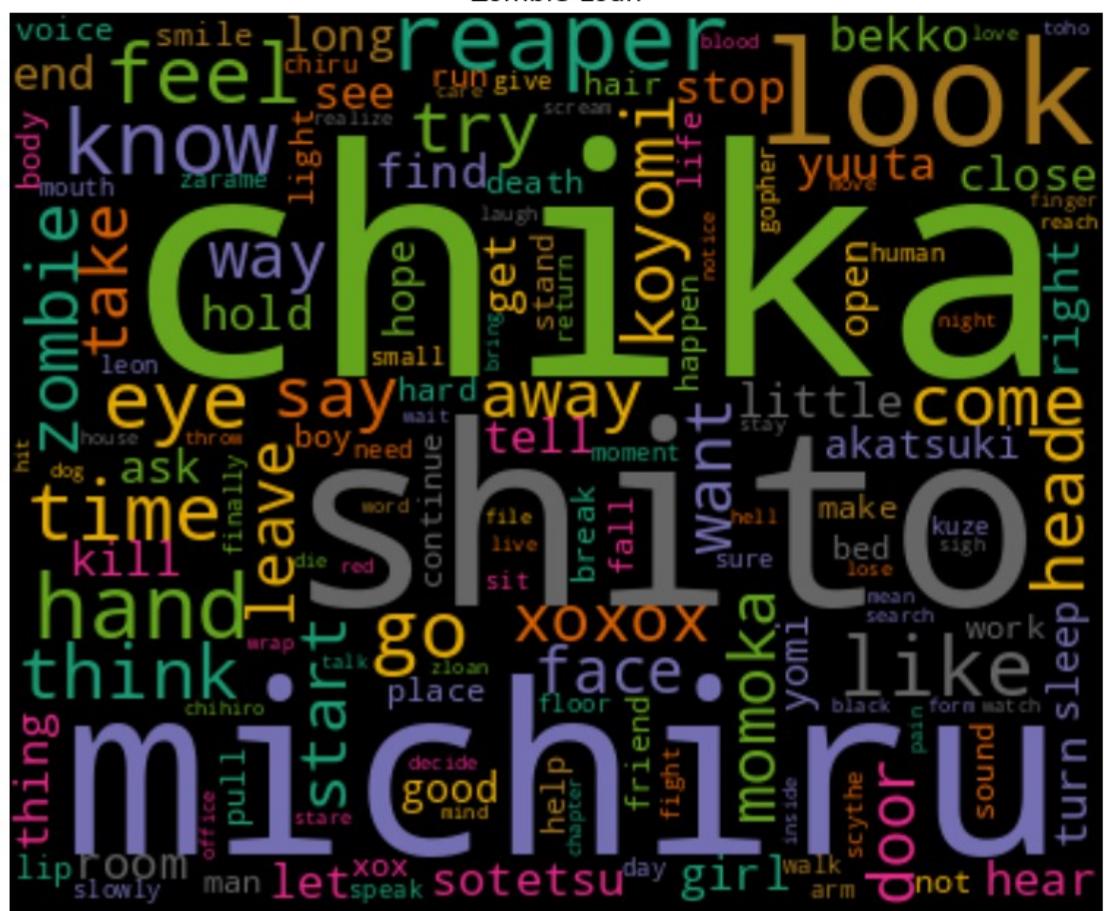


Zombie Prom

A dense word cloud composed of numerous Korean words and their English translations. The words are arranged in a grid-like pattern, with some words appearing multiple times in different colors. The words include: knife (knife), tell (create), go (dose), reply (space), want (condition), stay (let), junhong (sob), sob (unconscious), tell (create), suddenly (suddenly), kill (kill), mix (mix), mean (mean), deep (deep), fast (fast), stand (stand), jean (jean), life (life), voice (voice), file (file), person (person), wake (wake), walk (walk), start (start), like (type), explain (explain), hospital (hospital), stand (stand), wait (wait), lay (lay), long (long), this (this), tear (tear), talk (talk), slowly (slowly), bleed (bleed), stick (stick), huge (huge), downstair (downstairs), character (character), tonight (tonight), black (black), accident (accident), they (they), end (end), need (need), happen (happen), pain (pain), feel (feel), see (see), more (more), against (against), leave (leave), heart (heart), mind (mind), open (open), until (until), stop (stop), time (time), later (later), night (night), hell (hell), girl (girl), mind (mind), minute (minute), open (open), until (until), surgeon (surgeon), promise (promise), schedule (schedule), brown (brown), hair (hair), deserve (deserve), cry (cry), camera (camera), shock (shock), rest (rest), thing (thing), way (way), get (get), say (say), miss (miss), exactly (exactly), zelshan (zelshan), saw (saw), vintage (vintage), sayou (sayou), kid (kid), know (know), janghang (janghang), blood (blood), fuck (fuck), honest (honest), love (love), ass (ass), happen (happen), hold (hold), come (come), ghost (ghost), silent (silent), surger (surger), protect (protect), that (that), the (the), sure (sure), zeloshane (zeloshane), upi (upi), tomorrow (tomorrow), happen (happen), girlfriend (girlfriend), jagiyayou (jagiyayou), smile (smile).

## Zombie Survival Guide

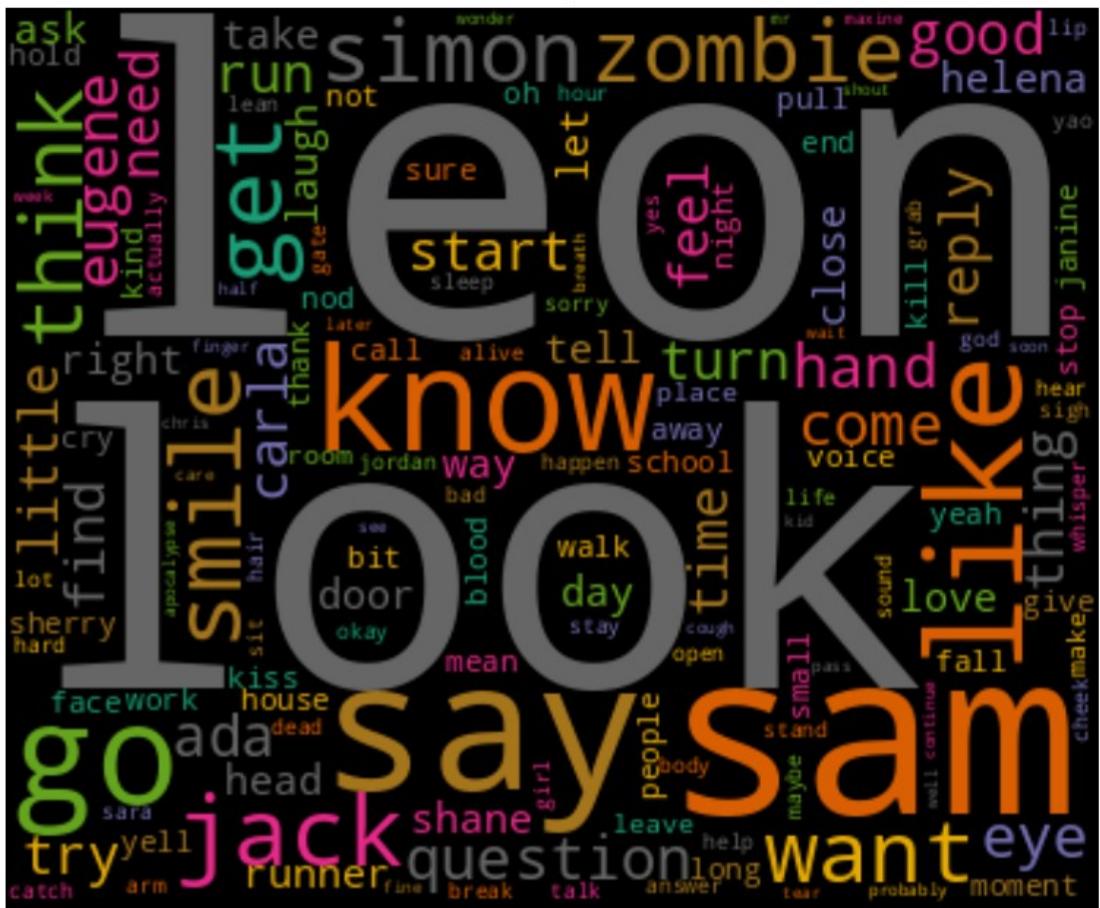
Zombie-Loan



Zombieland

mind stand story  
get hold young SURE love  
keep feel life big close face  
hope watch road lie open  
idea kid hand. meet hear  
tell door hand. drive laugh  
lot voice hard bit  
time stop second tal room eye  
help hard sorry fire  
make ya want world let right  
old well long rock mom start man  
well shake chapter run wait maybe  
make fall people shoot follow  
body day branson catch break  
try turn dead way zombielink  
away soon light give gun hit  
bad happen reach talk care  
rule Wichita mean place new see  
thing work head find ask hair  
tallahassee arm buck take sit

## Zombies, Run!



## Zone of The Enders



## Zoo Tycoon

Zoolander

A word cloud centered around the characters from the Zoolander commercial. The words are arranged in a grid-like structure, with larger words being the primary focus and smaller words appearing as secondary or tertiary elements.

**Derek:** proud, wrong, flow, right, der, young, know, son, bra jr, file, lip, new, wow, chill, admire, think.

**Hansel:** step, want, speechless, time, momentarily, end.

**Gertie:** strike, look, hair, pout, quickly, sec, finish, father.

**Others:** work, sheya, scraggly, hand, present, feel, relax, hanstupid.

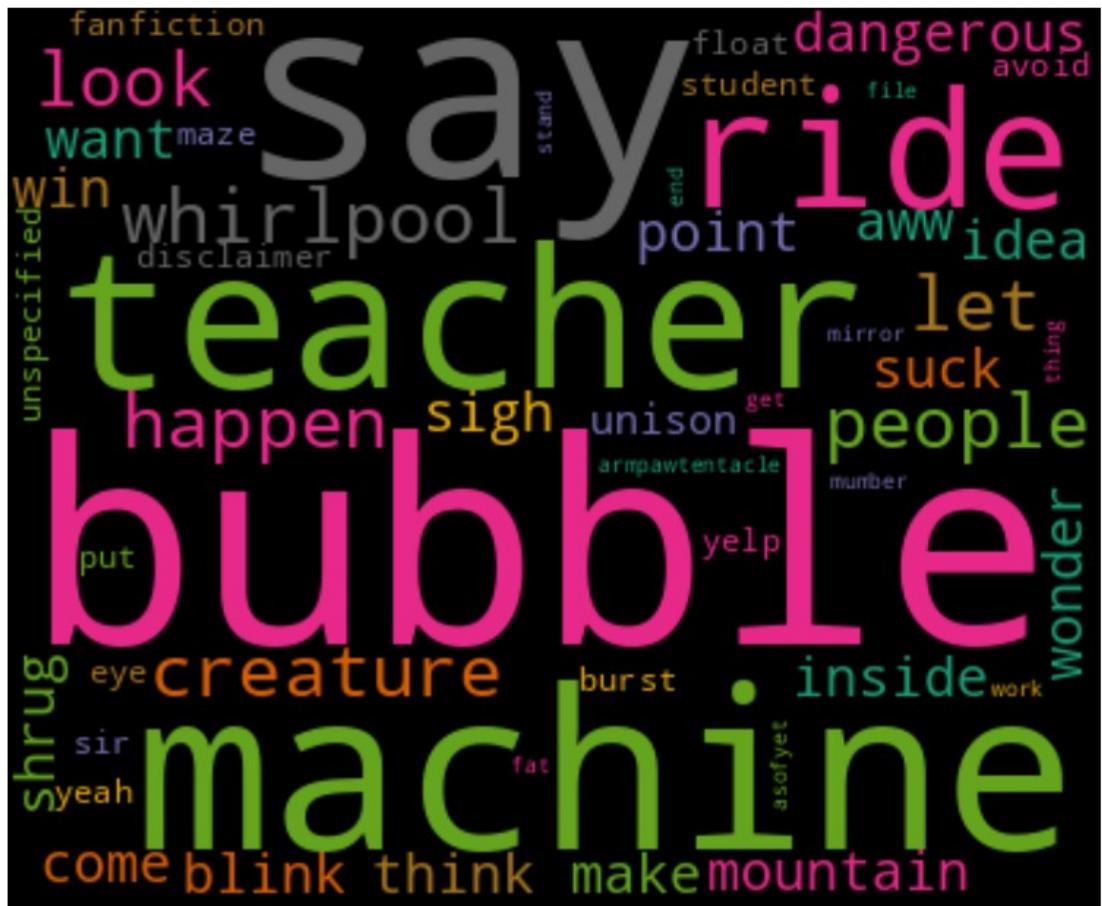
The background is black, and the words are in various colors including green, yellow, grey, pink, and purple.

## Zoom



Zombie Blondes





Zoop





```
# Saving the vectorized data to a csv
dtm_file_path = 'dtm_data.csv'
cleaned_data_path = 'cleaned_data.csv'

# Save DataFrame to CSV
df_dtm.to_csv(dtm_file_path, index=False)
df.to_csv(cleaned_data_path, index=False)
```