# ⌄ Case Study Part 2: Insights

This notebook is running on a Google Colab L4 CPU environment. In this section, the metrics inferred from the data in part 1 of the case study are clustered, analyzed and visualized.

```
1 %%capture
2 #!pip install numpy --upgrade #(uncomment install lines if running first time in google colab)
3 #!pip install scipy --upgrade
4 #!pip install gensim #(refresh kernel after installing to load new library and recomment installs)
```

```
 1 %%capture
 2 import pandas as pd
 3 import gensim #gensim is used for a latent text embedding model for clustering text by meaning
 4 from gensim.models import LdaMulticore
 5 from gensim.corpora import Dictionary
 6 from gensim.models.coherencemodel import CoherenceModel
 7 from matplotlib import pyplot as plt
 8 import matplotlib.patches as mpatches
 9 !pip install wordcloud
10 from wordcloud import WordCloud
```

```
1 from gensim.models.phrases import Phrases
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 #from sklearn.cluster import KMeans
4 #from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score, silhouette_score
5 import numpy as np
```

```
1 %%capture
2 !pip install datasets
3 from datasets import load_dataset
4 from google.colab import userdata
```

```
1 data = load_dataset("skeskinen/TinyStories-GPT4")
```

```
1 data_df = data['train'].to_pandas()
```

```
1 word_occurrences = pd.read_csv('word_occurrences.csv')
```

```
1 feature_occurrences = pd.read_csv('feature_occurrences.csv')
```

# ⌄ Statistical Analysis and Visualization

In this section the statistical metrics calculated from the training data is analyzed and visualized for making insights with regards to the stated objectives. Furthermore, given the complexity in understanding the level of literacy from the data - as no 'literacy level' column is within the data - a language model is used to classify the responses.

## ⌄ Mental Health

As for mental-health, goals may be defined in objective terms by inferring the emotions, subjects, topics, genres, or any commonly used keywords related to mental-health within the data. For this purpose, we can use the list of words and narrative features provided within the prompts provided to the AI in the training data for each row in the dataset. By counting the occurrences of the 'words' and 'features' columns used within the prompt, we can make some insights of the mental-health objectives.

Objective: Highlight narrative features that relate to mental health.
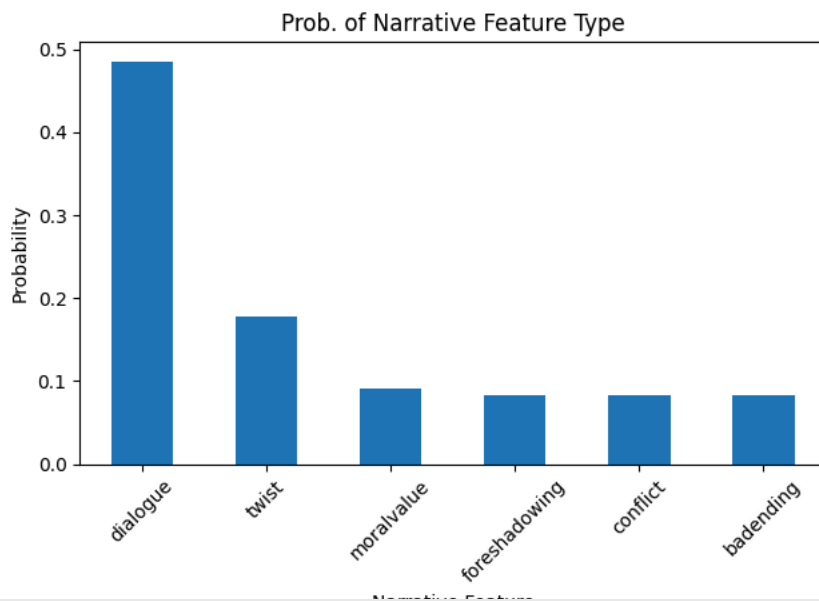
### ⌄ Narrative Features

Some of the 'feature narrative' words provided to the model within the prompt template's in the training dataset give insight into whether the story may contain any emotional subjects or topics related to mental health. A bar plot is shown that highlights the amount of each feature

narrative.

```
1 # converting the counts of feature occurrences to probability values
2 feature_occurence_probabilities = feature_occurrences['feature_occurrences'] / feature_occurrences['feature_occurrences'].su
```

```
1 feature_occurence_probabilities.index = feature_occurrences['feature']
```

```
1 ax = feature_occurence_probabilities.sort_values(ascending=False).plot(
2     kind='bar',
3     title='Prob. of Narrative Feature Type'
4 )
5
6 # Correctly set labels using set methods
7 plt.ylabel('Probability')
8 plt.xlabel('Narrative Feature')
9 plt.xticks(rotation=45)
10
11 # Optional: Adjust layout to prevent label cutoff
12 plt.tight_layout()
13
14 plt.show()
```



```
1 feature_occurrences.describe()
```

|  | feature_occurrences |
|---|---|
| count | 6.000000e+00 |
| mean | 5.059540e+05 |
| std | 4.859297e+05 |
| min | 2.503000e+05 |
| 25% | 2.507192e+05 |
| 50% | 2.624705e+05 |
| 75% | 4.730752e+05 |

**Mental Health Insight**

We can see from the bar chart that the majority of narratives have dialogue (i.e speech between characters in a story). Similarly, 'twists' and 'moralvalue' are 2nd and 3rd respectively, highlighting how the majority of prompts seek to generate interesting stories where characters interact and exemplify moral values. On the other end of the spectrum, the story prompts in the training data generate stories with narratives

like foreshadowing, conflict, and bad endings. As our target demographic is youth, it might be a good idea to do further investigation on story prompts that contain the 'conflict' and 'badending' narrative tags - these combined tags are contained in roughly **20%** of the prompts.

**Statistical Insight**

The description of the feature occurrences show that the average count of prompts with each narrative type is typically ~500,000 for each. However, the 50% quantile is ~260,000 showing that the majority of narrative types are roughly half as large as the mean. This implies that there are outlier narratives skewing the distribution up. The 'dialogue' and 'twist' narrative types are outliers, and are contained within ~50% and ~20% of the story prompts respectively, and combined contribute to ~70% of the story prompts. Similarly, the standard deviation value of ~480,000 shows how the variation in prompt counts is really high, further showing how the outliers ('dialogue', 'twist') are overshadowing the remaining narrative types ('moralvalue', 'foreshadowing', 'conflict', 'badending').

**Note**: Most of the prompts range in number of narrative tags, and may contain either none or all of the tags described above. As such, it is worth investigating how these tags correlate. For example, a combined 'badending' and 'conflict' narrative lacking a 'moralvalue' tag might lead to risky or insensitive storytelling.

## ⌄  Key Words

The key words column 'words' in the data are specific key vocabulary words provided to the AI in the story generation prompt. These words differ from the narrative feature words as the key words must be included in the generated story - whereas the narrative features are abstract storytelling features.

```
1 # converting the counts of feature occurrences to probability values
2 word_occurrence_probabilities = word_occurrences['word_occurrences'] / word_occurrences['word_occurrences'].sum()
```

```
1 word_occurrence_probabilities.index = word_occurrences['word']
```

```
1 word_occurrence_probabilities.describe()
```

|  | word_occurrences |
|---|---|
| count | 1603.000000 |
| mean | 0.000624 |
| std | 0.000422 |
| min | 0.000294 |
| 25% | 0.000311 |
| 50% | 0.000319 |
| 75% | 0.000854 |
| max | 0.002497 |

As the collection of vocabulary words within the dataset is large (~1604), the data must be clustered along the words' sentiment/meaning using a text embedding model. This way, we can reduce the dimensionality of the data while somewhat preserving the intended meaning of the collection. Given the individual probabilities for each word are very small (on the scale of 1/1000 likeliness), clustering to reduce the dimensionality will help understand the salient groups of words that are commonly paired within each vocab key collection.

```
1 # step 1: prepare the n-grams
2 data_df['words'].head() # every collection of words is already of length 3 (trigram format)
```

|  | words |
|---|---|
| 0 | [receive, opera, red] |
| 1 | [use, sheet, blue] |
| 2 | [relax, bus, uncomfortable] |
| 3 | [sail, cricket, wide] |
| 4 | [pray, pigeon, creative] |

Since the dataset is large, clustering can not be executed in one execution - even with Google Colabs biggest ram environment. Therefore, the data will be clustered in iterations via bootstrapping. In this way, the data will be clustered in batches and then aggregated.
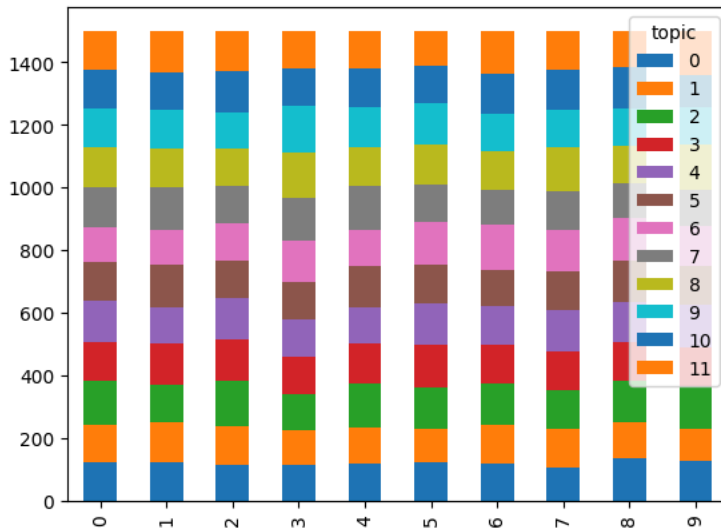
```python
1  # note: vectorizer must be initialized outside the loop
2  # in order to maintain the training data across batches
3
4  # Step 2: LDA Clustering the batched TF-IDF matrix with
5  def bootstrap_clustering(data, num_clusters, num_iterations, batch_size):
6      """
7      Performs bootstrapped clustering on the data, incorporating TF-IDF creation within the loop.
8
9      Args:
10         data: The data to cluster (the 'words' column of your DataFrame).
11         num_topics: The number of clusters in the LDA clustering.
12         num_iterations: The number of iterations.
13         batch_size: The batch size.
14
15     Returns:
16         A list of cluster assignments for each data point.
17     """
18
19     all_topic_assignments = []
20     batches = []
21     topic_assignment_batches = []
22     # Initialize lists to store cluster statistical metrics
23     coherence_scores = []
24     perplexity_scores = []
25
26     data_size = len(data)
27
28     for _ in range(num_iterations):
29         # Randomly sample a batch of data
30         batch_indices = np.random.choice(data_size, size=batch_size, replace=True)
31         batch_data = data.iloc[batch_indices]['words']  # Select the 'words' column
32         batches.append(batch_data)
33
34         # Prepare data for LDA
35         dictionary = Dictionary([word for word in batch_data.values])
36         corpus = [dictionary.doc2bow(text) for text in batch_data.values]
37
38         # Perform LDA on the batch
39         lda_model = LdaMulticore(corpus, num_topics=num_topics, id2word=dictionary, passes=15, workers=2)
40         batch_topic_assignments = [max(lda_model[doc], key=lambda item: item[1])[0] for doc in corpus]
41
42         topic_assignment_batches.append(batch_topic_assignments)
43         # Append the cluster assignments to the list
44         all_topic_assignments.extend(batch_topic_assignments)
45         # Calculate and store metrics for this iteration
46         coherence_model = CoherenceModel(model=lda_model, texts=batch_data.values, dictionary=dictionary, coherence='c_v')
47         coherence_score = coherence_model.get_coherence()
48         coherence_scores.append(coherence_score)
49
50         perplexity = lda_model.log_perplexity(corpus)
51         perplexity_scores.append(perplexity)
52
53     return all_topic_assignments, batches, topic_assignment_batches, coherence_scores, perplexity_scores
```

```python
1  num_topics = 12 # 12 is chosen to be able to visualize wordcloud in grid format
2  num_iterations = 10  # Adjust as needed
3  batch_size = 1500  # This number is close to the maximum compute
4
5  # Perform bootstrapped clustering
6  topic_assignments, batches, topic_assignment_batches, coherence_scores, perplexity_scores = bootstrap_clustering(
7      data_df, num_topics, num_iterations, batch_size
8  )
9
10 # Assign the cluster assignments to the original DataFrame
11 # word_trigram_df['cluster'] = cluster_assignments
```

```python
1  cluster_assignments_df = pd.DataFrame(topic_assignments, columns=['topic'])
2  # adding in a cluster iteration index
3  cluster_assignments_df['iteration'] = cluster_assignments_df.index // batch_size
```

```
1 cluster_counts_by_iteration = cluster_assignments_df.groupby('iteration')['topic'].value_counts()
```

```
1 # let's visualize the distribution of cluster counts by iteration by showing a bar chart with a legend using the iteration
2 cluster_counts_by_iteration.unstack().plot(kind='bar', stacked=True, legend=True)
```

<Axes: xlabel='iteration'>



Since the topics discovered via LDA clustering are abstract, classification of these topics into a understandable theme may be done using a large language model with a prompt to classify into the existing objectives of 'mental health' or 'creativity'.

```
 1 mental_health_topics = [
 2     "Dealing with Fear and Anxiety",
 3     "Coping with Sadness and Loss",
 4     "Managing Anger and Frustration",
 5     "Practicing Gratitude and Mindfulness",
 6     "Developing Healthy Relationships with Family",
 7     "Developing Healthy Relationships with Friends"
 8 ]
 9
10 creativity_topics = [
11     "Thinking Outside the Box",
12     "Embracing Imagination and Curiosity",
13     "Using Art and Expression",
14     "Music",
15     "Painting",
16     "Sculpting"]
```
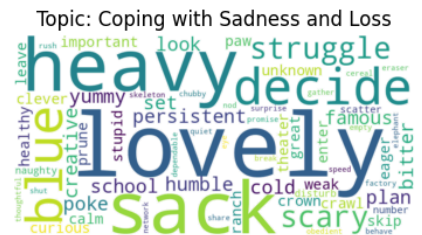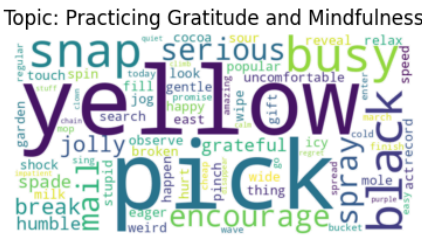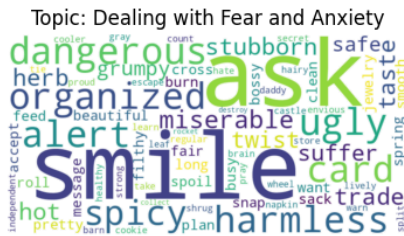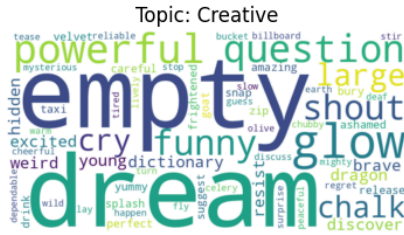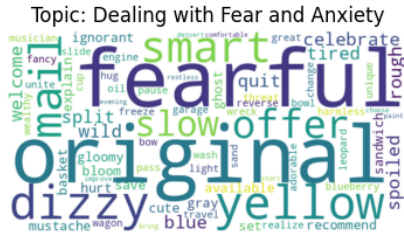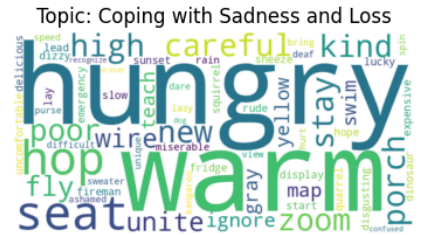
```
 1 # let's create a function to call openAI's chatgpt-4o with a prompt to classify topics
 2 import openai
 3 import os
 4
 5 client = openai.OpenAI(api_key=userdata.get('openai'))
 6 model = 'gpt-4o-mini'
 7 def classify_topics(mental_health_topics, creativity_topics, topic_df_words_counts):
 8     """ this takes a preset list of mental health and creativity topics,
 9         and classifies the topic_ids by using the data in the topic_df_words
10         note: there are 12 topic_ids, and 6 mental health topics, and 6 creativity topics
11         thus only one topic_id can be classified into one of the 2 categories
12     """
13     prompt_template = (
14         f"The following is a list of mental health topics: {mental_health_topics}",
15         f"The following is a list of creative topics: {creativity_topics}",
16         f"Use this list of words and their counts to classify the most likely topic: {topic_df_words_counts}",
17         "Note: Each topic in either list can be assigned only once. Only respond with the topic chosen and nothing else. The
18         "Topic: "
19     )
20     prompt = "\n".join(prompt_template)
21     response = client.chat.completions.create(
22         messages=[
23             {"role": "system", "content": "You are comfortable classifying the language used in stories by simply inferring
24             {"role": "user", "content": prompt}
```

```
25          ],
26          model = model,
27          temperature = 0.0
28      )
29      return response.choices[0].message.content
30
```

```
 1 # Step 3: creating a word cloud for all topics within a single iteration
 2 # add the cluster labels to the batch
 3 sample_batch = batches[0]
 4 batch_df = pd.DataFrame({'words': sample_batch, 'topic': topic_assignment_batches[0]})
 5 # Determine the number of rows and columns for the grid
 6 num_topics = len(batch_df['topic'].unique())
 7 num_cols = 3  # Adjust as needed
 8 num_rows = (num_topics + num_cols − 1) // num_cols
 9
10 # Create the figure and axes for the grid
11 fig, axes = plt.subplots(num_rows, num_cols, figsize=(15, 5 * num_rows))
12 i = 0
13 for topic_id in batch_df['topic'].unique():
14     topic_df = batch_df[batch_df['topic'] == topic_id].copy()
15
16     # Create a word cloud of the topics
17     wordcloud = WordCloud(width=800, height=400, background_color='white', max_words=100)
18     topic_df['words'] = topic_df['words'].apply(lambda x: ' '.join(x))
19     # Concatenate all word strings in topic_df['words']
20     wordcloud_text = ' '.join(topic_df['words'].values)
21     # Filter counts of all words in the wordcloud_text and then select top 100 with count > 1
22     wordcloud_text_counts = wordcloud.process_text(wordcloud_text)
23     wordcloud_text_counts = {k: v for k, v in sorted(wordcloud_text_counts.items(), key=lambda item: item[1], reverse=True)}
24     wordcloud_text_counts = dict(list(wordcloud_text_counts.items())[:100])
25     wordcloud_text_counts = pd.Series(wordcloud_text_counts, index=wordcloud_text_counts.keys())
26     # Filter to words with counts > 1
27     wordcloud_text_counts = wordcloud_text_counts[wordcloud_text_counts > 1]
28     wordcloud_text_filtered = ' '.join(wordcloud_text_counts.index)
29     # Classifying the topic_df['words']
30     topic = classify_topics(mental_health_topics, creativity_topics, wordcloud_text_counts.to_string())
31     topic_df['topic'] = topic
32     wordcloud.generate(wordcloud_text_filtered)
33     # Plot the wordcloud on the corresponding subplot
34     row = i // num_cols
35     col = i % num_cols
36     ax = axes[row, col]  # Get the current subplot
37     ax.imshow(wordcloud, interpolation='bilinear')
38     ax.axis('off')
39     ax.set_title(f'Topic: {topic}')  # Add a title/label
40     i += 1
41
```

Topic: Creative



Topic: Creative



Topic: Coping with Sadness and Loss



Topic: Dealing with Fear and Anxiety



Topic: Creative



Topic: Creative



Topic: Dealing with Fear and Anxiety



Topic: Mental Health



Topic: Practicing Gratitude and Mindfulness



Topic: Creative



Topic: Creative



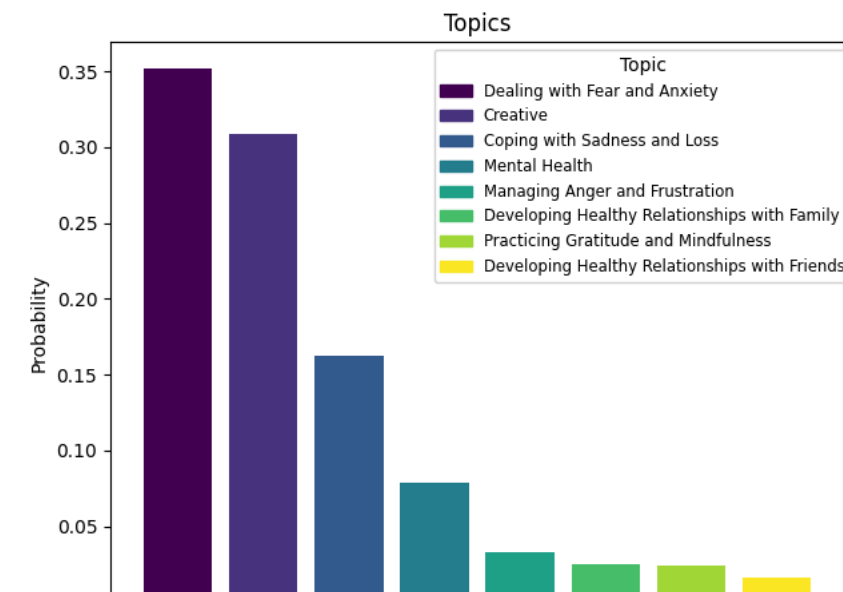Topic: Coping with Sadness and Loss

```python
1 # Step 4: classifying all the topic_ids across all bootstrap cluster iterations
2 all_classified_topic_batches = []
3 for i in range(len(topic_assignment_batches)):
4     batch_df = pd.DataFrame({'words': batches[i], 'topic': topic_assignment_batches[i]})
5     for topic_id in batch_df['topic'].unique():
6         topic_df = batch_df[batch_df['topic'] == topic_id].copy()
7         topic_df['words'] = topic_df['words'].apply(lambda x: ' '.join(x))
8         # Concatenate all word strings in topic_df['words']
9         wordcloud_text = ' '.join(topic_df['words'].values)
10        # Filter counts of all words in the wordcloud_text and then select top 100 with count > 1
11        wordcloud_text_counts = wordcloud.process_text(wordcloud_text)
12        wordcloud_text_counts = {k: v for k, v in sorted(wordcloud_text_counts.items(), key=lambda item: item[1], reverse=Tr
13        wordcloud_text_counts = dict(list(wordcloud_text_counts.items())[:100])
14        wordcloud_text_counts = pd.Series(wordcloud_text_counts, index=wordcloud_text_counts.keys())
15        # Filter to words with counts > 1
16        wordcloud_text_counts = wordcloud_text_counts[wordcloud_text_counts > 1]
17        topic = classify_topics(mental_health_topics, creativity_topics, wordcloud_text_counts.to_string())
18        topic_df['topic'] = topic
19        all_classified_topic_batches.append(topic_df)
```

```python
1 # concatenating all tables for summary
2 all_classified_topic_batches_df = pd.concat(all_classified_topic_batches)
```

```python
1 def plot_topic_counts(df, title="Topics", xlabel="Topic", ylabel="Probability"):
2     # Getting topic counts and labels
3     topic_counts = df['topic'].value_counts(normalize=True)
4     labels = topic_counts.index.tolist()
5     counts = topic_counts.values.tolist()
6
7     # Plotting the bars with distinct colors
8     # Generating distinct colors using a colormap
9     num_bars = len(labels)
10    cmap = plt.get_cmap('viridis', num_bars)  # Choose a colormap
11
12    bars = plt.bar(labels, counts, color=[cmap(i) for i in range(num_bars)])
13
14    # Setting chart elements
15    plt.title(title)
16    # Hiding x-ticks and labels
17    plt.xticks([])
18    plt.ylabel(ylabel)
19    #plt.xticks(rotation=45, ha='right')  # Rotate x-axis labels for better visibility
20    plt.tight_layout()  # Prevents labels from overlapping
21
22    # Generating legend handles with x-axis labels as legend entries
23    legend_handles = [mpatches.Patch(color=cmap(i), label=label) for i, label in enumerate(labels)]
24
25    # Smaller legend with adjusted position
26    plt.legend(handles=legend_handles, title=xlabel, loc='upper right', fontsize='small', bbox_to_anchor=(1.02, 1))
27
28    plt.show()
```

```python
1 plot_topic_counts(all_classified_topic_batches_df)
```

## Topics



Mental Health Insight

The classification of the topic id's found via clustering of the key words using a Latent Dirichlet Allocation model for latent text embedding analysis is shown above. It is surprising that a majority of the classified topic labels fall under the 'mental health' category. However, despite the language model not classifying much of the sub-categories within the 'creativity' category, it still categorized much of the clusters as belonging to the general 'creative' label.
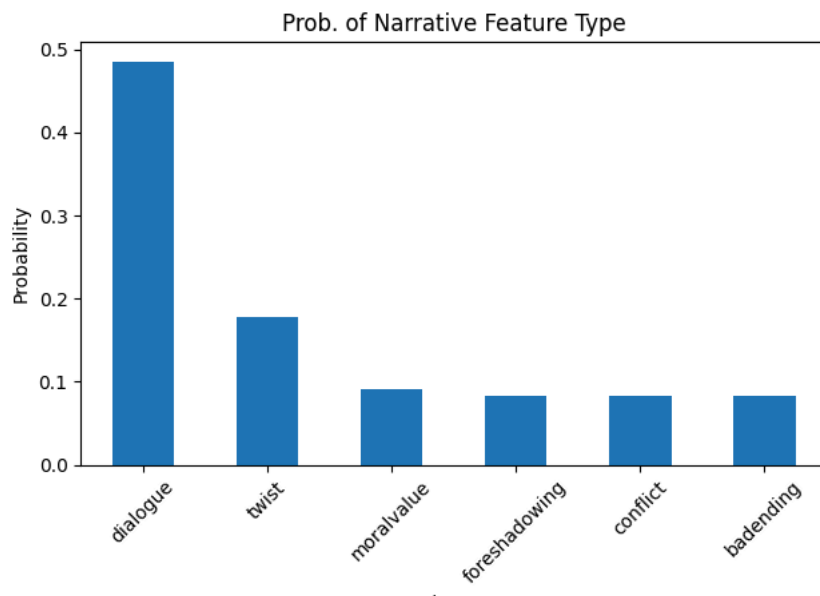
Statistical Insight

## ∨  Creativity

Objectively speaking, creativity is a hard goal to define as it can be objectively defined in a multitude of way, as are the aforementioned topics of literacy and mental-health. However, many people might agree that creativity is somehow unique. Therefore, it may be possible to define the goal of creativity by understanding the level of variance in the models responses to similar prompts.
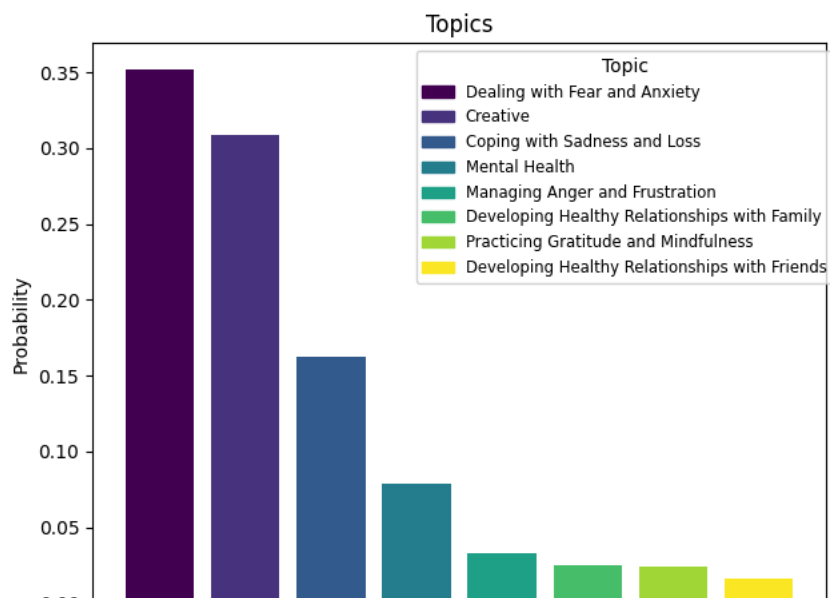
## ∨  Narrative Features

```
1 ax = feature_occurence_probabilities.sort_values(ascending=False).plot(
2     kind='bar',
3     title='Prob. of Narrative Feature Type'
4 )
5
6 # Correctly set labels using set methods
7 plt.ylabel('Probability')
8 plt.xlabel('Narrative Feature')
9 plt.xticks(rotation=45)
10
11 # Optional: Adjust layout to prevent label cutoff
12 plt.tight_layout()
13
14 plt.show()
```

## Prob. of Narrative Feature Type



### Key Words

```
1 plot_topic_counts(all_classified_topic_batches_df)
```

## Topics



### Literacy

WonderWords' literacy goals may be defined in objective terms by classifying whether the responses are at a lower or a higher reading level. As our application is targeting a youth demographic, utilizing the categorization system used in most libraries and school systems will help illustrate whether the training data is biased towards a specific set of reading levels.
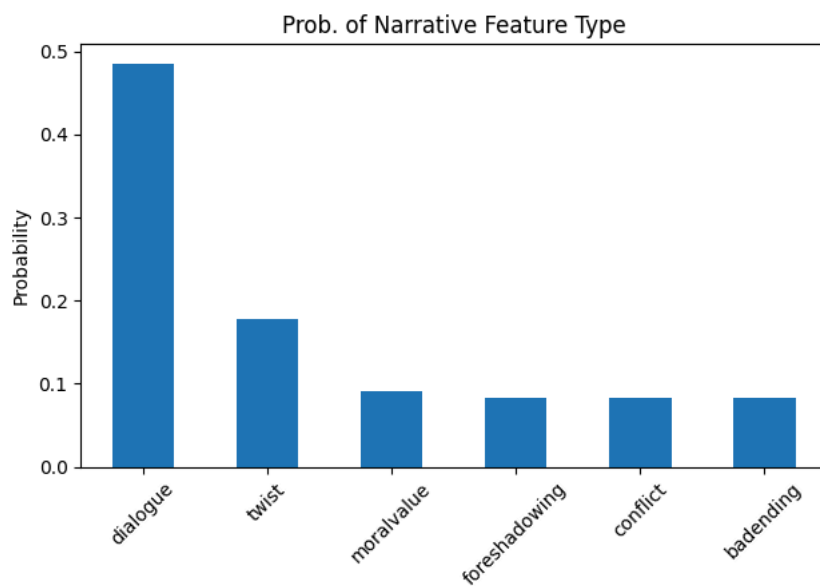
### Narrative Features

```
1 ax = feature_occurence_probabilities.sort_values(ascending=False).plot(
2     kind='bar',
3     title='Prob. of Narrative Feature Type'
4 )
5
6 # Correctly set labels using set methods
```

```
 7 plt.ylabel('Probability')
 8 plt.xlabel('Narrative Feature')
 9 plt.xticks(rotation=45)
10
11 # Optional: Adjust layout to prevent label cutoff
12 plt.tight_layout()
13
14 plt.show()
```



Prob. of Narrative Feature Type

## Key Words

```
 1 literacy_levels = [
 2     "Kindergarten (5-6)",
 3     "1st Grade (6-7)",
 4     "2nd Grade (7-8)",
 5     "3rd Grade (8-9)",
 6     "4th Grade (9-10)",
 7     "5th Grade (10-11)",
 8     "6th Grade (11-12)",
 9     "7th Grade (12-13)",
10     "8th Grade (13-14)",
11     "9th Grade (14-15)",
12     "10th Grade (15-16)",
13     "11th Grade (16-17)",
14     "12th Grade (17-18)"
15 ]
```

```
 1 def classify_literacy(literacy_levels, topic_df_words_counts):
 2     """ this takes a preset list of literacy levels,
 3         and classifies the topic_ids by using the data in the topic_df_words
 4         note: there are 12 topic_ids and 13 corresponding literacy levels.
 5     """
 6     prompt_template = (
 7         f"The following is a list of education levels and corresponding ages: {literacy_levels}",
 8         f"Use this list of words and their probabilities to classify the most likely level for this set of words: {topic_df_
 9         "Note: Each topic in either list can be assigned only once. Only respond with the topic chosen and nothing else.",
10         "Nearest literacy level: "
11     )
12     prompt = "\n".join(prompt_template)
13     response = client.chat.completions.create(
14         messages=[
15             {"role": "system", "content": "You are comfortable classifying the language used in stories by simply inferring
16             {"role": "user", "content": prompt}
17         ],
18         model = model,
19         temperature = 0.0
20     )
21     return response.choices[0].message.content
```
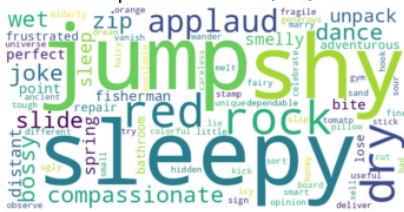
```
 1 # Step 3: creating a word cloud for all topics within a single iteration
 2 # add the cluster labels to the batch
 3 sample_batch = batches[0]
 4 batch_df = pd.DataFrame({'words': sample_batch, 'topic': topic_assignment_batches[0]})
 5 # Determine the number of rows and columns for the grid
 6 num_topics = len(batch_df['topic'].unique())
 7 num_cols = 3  # Adjust as needed
 8 num_rows = (num_topics + num_cols - 1) // num_cols
 9
10 # Create the figure and axes for the grid
11 fig, axes = plt.subplots(num_rows, num_cols, figsize=(15, 5 * num_rows))
12 i = 0
13 for topic_id in batch_df['topic'].unique():
14     topic_df = batch_df[batch_df['topic'] == topic_id].copy()
15
16     # Create a word cloud of the topics
17     wordcloud = WordCloud(width=800, height=400, background_color='white', max_words=100)
18     topic_df['words'] = topic_df['words'].apply(lambda x: ' '.join(x))
19     # Concatenate all word strings in topic_df['words']
20     wordcloud_text = ' '.join(topic_df['words'].values)
21     # Filter counts of all words in the wordcloud_text and then select top 100 with count > 1
22     wordcloud_text_counts = wordcloud.process_text(wordcloud_text)
23     wordcloud_text_counts = {k: v for k, v in sorted(wordcloud_text_counts.items(), key=lambda item: item[1], reverse=True)}
24     wordcloud_text_counts = dict(list(wordcloud_text_counts.items())[:100])
25     wordcloud_text_counts = pd.Series(wordcloud_text_counts, index=wordcloud_text_counts.keys())
26     # Filter to words with counts > 1
27     wordcloud_text_counts = wordcloud_text_counts[wordcloud_text_counts > 1]
28     wordcloud_text_filtered = ' '.join(wordcloud_text_counts.index)
29     # Classifying the topic_df['words']
30     topic = classify_literacy(literacy_levels, wordcloud_text_counts.to_string())
31     topic_df['topic'] = topic
32     wordcloud.generate(wordcloud_text_filtered)
33     # Plot the wordcloud on the corresponding subplot
34     row = i // num_cols
35     col = i % num_cols
36     ax = axes[row, col]  # Get the current subplot
37     ax.imshow(wordcloud, interpolation='bilinear')
38     ax.axis('off')
39     ax.set_title(f'Topic: {topic}')  # Add a title/label
40     i += 1
```

Topic: 2nd Grade (7-8)



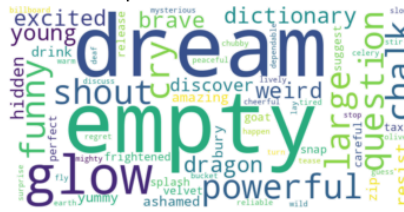Topic: 3rd Grade (8-9)



Topic: 2nd Grade (7-8)



Topic: 4th Grade (9-10)



Topic: 3rd Grade (8-9)



Topic: 4th Grade (9-10)



Topic: 4th Grade (9-10)



Topic: 4th Grade (9-10)



Topic: 2nd Grade (7-8)

```
1 # Step 4: classifying all the topic_ids across all bootstrap cluster iterations
2 all_classified_literacy_batches = []
3 for i in range(len(topic_assignment_batches)):
4     batch_df = pd.DataFrame({'words': batches[i], 'topic': topic_assignment_batches[i]})
5     for topic_id in batch_df['topic'].unique():
6         topic_df = batch_df[batch_df['topic'] == topic_id].copy()
7         topic_df['words'] = topic_df['words'].apply(lambda x: ' '.join(x))
8         # Concatenate all word strings in topic_df['words']
9         wordcloud_text = ' '.join(topic_df['words'].values)
10        # Filter counts of all words in the wordcloud_text and then select top 100 with count > 1
11        wordcloud_text_counts = wordcloud.process_text(wordcloud_text)
12        wordcloud_text_counts = {k: v for k, v in sorted(wordcloud_text_counts.items(), key=lambda item: item[1], reverse=Tr
13        wordcloud_text_counts = dict(list(wordcloud_text_counts.items())[:100])
14        wordcloud_text_counts = pd.Series(wordcloud_text_counts, index=wordcloud_text_counts.keys())
15        # Filter to words with counts > 1
16        wordcloud_text_counts = wordcloud_text_counts[wordcloud_text_counts > 1]
17        topic = classify_literacy(literacy_levels, wordcloud_text_counts.to_string())
18        topic_df['topic'] = topic
```