

In [1]:

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6
7 warnings.filterwarnings("ignore")

```

In [2]:

```

1 movies = pd.read_csv('movies.csv')
2 credits = pd.read_csv('credits.csv')

```

In [3]:

```
1 movies.shape, credits.shape
```

Out[3]:

```
((4803, 20), (4803, 4))
```

In [4]:

```
1 movies.head(2)
```

Out[4]:

	budget	genres	homepage	id	keywords	original_
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id":...	
1	300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "...	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "na...	

In [5]:

```
1 movies["genres"][0]
```

Out[5]:

```
'[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'
```

In [6]:

```
1 import json
```

In [7]:

```
1 json.loads(movies["genres"][0])
```

Out[7]:

```
[{'id': 28, 'name': 'Action'},  
 {'id': 12, 'name': 'Adventure'},  
 {'id': 14, 'name': 'Fantasy'},  
 {'id': 878, 'name': 'Science Fiction'}]
```

In [8]:

```
1 json.loads(movies["production_companies"][0])
```

Out[8]:

```
[{'name': 'Ingenious Film Partners', 'id': 289},  
 {'name': 'Twentieth Century Fox Film Corporation', 'id': 306},  
 {'name': 'Dune Entertainment', 'id': 444},  
 {'name': 'Lightstorm Entertainment', 'id': 574}]
```

In [9]:

```
1 json.loads(movies["production_countries"][0])
```

Out[9]:

```
[{'iso_3166_1': 'US', 'name': 'United States of America'},  
 {'iso_3166_1': 'GB', 'name': 'United Kingdom'}]
```

In [10]:

```
1 credits.head(2)
```

Out[10]:

	movie_id	title	cast	crew
0	19995	Avatar	[{"cast_id": 242, "character": "Jake Sully", "...	[{"credit_id": "52fe48009251416c750aca23", "de...
1	285	Pirates of the Caribbean: At World's End	[{"cast_id": 4, "character": "Captain Jack Spa...	[{"credit_id": "52fe4232c3a36847f800b579", "de...

In [11]:



```
1 json.loads(credits["cast"][0])[:2]
```

Out[11]:

```
[{'cast_id': 242,
  'character': 'Jake Sully',
  'credit_id': '5602a8a7c3a3685532001c9a',
  'gender': 2,
  'id': 65731,
  'name': 'Sam Worthington',
  'order': 0},
 {'cast_id': 3,
  'character': 'Neytiri',
  'credit_id': '52fe48009251416c750ac9cb',
  'gender': 1,
  'id': 8691,
  'name': 'Zoe Saldana',
  'order': 1}]
```

In [12]:



```
1 json.loads(credits["crew"][0])[:2]
```

Out[12]:

```
[{'credit_id': '52fe48009251416c750aca23',
  'department': 'Editing',
  'gender': 0,
  'id': 1721,
  'job': 'Editor',
  'name': 'Stephen E. Rivkin'},
 {'credit_id': '539c47ecc3a36810e3001f87',
  'department': 'Art',
  'gender': 2,
  'id': 496,
  'job': 'Production Design',
  'name': 'Rick Carter'}]
```

In [13]:



```
1 # Data Merging
```

In [14]:



```
1 movie_credits = pd.merge(movies,credits, left_on="id" ,right_on = "movie_id")
2 movie_credits.shape
```

Out[14]:

```
(4803, 24)
```

In [15]:

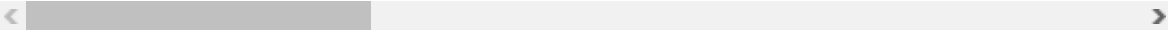


```
1 movie_credits.head(2)
```

Out[15]:

	budget	genres	homepage	id	keywords	original_
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": 1464, "name": "culture clash"}]	
1	300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "name": "ocean"}]	

2 rows × 24 columns



In [16]:



```
1 movie_credits.columns
```

Out[16]:

```
Index(['budget', 'genres', 'homepage', 'id', 'keywords', 'original_language',
      'original_title', 'overview', 'popularity', 'production_companies',
      'production_countries', 'release_date', 'revenue', 'runtime',
      'spoken_languages', 'status', 'tagline', 'title_x', 'vote_average',
      'vote_count', 'movie_id', 'title_y', 'cast', 'crew'],
      dtype='object')
```

In [17]:



```
1 movie_credits["release_date"] = pd.to_datetime(movie_credits["release_date"])
2 movie_credits["release_date"].head()
```

Out[17]:

```
0    2009-12-10
1    2007-05-19
2    2015-10-26
3    2012-07-16
4    2012-03-07
Name: release_date, dtype: datetime64[ns]
```

In [18]:



```
1 release_year = movie_credits["release_date"].dt.year.value_counts(dropna = False)
```

In [19]:



```
1 # need to filter the data so that you only retain records for years where there were
2 release_year = release_year[release_year >= 100].sort_index()
3 release_year.index = release_year.index.astype(int)
4 release_year
```

Out[19]:

1997	112
1998	133
1999	171
2000	166
2001	183
2002	203
2003	169
2004	204
2005	217
2006	237
2007	195
2008	227
2009	247
2010	225
2011	223
2012	208
2013	231
2014	238
2015	216
2016	104

Name: release\_date, dtype: int64

In [20]:

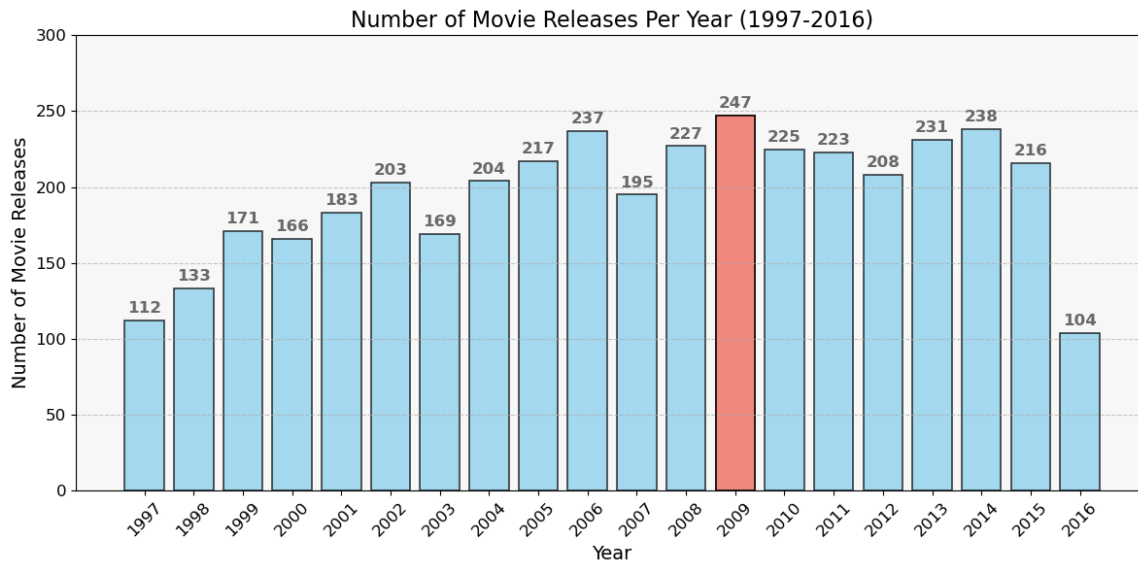


```
1 years = release_year.index.to_list()
2 movie_counts = release_year.values.tolist()
```

In [21]:



```
1 import matplotlib.pyplot as plt
2
3 # Find the year with the highest movie releases
4 max_count_year = years[movie_counts.index(max(movie_counts))]
5
6 # Create a stylish bar chart
7 plt.figure(figsize=(12, 6))
8 plt.bar(years, movie_counts, color='skyblue', alpha=0.75, edgecolor='black', linewidth=1)
9 plt.xlabel('Year', fontsize=14)
10 plt.ylabel('Number of Movie Releases', fontsize=14)
11 plt.title('Number of Movie Releases Per Year (1997-2016)', fontsize=16)
12 plt.xticks(rotation=45, fontsize=12)
13 plt.yticks(fontsize=12)
14 plt.grid(axis='y', linestyle='--', alpha=0.7)
15 plt.grid(axis='x', linestyle='')
16
17 # Highlight the bar for the year with the highest movie releases
18 highlighted_bar = years.index(max_count_year)
19 plt.bar(years[highlighted_bar], movie_counts[highlighted_bar], color='salmon', alpha=0.75)
20
21 # customised y-ticks range
22 plt.yticks(range(0, 301, 50), fontsize=12)
23
24
25 # Adding data labels
26 for year, count in zip(years, movie_counts):
27     plt.text(year, count + 5, str(count), ha='center', fontsize=12, fontweight='bold')
28
29
30 # Customize tick labels
31 plt.xticks(years, rotation=45, fontsize=12)
32 plt.yticks(fontsize=12)
33
34
35 # Adding a background color
36 ax = plt.gca()
37 ax.set_facecolor('#f7f7f7')
38
39 # Display the chart with improved aesthetics
40 plt.tight_layout()
41 plt.show()
```



In [22]:

```
1 movies["popularity"].agg(func=["min", "max", "mean"])
```

Out[22]:

```
min      0.000000
max     875.581305
mean     21.492301
Name: popularity, dtype: float64
```

## Budget vs. Revenue Scatter Plot:

Create a scatter plot to visualize the relationship between budget and revenue for the movies. This can help identify if there's a correlation between the two.

In [23]:

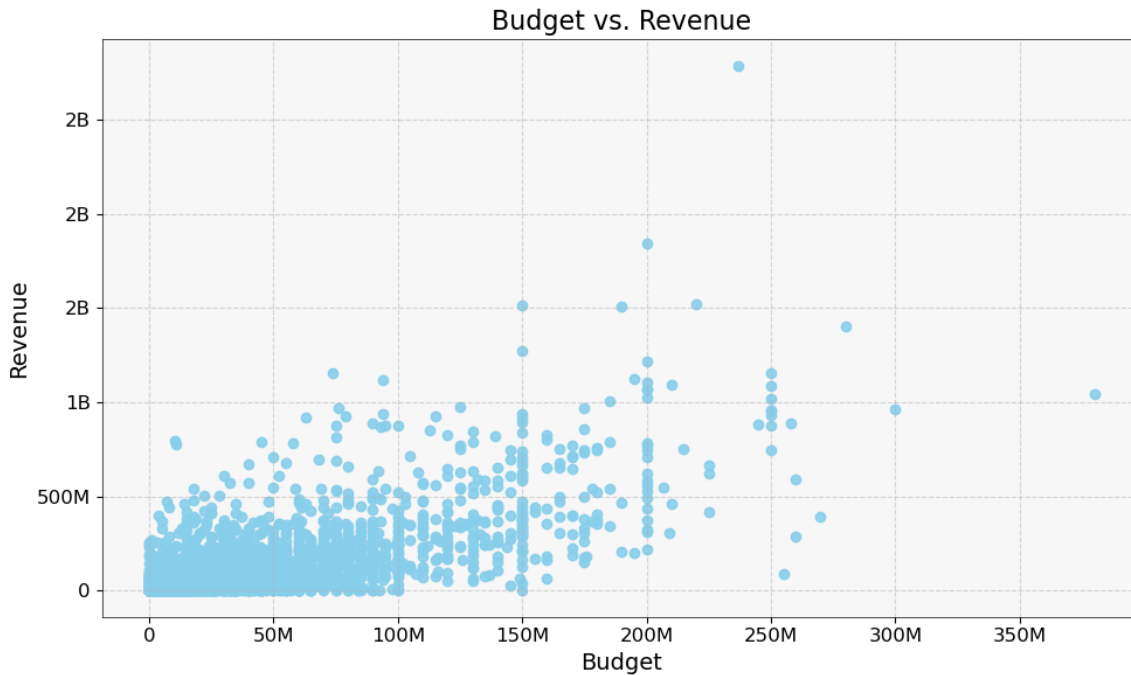
```
1 budget = movie_credits["budget"]
2 revenue = movie_credits["revenue"]
```

In [24]:



```
1 # Function to format budget values as alphanumeric units
2 def format_budget(value, pos):
3     if value >= 1e9:
4         return f'{value / 1e9:.0f}B'
5     elif value >= 1e6:
6         return f'{value / 1e6:.0f}M'
7     elif value >= 1e3:
8         return f'{value / 1e3:.0f}K'
9     else:
10        return f'{value:.0f}'
11
12
13 # Create a scatter plot
14 plt.figure(figsize=(10, 6))
15 plt.scatter(budget, revenue, color='skyblue', alpha=0.9)
16 plt.xlabel("Budget", fontsize=14)
17 plt.ylabel("Revenue", fontsize=14)
18 plt.title("Budget vs. Revenue", fontsize=16)
19
20 # Customize the appearance
21 plt.grid(True, linestyle='--', alpha=0.5)
22 plt.xticks(fontsize=12)
23 plt.yticks(fontsize=12)
24 plt.gca().spines['top'].set_linewidth(0.5)
25 plt.gca().spines['right'].set_linewidth(0.5)
26 plt.gca().spines['bottom'].set_linewidth(0.5)
27 plt.gca().spines['left'].set_linewidth(0.5)
28
29 # Add a background color
30 plt.gca().set_facecolor('#f7f7f7')
31
32 # Update x-axis and y-axis tick labels using the format_budget function
33 plt.gca().get_xaxis().set_major_formatter(plt.FuncFormatter(format_budget))
34 plt.gca().get_yaxis().set_major_formatter(plt.FuncFormatter(format_budget))
35
36
37 plt.tight_layout()
38 plt.show()
39
```





## Genre Distribution Bar Chart:

Create a bar chart to show the distribution of movie genres. You can count the number of movies in each genre category and display it as a bar chart.

In [25]:

```
1 movie_credits["genres"][0]
```

Out[25]:

```
'[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'
```

In [26]:

```
1 genre_data = movie_credits["genres"]
2 genre_data
```

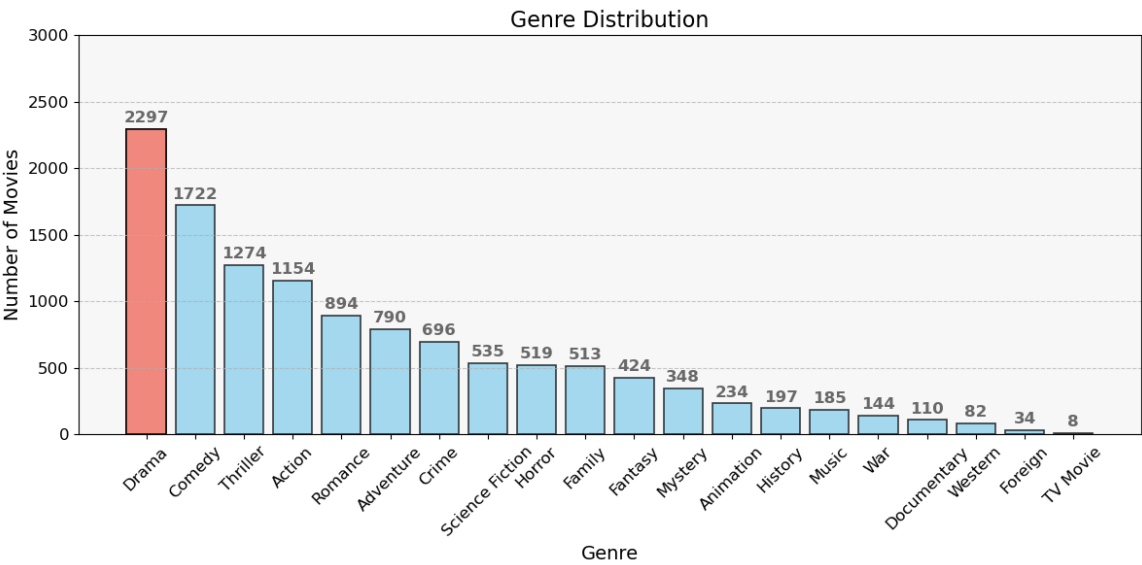
Out[26]:

```
0 [{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
1 [{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
2 [{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
3 [{"id": 28, "name": "Action"}, {"id": 80, "name": "Documentary"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
4 [{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
...
4798 [{"id": 28, "name": "Action"}, {"id": 80, "name": "Documentary"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
4799 [{"id": 35, "name": "Comedy"}, {"id": 10749, "name": "Drama"}, {"id": 18, "name": "Romance"}, {"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
4800 [{"id": 35, "name": "Comedy"}, {"id": 18, "name": "Romance"}, {"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]
4801 []
4802 [{"id": 99, "name": "Documentary"}]
Name: genres, Length: 4803, dtype: object
```

In [27]:



```
1 # Create a DataFrame from the sample data
2 df = pd.DataFrame({'genres': genre_data})
3
4 # Function to extract genre names from JSON data
5 def extract_genres(genre_json):
6     genres = json.loads(genre_json)
7     return [genre['name'] for genre in genres]
8
9 # Apply the extract_genres function to each row in the DataFrame
10 df['genre_names'] = df['genres'].apply(extract_genres)
11
12 # Flatten the list of genre names
13 all_genres = [genre for genres in df['genre_names'] for genre in genres]
14
15 # Count the occurrences of each genre
16 genre_counts = pd.Series(all_genres).value_counts()
17
18 # Find the most common genre
19 most_common_genre = genre_counts.idxmax()
20
21 # Create a stylish bar chart
22 plt.figure(figsize=(12, 6))
23 plt.bar(genre_counts.index, genre_counts.values, color='skyblue', alpha=0.75, edgecol
24 plt.xlabel('Genre', fontsize=14)
25 plt.ylabel('Number of Movies', fontsize=14)
26 plt.title('Genre Distribution', fontsize=16)
27 plt.xticks(rotation=45, fontsize=12)
28 plt.yticks(fontsize=12)
29 plt.grid(axis='y', linestyle='--', alpha=0.7)
30 plt.grid(axis='x', linestyle='')
31
32 # Highlight the most common genre
33 plt.bar(most_common_genre, genre_counts[most_common_genre], color='salmon', alpha=0.9
34
35 # customised y-ticks range
36 plt.yticks(range(0, 3001, 500), fontsize=12)
37
38 # Adding data labels
39 for genre, count in zip(genre_counts.index, genre_counts.values):
40     plt.text(genre, count + 50, str(count), ha='center', fontsize=12, fontweight='bo
41
42 # Customize tick labels
43 plt.xticks(fontsize=12)
44
45 # Adding a background color
46 ax = plt.gca()
47 ax.set_facecolor('#f7f7f7')
48
49 # Display the chart with improved aesthetics
50 plt.tight_layout()
51 plt.show()
```



Release Date Time Series:

Create a time series plot to show how the number of movie releases has changed over time. You can group the data by year or month and plot the count of movies released.

In [28]:

```
1 movie_credits["release_date"]
```

Out[28]:

```
0      2009-12-10
1      2007-05-19
2      2015-10-26
3      2012-07-16
4      2012-03-07
...
4798   1992-09-04
4799   2011-12-26
4800   2013-10-13
4801   2012-05-03
4802   2005-08-05
Name: release_date, Length: 4803, dtype: datetime64[ns]
```

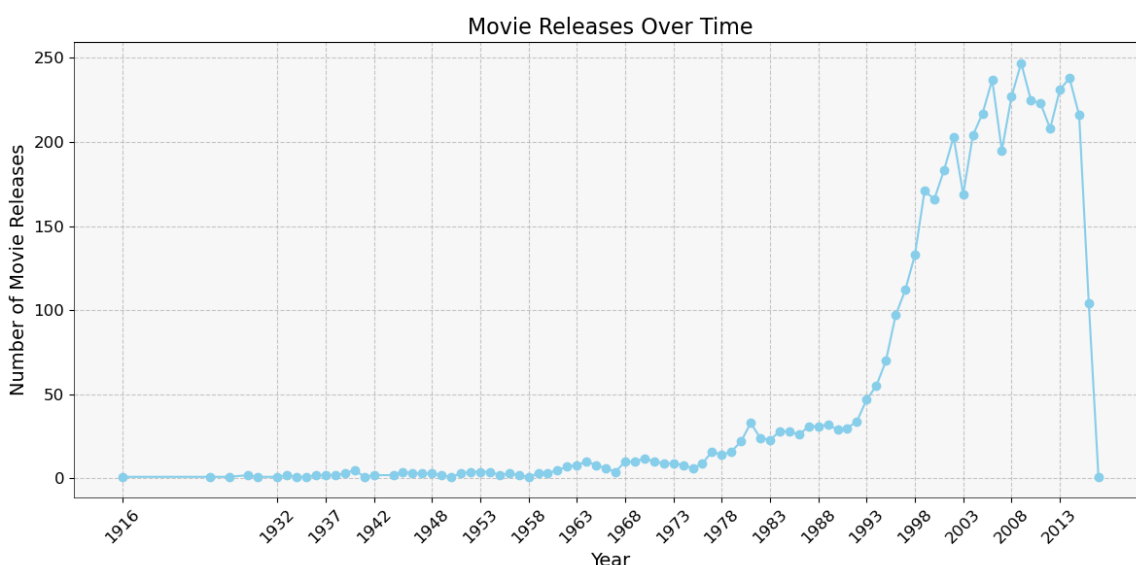
In [29]:



```

1 release_dates = movie_credits["release_date"]
2
3 # Extract the release years from the dates
4 release_years = release_dates.dt.year
5
6 # Count the number of movie releases per year
7 movie_count_per_year = release_years.value_counts().sort_index()
8
9 # Create a time series plot
10 plt.figure(figsize=(12, 6))
11 plt.plot(movie_count_per_year.index, movie_count_per_year.values, marker='o', color='
12 plt.xlabel('Year', fontsize=14)
13 plt.ylabel('Number of Movie Releases', fontsize=14)
14 plt.title('Movie Releases Over Time', fontsize=16)
15
16 # Customize the appearance
17 plt.xticks(rotation=45, fontsize=12)
18 plt.yticks(fontsize=12)
19 plt.grid(True, linestyle='--', alpha=0.7)
20 plt.gca().spines['top'].set_visible(0.5)
21 plt.gca().spines['right'].set_visible(0.5)
22 plt.gca().spines['bottom'].set_linewidth(0.5)
23 plt.gca().spines['left'].set_linewidth(0.5)
24
25 # Adding a background color
26 plt.gca().set_facecolor('#f7f7f7')
27
28 # Limit the number of x-axis ticks
29 plt.xticks(movie_count_per_year.index[::5]) # Show every 5 years
30
31 plt.tight_layout()
32 plt.show()
33

```



## Language Distribution Pie Chart:

Visualize the distribution of original languages in your dataset using a pie chart. It can show which languages are most common.

In [30]:

```
1 df = movie_credits[["spoken_languages"]]
2 df.head()
```

Out[30]:

	spoken_languages
0	[{"iso_639_1": "en", "name": "English"}, {"iso...
1	[{"iso_639_1": "en", "name": "English"}]
2	[{"iso_639_1": "fr", "name": "Fran\u00e7ais"},...
3	[{"iso_639_1": "en", "name": "English"}]
4	[{"iso_639_1": "en", "name": "English"}]

In [31]:

```
1 # Function to extract language names from JSON data
2 def extract_languages(language_json):
3     languages = json.loads(language_json)
4     return [lang['name'] for lang in languages]
```

In [32]:

```
1 all_languages = df['spoken_languages'].apply(extract_languages)
2 all_languages.head()
```

Out[32]:

```
0      [English, Espa\u00f1ol]
1      [English]
2  [Fran\u00e7ais, English, Espa\u00f1ol, Italiano, Deutsch]
3      [English]
4      [English]
Name: spoken_languages, dtype: object
```

In [33]:

```
1 final_languages = []
2
3 for i in all_languages.values:
4     final_languages.extend(i)
```

In [34]:



```
1 final_languages[:10]
```

Out[34]:

```
['English',  
 'Español',  
 'English',  
 'Français',  
 'English',  
 'Español',  
 'Italiano',  
 'Deutsch',  
 'English',  
 'English']
```

In [35]:



```
1 # Count the occurrences of each Language  
2 language_counts = pd.Series(final_languages).value_counts()  
3  
4 # Filter Languages with more than 100 occurrences  
5 language_counts = language_counts[language_counts > 100]  
6 language_counts
```

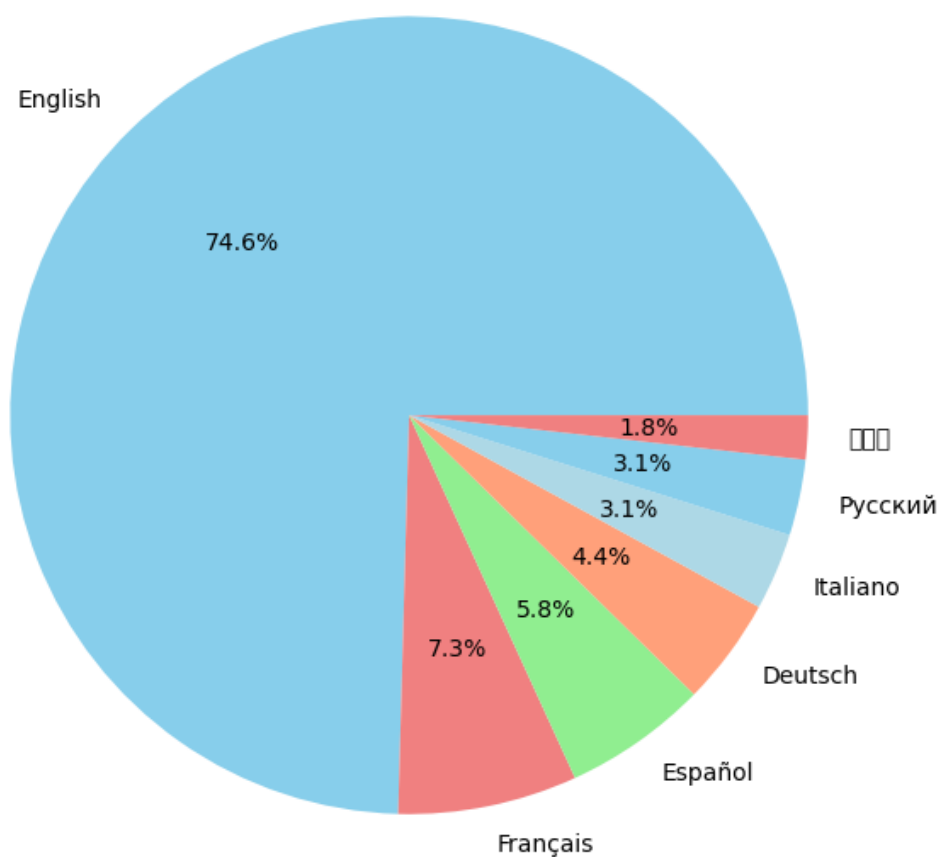
Out[35]:

```
English      4485  
Français    437  
Español      351  
Deutsch      262  
Italiano     188  
Русский      185  
普通话       107  
dtype: int64
```

In [36]:

```
1 # Create a pie chart
2 plt.figure(figsize=(6, 6))
3 plt.pie(language_counts, labels=language_counts.index, autopct='%1.1f%%', colors=['sk
4 plt.title('Language Distribution (Languages with >100 Occurrences)', fontsize=16)
5
6 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
7
8 plt.tight_layout()
9 plt.show()
```

## Language Distribution (Languages with >100 Occurrences)



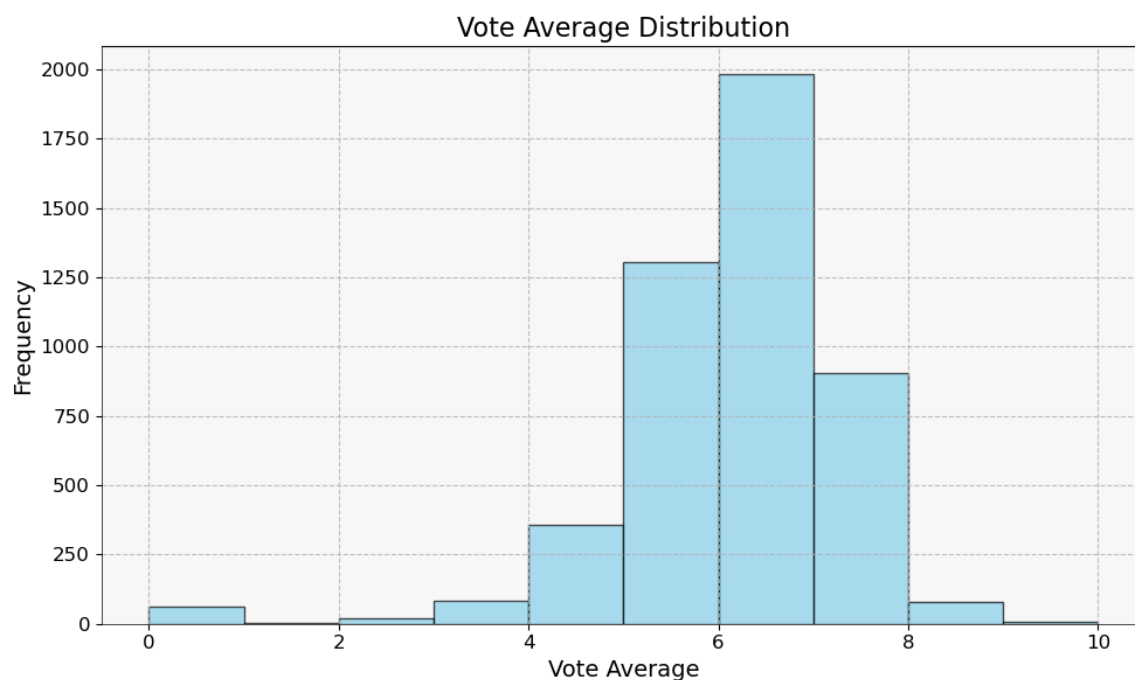
## Vote Average Distribution Histogram:

Create a histogram to show the distribution of vote averages for movies. This can help you understand the general sentiment or rating distribution of the movies.

In [37]:



```
1 vote_averages = movie_credits["vote_average"].values
2
3 # Create a histogram
4 plt.figure(figsize=(10, 6))
5 plt.hist(vote_averages, bins=10, color='skyblue', edgecolor='black', alpha=0.7)
6 plt.xlabel("Vote Average", fontsize=14)
7 plt.ylabel("Frequency", fontsize=14)
8 plt.title("Vote Average Distribution", fontsize=16)
9
10 # Customize the appearance
11 plt.grid(True, linestyle='--', alpha=0.8)
12 plt.xticks(fontsize=12)
13 plt.yticks(fontsize=12)
14 plt.gca().spines['top'].set_visible(0.5)
15 plt.gca().spines['right'].set_visible(0.5)
16 plt.gca().spines['bottom'].set_linewidth(0.5)
17 plt.gca().spines['left'].set_linewidth(0.5)
18
19 # Add a background color
20 plt.gca().set_facecolor('#f7f7f7')
21
22 plt.tight_layout()
23 plt.show()
```



## Tagline Word Cloud:

Generate a word cloud from movie taglines to visualize the most common words or phrases used in movie marketing.





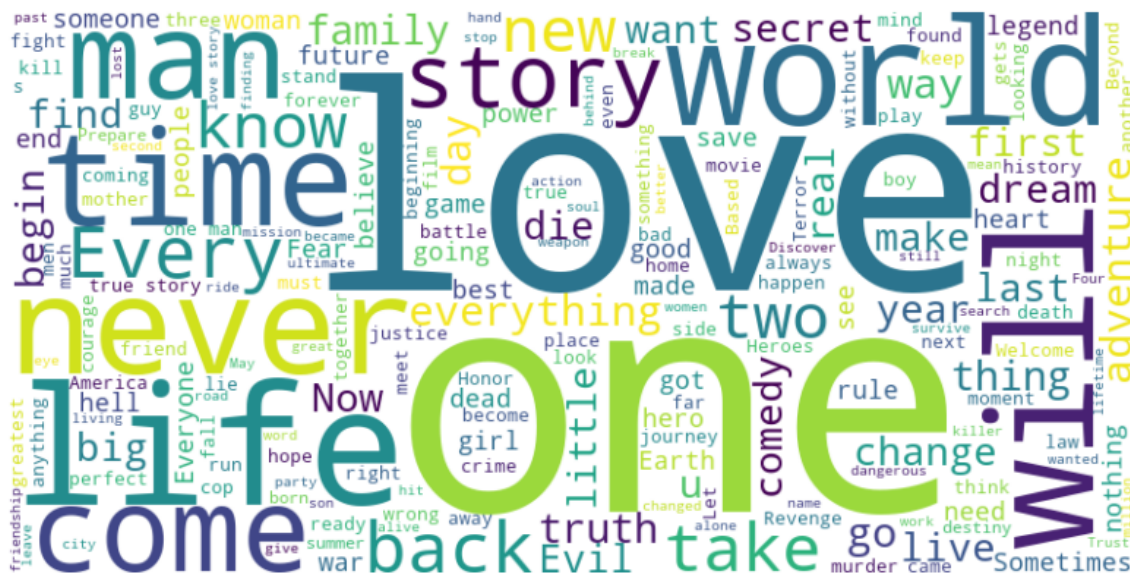
```
1 movie_credits["tagline"].dropna()
```

```
0          Enter the World of Pandora.
1      At the end of the world, the adventure begins.
2          A Plan No One Escapes
3          The Legend Ends
4      Lost in our world, found in another.

...
4795      Sometimes you've got to break the rules
4796      What happens if it actually works?
4798      He didn't come looking for trouble, but troubl...
4799      A newlywed couple's honeymoon is upended by th...
4801      A New Yorker in Shanghai
Name: tagline, Length: 3959, dtype: object
```



```
1 from wordcloud import WordCloud
2
3
4 taglines = movie_credits["tagline"].dropna()
5
6 # Combine all taglines into a single text
7 taglines_text = " ".join(taglines)
8
9 # Generate the word cloud
10 wordcloud = WordCloud(width=800, height=400, background_color='white').generate(taglines_text)
11
12 # Display the word cloud
13 plt.figure(figsize=(10, 5))
14 plt.imshow(wordcloud, interpolation='bilinear')
15 plt.axis('off') # Turn off the axis labels
16
17 plt.show()
18
```



In [40]:



```
1 taglines[taglines.str.contains("love")]
```

Out[40]:

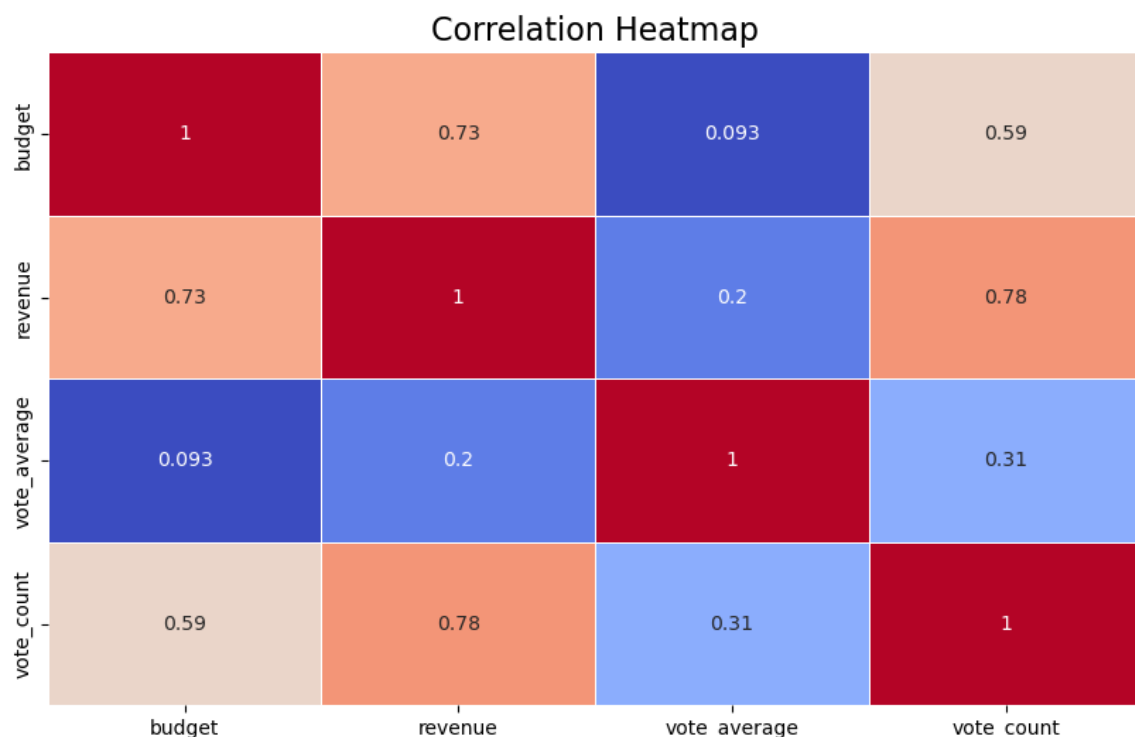
```
11      For love, for hate, for justice, for revenge.
110    It takes a moment to change history. It takes ...
124    Only the act of true love will thaw a frozen h...
145    For passion. For honor. For destiny. For victo...
164      The faces you love. The action you expect.
      ...
4567      There are two sides to every love story.
4603      Sometimes, love really is a bitch.
4631      It's not who you love. It's how.
4655      Everybody needs some bunny to love.
4696    A (sort of) love story between two guys over a...
Name: tagline, Length: 172, dtype: object
```

## Correlation Heatmap:

Create a heatmap to visualize the correlations between numerical columns like budget, revenue, vote average, and vote count. This can help identify strong relationships between variables.

In [41]:

```
1 df = movie_credits[["budget", "revenue", "vote_average", "vote_count"]]
2
3 # Create a correlation matrix
4 correlation_matrix = df.corr()
5
6 # Create a heatmap
7 plt.figure(figsize=(10, 6))
8 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', linewidths=0.5, cbar = False)
9 plt.title('Correlation Heatmap', fontsize=16)
10 plt.show()
11
```



15. **Overview Word Cloud:** Generate a word cloud from movie overviews to visualize the most common words or themes in movie descriptions.

In [42]:

```
1 movie_credits["overview"].isnull().sum()
```

Out[42]:

3



In [ ]:



1	
---	--