
HUMAN-INTERPRETABLE CLUSTERING OF SHORT-TEXT USING LARGE LANGUAGE MODELS

Justin K. Miller
School of Physics
University of Sydney
Camperdown, NSW 2006
justin.k.miller@sydney.edu.au

Tristram J. Alexander
School of Physics
University of Sydney
Camperdown, NSW 2006
tristram.alexander@sydney.edu.au

May 14, 2024

ABSTRACT

Large language models have seen extraordinary growth in popularity due to their human-like content generation capabilities. We show that these models can also be used to successfully cluster human-generated content, with success defined through the measures of distinctiveness and interpretability. This success is validated by both human reviewers and ChatGPT, providing an automated means to close the ‘validation gap’ that has challenged short-text clustering. Comparing the machine and human approaches we identify the biases inherent in each, and question the reliance on human-coding as the ‘gold standard’. We apply our methodology to Twitter bios and find characteristic ways humans describe themselves, agreeing well with prior specialist work, but with interesting differences characteristic of the medium used to express identity.

1 Introduction

Short text is playing an increasingly important role in human expression and interaction, due to the widespread use of social media platforms and messaging services such as X (formerly Twitter), Weibo, WhatsApp, Instagram and Reddit. The enormous quantities of data produced by users of these platforms holds the promise of not just real-time identification of events [1] and current opinions [2], but also a deeper understanding of the drivers of information flow between the users [3]. A first step in engaging with the large data sets is to typically reduce the complexity, by clustering the text data into similar groups [4]. However, short text clustering is challenging, due to the limited contextual information available in a single piece of text, and the low incidence of word co-occurrence between pieces of text [5, 6].

The possible applications of success has led to a focus in the machine learning community on clustering, with an increasing number of methods developed to provide a deeper understanding of large collections of short text data [7, 8]. However, a major criticism facing current clustering approaches is that the resulting clusters appear largely uninterpretable by humans [9]. In addition, automated metrics used to quantify clustering success appear to be poorly correlated with interpretability [10]. There is thus a two-fold challenge, the development of an approach which can lead to interpretable clusters of short text data, and an automated approach to stand as a proxy for humans interpreting the resulting clusters.

Approaches to the clustering of short text can be placed under three broad categories: frequency-based approaches, embedding methods, and the use of deep learning techniques [7]. Traditionally, approaches to text involve looking for topics which span the range of documents, in a process known as topic modelling [8]. Typically this involves building a probabilistic model based on observed word frequencies [11], and then for short text, assuming that each document belongs to only one topic [8]. Documents in a given topic are then taken to form a cluster [12]. The most commonly employed topic modelling approach is latent Dirichlet allocation [11], due to its ease of use, however the severe sparsity of short text data is known to lead to poor performance [6].

Embedding approaches instead seek to represent the text data in a vector space [7], where standard clustering techniques such as k-means or Gaussian mixture modelling can be applied [13]. We consider in this class methods such as doc2vec [14], which establish an embedding space by training on the data set of interest. This is in contrast to embeddings created using deep learning methods, such as the Bidirectional Encoder Representations in Transformers (BERT) [15, 16]. In a transformer-based approach, the model is built by training on a large external data set, and then the data set of interest is embedded in the vector space created using the training set.

While transformers have shown early promise for text clustering [16], they have seen spectacular recent success in text generation, thanks to large language models such as ChatGPT [17]. ChatGPT has been rapidly embraced by the public, but large language models more broadly have also seen applications in areas as diverse as finance [18], health [19], law [20] and academia [21]. The transformer-based architectures underpinning these models typically have billions of parameters and are trained on hundreds of GB of text data [22]. While the focus of these models has been on producing human-like content [23], more quantitative applications are starting to emerge, such as sentiment analysis [24] and text annotation [25]. The extraordinary capabilities of these models to produce extended examples of human-like text raises the prospect that they might also be used to group together similar pieces of text, and potentially even be used to interpret the resulting groups.

In this work we use large language models (LLMs) to cluster short text and to interpret the resulting clusters. Out of the many possible clustering techniques available, we compare the LLM with LDA, as a prominent example of topic modelling, and doc2vec as a non-deep-learning approach. Despite the limitations of LDA the method is still widely used [26] and was recently employed in a comparative work on short text clustering methods [12]. Doc2vec was used in the same comparative work, and found to provide the best results when judged using a set of ‘gold standard’ labels of short text documents [12]. In contrast, we consider the case of unlabelled data, and seek to quantify clustering success through human interpretability.

A complication which arises when framing success in cluster creation using human interpretability, is that humans have biases and limitations when faced with the tasks of clustering (also called categorization, the identification of unlabelled groups of similar objects) and classification (labelling of already existing clusters) [27]. Humans appear to cluster preferentially along only one selected dimension in high dimensional data [28], and the dimension chosen depends on the experiences and biases of an individual performing the clustering [29]. Human-labelled data is typically taken as a base truth, or gold standard, however these limitations suggest that, at least for cluster labels, there may be no such thing as a gold standard. One of the goals of this work is to determine a measure of clustering success in the presence of this ambiguity.

We use as a test case for human-interpretable automated short-text clustering a domain familiar from lived experience, that of human identity. The nature of how humans self-identify is complex, however, twelve categories have been found in analysis of identity-related nouns in the English language: occupation, political, kinship, religion, biosocial, ethnicity, corporal, leisure, sexuality, stigma, esteemed, and other [30]. When seeking to classify identity based on short text, we expect there to be substantial ambiguity, making classification challenging for automated approaches and for human reviewers. However, success in this domain promises insights into social media users. Knowledge of identity can be predictive of a person’s behaviour in many areas of life. For instance, political identity can predict a person’s beliefs and actions, such as believing in conspiracy theories or participating in activism for climate change [31, 32]. Humans are interacting to an ever larger degree via social media, and these interactions are typically mediated by short text. With our methodology we show that it is possible to infer information about the people producing the text, when performed at the large data scale.

We focus specifically on a short text data set consisting of ‘Twitter bios’, which contains user responses to the general prompt ‘Describe yourself’. Users have 160 characters to respond to this prompt, resulting in highly variable short text. We begin by using the three methods of LDA, doc2vec and an LLM to find clusters in this data set. We then ask human reviewers to interpret and rate the resulting clusters. We examine what constitutes a ‘good’ cluster in light of the cluster characteristics and the human reviews and look for possible metrics that can provide a quantitative measure of success. We then turn to the possibility of using ChatGPT to act as a proxy for the human reviewers. We examine the results of both the machine and human approaches in light of what is known about human identity and conclude with a proposal of a method for finding and interpreting short text clusters.

2 Results

We use a short-text dataset consisting of 38,639 Twitter user bios from users who used the keywords ‘trump’ or ‘realDonaldTrump’ on September 3rd and 4th 2020 (see Methods for further details). This limits the domain of users to those with some interest in US politics. In Fig. 1 we get an overview of the nature of the bios by examining the top fifty word co-occurrences in the bios. In this ‘top-word co-occurrence graph’ the sizes of the nodes represent the

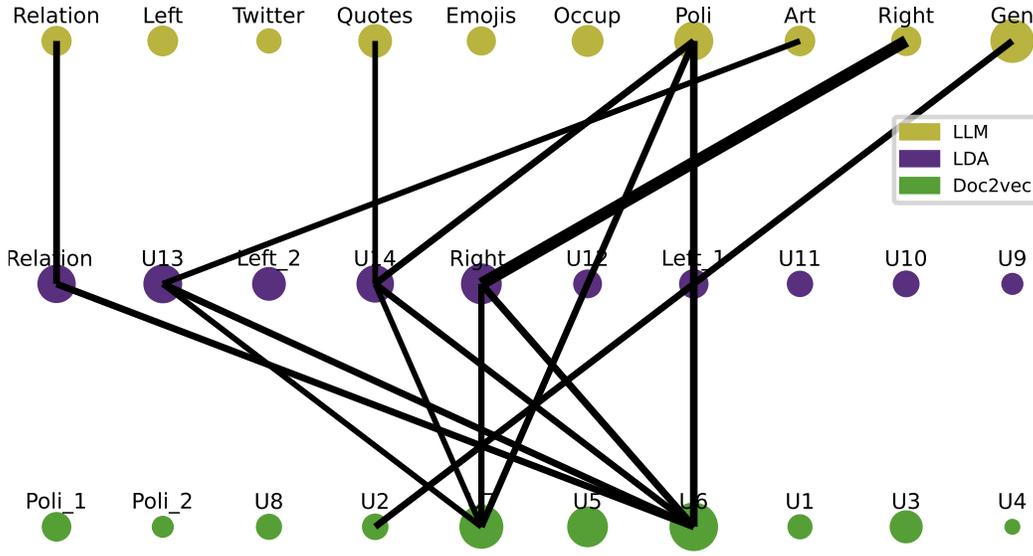


Figure 2: A cluster overlap graph in which nodes represent the clusters by model (Top row: LLM; Middle row: LDA; Bottom row: doc2vec). Lines represent overlap between cluster pairs, as measured by the number of bios in common in a pair. A threshold of 1190 overlapping bios is set for the appearance of a line, which represents the top 5% of cluster overlaps. Thicker lines represent greater numbers of bios in common. Maximum overlap is 2,638 between the ‘Right’ cluster of the LLM and of LDA. Cluster names are given by the authors (see text).

the dataset, only one cluster appears to show any appreciable overlap between the methods, as indicated by the thicker line linking a cluster in the LDA and LLM methods. This appears to be for politically right-leaning users. Overall, the methods appear to have created largely distinct partitions of the short-text data.

A key measure of interpretability is the ability of humans to easily name a cluster [9] and see agreement between top words of a cluster and sample documents. Coherence is a commonly-used automated metric, however its performance as a useful measure has come under recent scrutiny [33]. Mindful of earlier guidelines on human interpretability of clusters [9] we recruited 39 participants to rate: the coherence of a set of 20 sample bios chosen randomly from a cluster, the coherence of the top 10 words from each cluster, and the coherence between the sample bios and the top words (see Methods for details). We also asked the reviewers to name each cluster (if possible) and rate their confidence that the chosen name correctly describes the cluster. As a null model we included also randomly generated clusters, with top words and sample bios from these random clusters provided to the reviewers. The results of the model comparison can be seen in Fig. 3. As the human reviewers have rated their responses on a likert scale, we use ordinal regression [34] with respect to the baseline found for the random set, which is set to 0 (see Methods for further details). We see that the LLM has performed consistently better than the other two methods on all metrics. LDA performed consistently better than random, while doc2vec performed consistently worse. Inspection of the doc2vec clusters reveals that the clustering approach has indeed made use of features that are largely invisible to human reviewers, such as a cluster of bios using low-frequency words. Interestingly doc2vec was found to be successful at clustering when a gold standard is present [12], but in this case there is no gold standard. The results are consistent with the finding that traditional clustering approaches produce clusters that are often uninterpretable by humans [9].

When we look at the performance by cluster in Fig. 4, we find the results are more nuanced. In Fig. 4(a) we see that while the reviewers were usually confident in naming clusters produced by the LLM, some of the clusters were less clear, and one cluster could not be named at all. The cluster names provided on the x -axis have been created by the authors based on observation of the reviewer-provided names and a deeper review of the clusters (consistent with the names used in Fig. 2, see Methods for further details). The ‘General’ cluster could not be named, and the ‘Twitter’, ‘Quotes’ and ‘Emojis’ clusters were also challenging for the reviewers. We examine the names provided by reviewers further below. Interestingly, as seen in Fig. 4(b), coherence of top words, as rated by the reviewers, is not a good predictor of confidence in naming, with little distinction detected between the models, and even the randomly generated clusters rated highly. Looking at Fig. 4(c,d), the LLM appears to have aggregated bios such that they appear more coherent to the reviewers, as well as producing clusters with greater coherence between the bios and

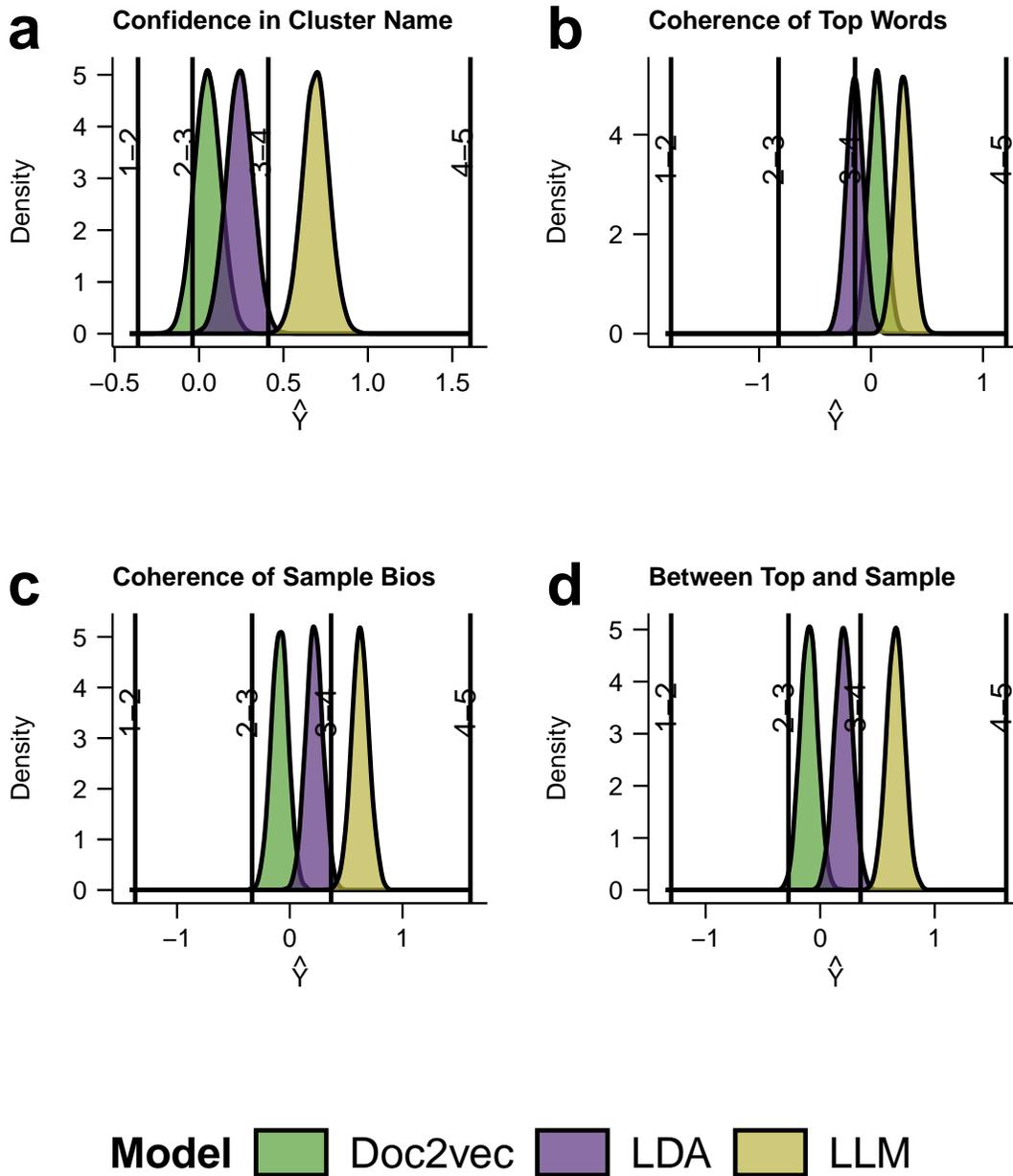


Figure 3: Ordinal Regression Analysis of Clusters created by LDA, Doc2vec, and LLM illustrating the outcomes of an ordinal regression model applied to the ratings of 39 reviewers. The reviewers assessed Twitter bio clusters and were asked to rate the coherence of the clusters across four categories: (a) Confidence in Cluster Name, (b) Coherence of Top Words, (c) Coherence of Sample Bios, and (d) Coherence between the top words and sample bios. Each panel (a-d) represents the density distribution (\hat{Y}) of ratings in each category, showcasing the statistical modeling of ordered categorical data. The reviewer ratings of the random clusters are set to 0 so any value greater than 0 is performing better than random. The vertical lines identify the transitions between values on the likert scales used by the reviewers. We see that the LLM has consistently been scored at 4 on all measures.

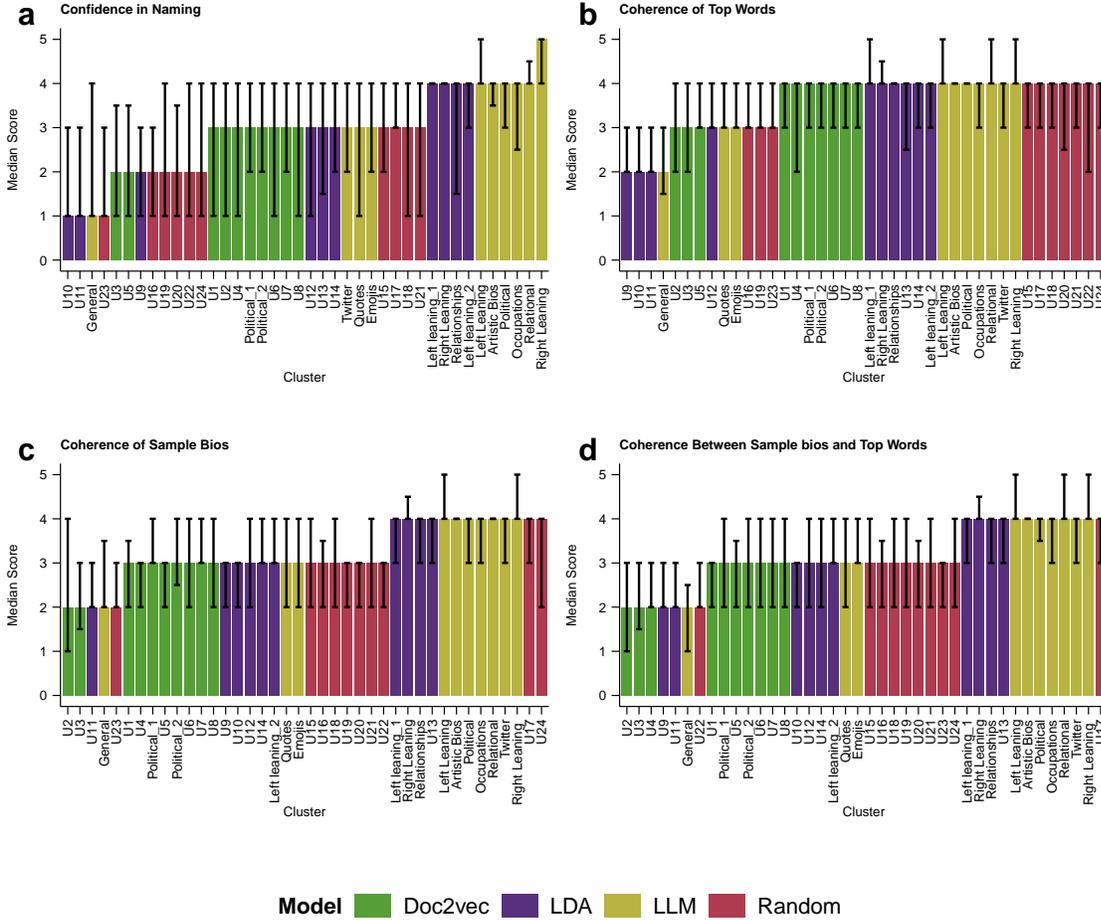


Figure 4: The median reviewer score for each cluster across the four categories: (a) Confidence in Cluster Name, (b) Coherence of Top Words, (c) Coherence of Sample Bios, and (d) Coherence between the top words and sample bios. Error bars represent the first and third quartile scores. All clusters with the U designation were unable to be named by the authors of this paper.

the top words. This suggests that top words can play a role in naming, as is common practice with topic modelling [9], however the high ratings also assigned to random clusters suggests that human estimations of coherence do not provide a measure of clustering success.

While interpretability has been identified as important, we suggest that clusters should also be distinctive. A subset of a randomly generated cluster may appear coherent, and even interpretable, to a human reviewer, due to the well-known human bias of seeing patterns in noise [35]. However, when dealing with large randomly-generated clusters, any fluctuations are expected to average out and the clusters will appear similar to each other. We need a measure to capture this degree of distinctiveness between clusters. To do so, we make use of the method of keyword analysis from corpus linguistics. We identify the number of keywords special to a cluster (see Methods) and plot this versus human confidence in naming, see Fig. 5. We find that all the clusters produced by the LLM have large numbers of distinctive keywords, even the clusters which human reviewers were unable to interpret (such as the ‘General’ cluster). In contrast, all the randomly generated clusters have zero, or close to zero, distinctive keywords, as expected for a complete mixing of the data set. The methods of LDA and doc2vec lead to some clusters with distinctive keywords, but much less distinctiveness than the LLM approach. We find therefore that the large language model has allowed for the creation of not just more human interpretable clusters, but also clusters which are more distinctive relative to each other.

Ideally, we would have an automated approach in which clusters are created and then a metric is used to quantify clustering success. We have identified an automated approach for measuring distinctiveness, we would also like to

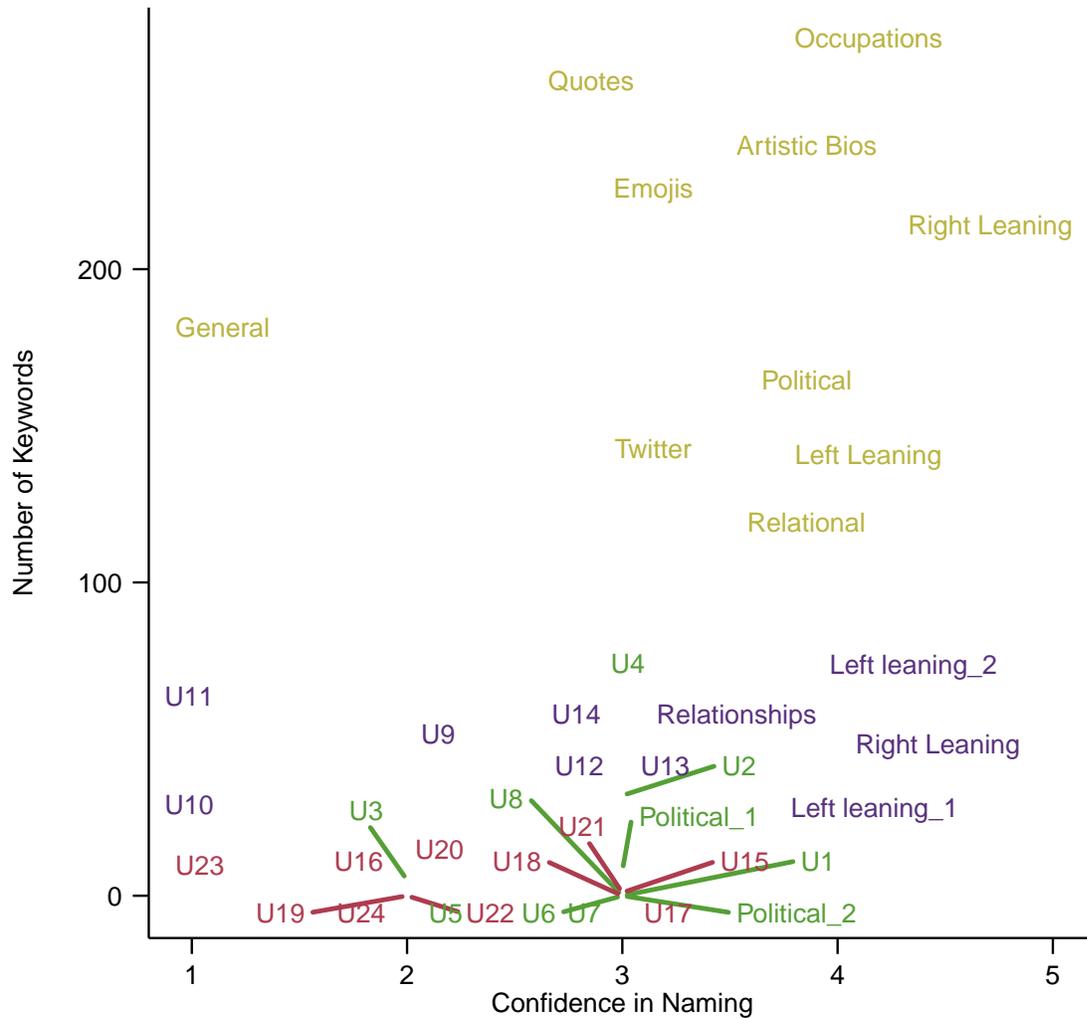


Figure 5: The number of keywords in each cluster, where keywords are determined by comparing the expected frequency of words in a cluster against their actual frequency, using a Bayesian factor greater than 10 to assess if the difference is statistically significant. The names of each cluster are the names given by the authors and the prefix U indicates that a cluster was not able to be named. Colors identify the clustering model used and are consistent with the color code used in Fig. 4. The LLM clusters have significantly more keywords than the clusters found using the other methods.

have a metric which correlates with human interpretability. The problems facing the use of automated metrics have been widely identified [10, 36, 9, 33, 8], with the popular metric of coherence singled out [9, 33]. While there are numerous methods used to calculate coherence [9] we choose two of the most common: CV coherence and UMASS coherence [37, 38, 33]. We also take advantage of the LLM’s ability to convert text into a vector space and use cosine similarity, Euclidean distance and the silhouette score to provide a measure of clustering success. In Fig. 6 we examine the correlation between these metrics and the ratings provided by the reviewers. We see that the correlation is generally poor. The Mean Standard Deviation metric, which serves as a proxy for cluster width in the embedded space, shows the best correlation with the human ratings, followed by the silhouette score. The coherence measures show little correlation with the human measures, including the human ratings of coherence. We find that these coherence measures also show poor correlation with each other (see Fig. S1), so failing a posited criterion of interpretability, that different coherence measures should correlate [9]. Our results are therefore consistent with findings elsewhere that, despite being widely used, the metric of coherence is not useful when measuring human interpretability [33]. We note that the number of keywords also does not correlate with the human measures, however as discussed earlier, we argue this metric provides a measure orthogonal to interpretability, and a lack of correlation is consistent with this.

So far we have focused on the quantification of the clusters provided by the reviewers and the automated metrics, however typically it is nature of the classification which is most interesting when examining clusters. The consistency of names provided by reviewers is also a measure of interpretability [9]. To this end, we turn now to the names of the clusters. The names given to some of the clusters in Fig. 4 are given by the authors based on top words and an extended set of bios in each cluster, as well as the names given by the reviewers. However more interesting is to see the consistency in naming provided by the reviewers. In Fig. 7(a) we show the most popular words appearing by cluster in the reviewer LLM cluster names. The x-axis tick labels are again the researcher-given names for the associated clusters. We see that the Right Leaning cluster has been consistently named by reviewers to be a ‘trump’ related cluster, and the Relational cluster is similarly annotated by the word ‘family’. This is consistent with the coarse clustering approach shown in Fig. 1, where both right-leaning and relational clusters were evident in the top words. The existence of ‘Relational’, ‘Occupation’ and ‘Political’ self-identities are consistent with the broader study of human identity [30]. Other clusters are less obvious and are perhaps specific to the medium of the social network. A sizeable proportion of Twitter account bios reference the use of Twitter as a primary activity, and were picked up and named by the reviewers. Others however appear to have been challenging for the reviewers to name, such as the ‘General’, ‘Quotes’ and ‘Emojis’ clusters. We seek now to identify whether these were difficult to name because they are poorly-defined clusters, or because the classification task was difficult due to the limited information available to the reviewers.

To explore this problem, we contrast the reviewer names, Fig. 7(a), with those produced by ChatGPT-4.0, Fig. 7(b). We find that ChatGPT provides names for all the clusters, unlike many of the human reviewers who selected ‘None’ for some of the clusters. The names ChatGPT has given look consistent with the names provided by the reviewers, and those used by the researchers. ‘General’ and ‘Quotes’ have both been given appropriate names. This suggests that the human reviewers struggled to identify the underlying threads in the clusters of bios. This is consistent with a bias identified in human classification, in which humans tend to seek a low dimensional representation for data when they are asked to perform a classification task (e.g. choosing to focus only on expressions of political ideology) [29]. ‘Quotes’ do not easily fit into this representation, although we note that some reviewers did detect this feature in the set of bios. The ‘General’ class has been picked up by ChatGPT, and such a class is known to exist in expressions of human identity [30], but this class was overwhelmingly left unnamed by the reviewers. Interestingly, the ‘Emojis’ cluster was named by ChatGPT according to its content rather than the medium. This is perhaps unsurprising, given the tool focuses on word usage, and so has provided names based on the content rather than the form of the content. We conclude therefore that clusters that weren’t named by the reviewers, could be named by ChatGPT, suggesting human bias rather than limited information may be involved, however we also found that when faced with finding a ‘meta’ name for a cluster, the humans performed better than ChatGPT.

We are now ready to implement an automated approach for quantifying clustering success. In particular we wish to quantify the performance of ChatGPT with respect to the human reviewers. We simulate this by asking ChatGPT to name clusters multiple times and using the resulting names of both the human reviewers and ChatGPT to calculate consistency (a measure of interpretability) and distinctiveness of the names. These measures are calculated using the two corpora for each cluster: the names used by the reviewers and the names used by ChatGPT. Consistency is measured by counting the number of times the most frequently used word appears in each name, while distinctiveness is calculated using the mean Jensen Shannon divergence (see Methods). We see in Fig. 8 that humans and ChatGPT are broadly consistent with each other. The LLM clusters are given the most distinctive names and are also the most consistently interpreted. We see however some interesting differences between the human and machine approaches. ChatGPT is more consistent in its naming between runs, which is unsurprising given it is a single large language model as compared to a set of different human reviewers. This leads to a set of higher interpretability scores when compared

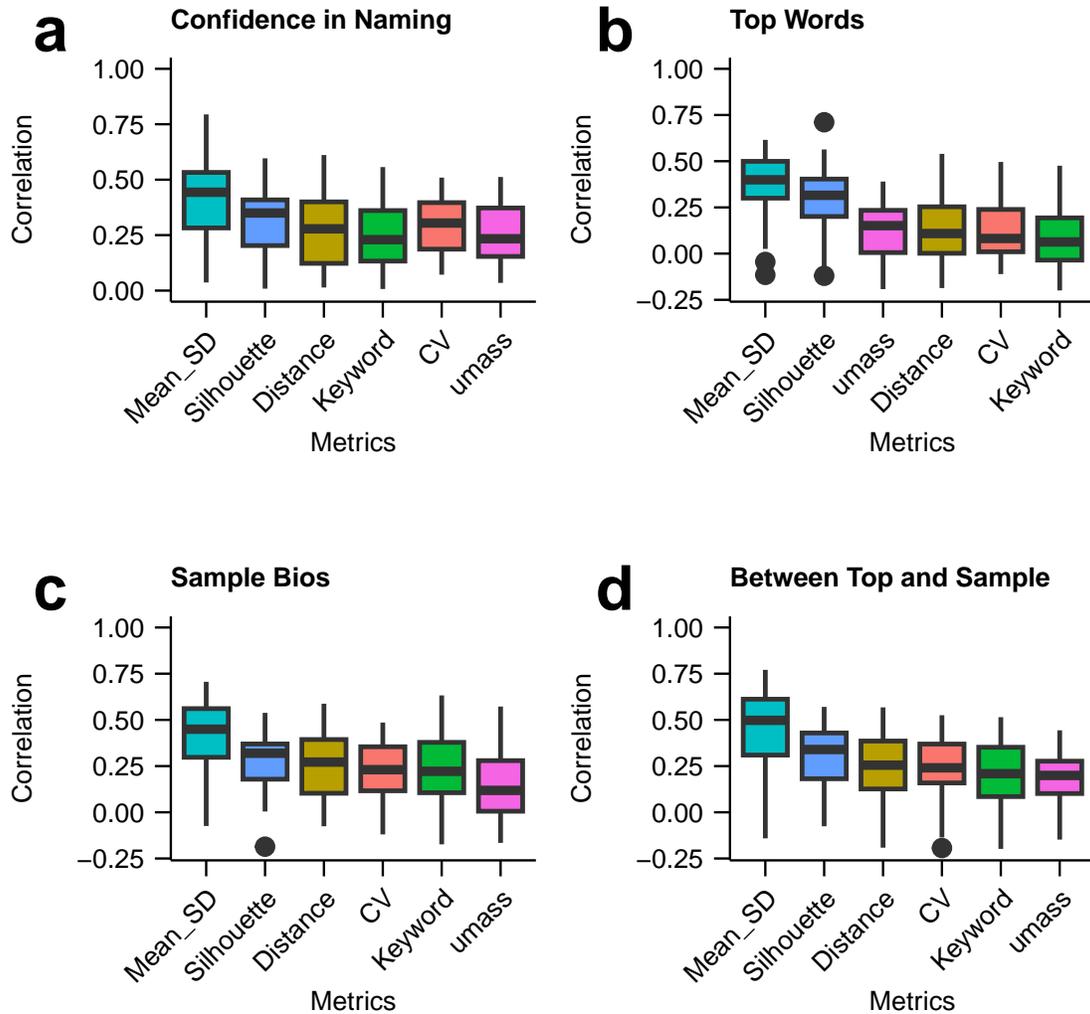


Figure 6: Boxplots showing the Spearman Rank Correlation between the reviewer provided ratings and the six automated metrics of mean cluster standard deviation, silhouette score, Euclidean distance, number of keywords, CV coherence and UMass coherence (x-axis labels), with reviewer ratings on (a) Confidence in naming; (b) Coherence of top words; (c) Coherence of sample bios and (d) Coherence between top words and sample bios. Coherence and keywords correlate poorly with reviewer ratings. Mean standard deviation appears to provide the best correlation with the ratings. The large variability in correlation across reviewers is evident, with outliers identified with solid circles, i.e., some reviewers correlated well with some measures, while others showed no correlation or negative correlation for some measures.

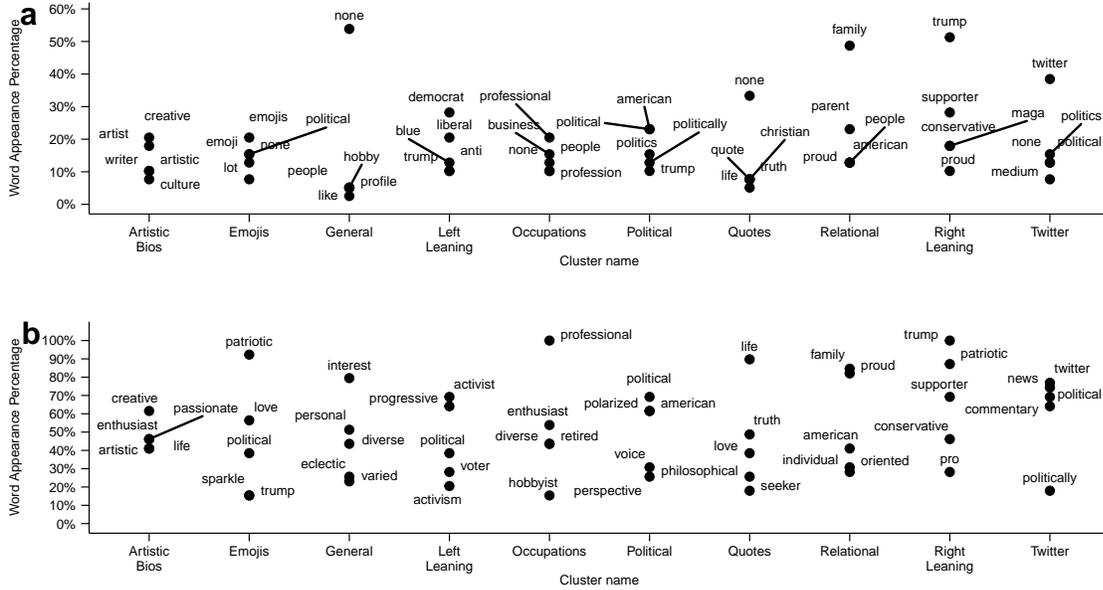


Figure 7: The top five words by fraction of appearance used by (a) reviewers and (b) ChatGPT to name the clusters created by the LLM. Along the x-axis are the names given to each cluster by the authors of this paper. We see that the words used to describe the clusters are largely consistent between ChatGPT and the reviewers, however there are cluster-dependent distinctions revealing human and machine limitations as discussed in the text.

to the human scores. However, while the LLM has scored well in both absolute and relative terms using ChatGPT, the random clusters have also performed well in absolute terms. This is likely a consequence of the underlying political skew to the dataset. ChatGPT appears to have better identified this underlying bias in the dataset, by producing lower scores in the distinctiveness measure for the random clusters when compared with the scores given by the human reviewers. An interesting future direction would be to consider datasets with less of an underlying signal.

Discussion

Our work reveals an interesting contrast between reviewer confidence in naming, and the interpretability measure based on naming consistency. When we ask human reviewers to rate their confidence that the name they have provided for a cluster is representative of the cluster as a whole, we see that the LLM-based clusters performed significantly better than the other clustering approaches. However, when we look at consistency between names as provided by the reviewers, we see that the methods are rated more closely, though the LLM clusters again performed the best. This observation is consistent with the findings by Doogan and Buntine, that looking at the results of reviewers completing a labelling task is more informative than asking reviewers about their ability with the task [9]. Focusing on labels also allows us to develop an automated approach based on ChatGPT. Asking ChatGPT to rate its confidence in naming would seem to be of questionable value, but asking ChatGPT to repeatedly provide a name allows us to obtain a measure of interpretability, using the tool’s strength in text comprehension. The top-word consistency measure we have introduced can be seen as a proxy for the inter-coder reliability measure suggested as a means to quantify interpretability [9].

We have chosen to base our distinctiveness measure on names provided, however an alternative is to use the complete cluster information, as was done in our keyword analysis. Complete information better distinguishes randomly created clusters and the separation of clusters in a semantic space, however a measure based on names better reflects the human perception of the clusters. We have chosen this human-centric approach here, however all the documents in a cluster could be combined and a text dissimilarity measure used to quantify the distinctiveness of the clusters [39].

We see that ChatGPT produces results that are consistent with human reviewers, though with some interesting distinctions. ChatGPT found greater variability between the clustering approaches, and appeared to be significantly better at distinguishing the random clusters than the human reviewers. This suggests that ChatGPT is less prone to the biases

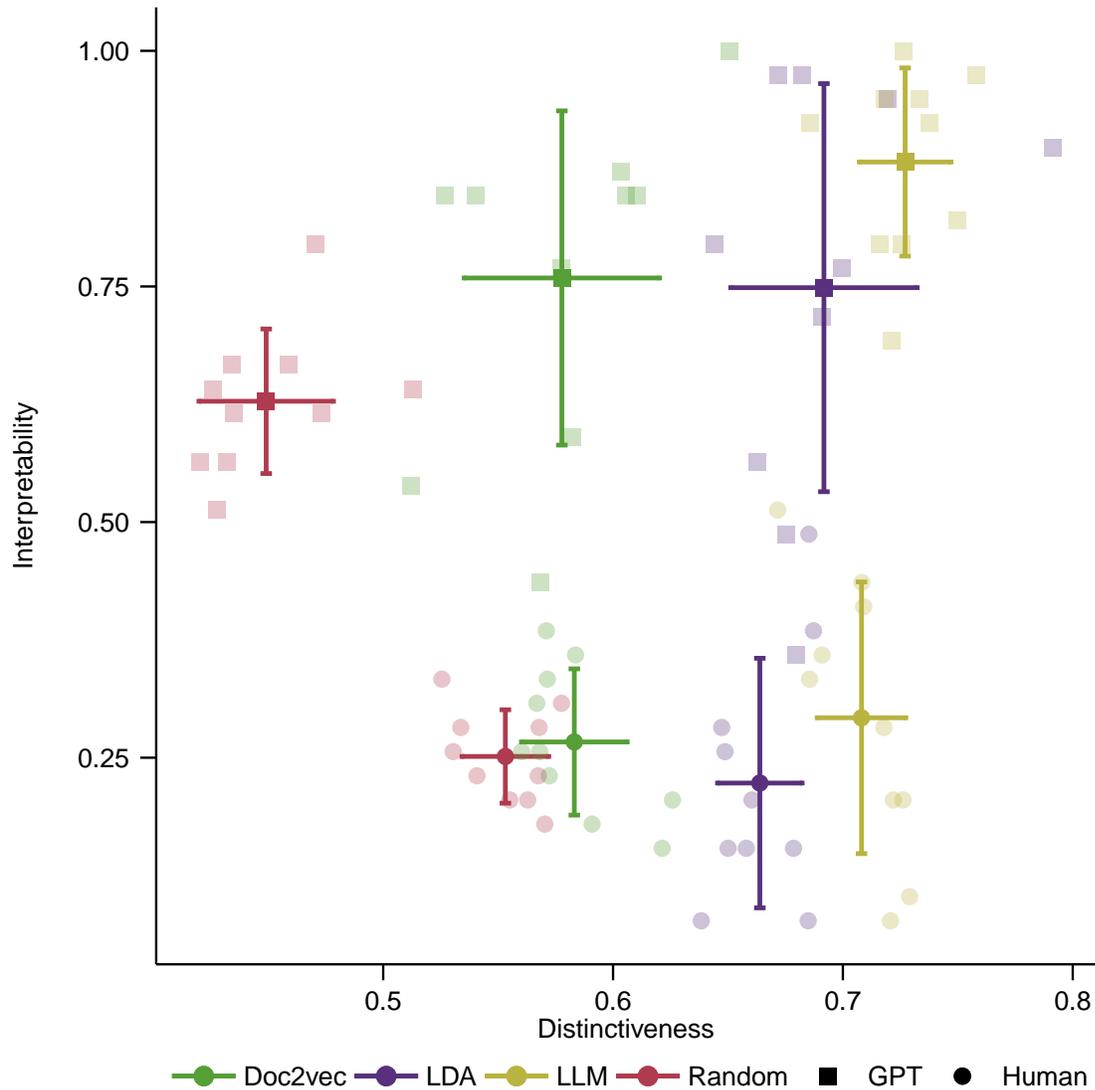


Figure 8: The interpretability and distinctiveness for each cluster as determined by the names given by human reviewers (circles) and ChatGPT (squares) for the four different clusterings (Doc2Vec, LDA, Random and LLM). The fainter symbols are results by cluster, the more solid symbols with error bars indicate the averages. We can see that the LLM produces the most interpretable and distinctive clusters and that ChatGPT is broadly consistent with the human reviewers, though with significant differences, as discussed in the text.

known to affect humans when naming clusters, such as the temptation to see patterns in noise. However, we should note that the data set we used had a strong political orientation, such that reviewers naming random clusters as ‘political’ is reflective of an underlying signal rather than noise. This signal however is present in all clusters, and ChatGPT appears to have been more consistent in detecting this inherent bias and so overall giving a lower distinctiveness measure to the random clusters than that given by the human reviewers.

The human difficulty of distinguishing a signal from noise is the reason we haven’t implemented a commonly suggested approach for measuring interpretability: topic or word intrusion [10]. The combination of short text and the inherent ambiguity in the clusters due to the absence of a gold standard, make it difficult for a human reviewer to identify an ‘intruder’ whether at the topic or document level. This is particularly the case with short text bios. The information provided by a user may conceivably belong to multiple clusters, with the optimal cluster difficult for a human to determine. For instance, a bio might have a clear political orientation and clear references to family. This ambiguity reflects the multifaceted nature of human identity [30]. However to be of value, the clusters that are found should reflect some deeper divide in the data.

We find the following groupings of bios on Twitter: artistic, occupation, relational, left-leaning, right-leaning, political, general, quotes, emojis, Twitter. This is contrasted with the earlier findings on general classes of identity: occupation, political, kinship, religion, biosocial, ethnicity, corporal, leisure, sexuality, stigma, esteemed, and other [30]. We see some overlap with this earlier classification, with users for instance identifying via their family connections or occupation. We find greater sub-division in the ‘political’ identity, which is unsurprising given the political nature of the dataset. The artistic bios would seem to overlap either with ‘leisure’ or ‘occupation’ so also appear consistent. The ‘general’ cluster was not picked up by the human reviewers, but was found by ChatGPT, and the existence of this class is consistent with the ‘other’ class found earlier. Interestingly, we find distinct classes which appear to be particular to the social media platform. Large numbers of users chose to use emojis to express themselves, making use of this additional avenue for expression (the ‘emojis’ cluster). Others self-identified with reference to the platform itself, appearing to clearly separate their online and offline identities (the ‘Twitter’ cluster). The remaining ‘Quotes’ cluster appears to be an amalgam of identity classes. It had many quotes from the bible, which is consistent with the ‘religion’ class of self-identity, but there were also non-religious quotes used. This cluster, like the emojis cluster, therefore appears to be more a grouping of means of expression rather than particular content (so, a ‘meta’ cluster). Whether people who choose to express themselves using emojis or quotes show other characteristics in common is an interesting direction for further research.

We see that the large language model approach to text clustering has revealed interesting, and human-interpretable, clusters. Based on the immense training sets underpinning such models, more sophisticated patterns have been found, such as the use of quotes or emojis, in addition to more clearly identifiable classes of human self-identity. These clusters also appear to show expected relationships with each other when projected into a two-dimensional space (see Fig. S2), with occupations sitting near artistic bios, but far from political bios. The commonly used approach of LDA appears to have also performed well in producing distinctive and interpretable clusters, however it suffers from a repeatability failure. Clusterings using different random seeds produce quite different partitions of the data, limiting the ability to make general comments about a particular observation (see Fig. S4). This appears to be less of a problem when using the doc2vec method (see Fig. S5), however doc2vec produces clusters that are difficult for a human to interpret (such as grouping bios that use rare words together).

Based on our results, the recommended approach in short-text clustering is to use a large language model to obtain an embedding which can be clustered. For automated analysis of the resulting cluster we have found that ChatGPT has performed as well, or better, than human reviewers. Where ChatGPT performed poorly was in the identification of means of expression rather than content (i.e. meta elements such as the use of quotes or emojis), however on balance the model has still performed better than human reviewers when faced with these clusters. We suggest therefore that LLMs may be used in both the creation and validation of clusters, providing a means to close the well-identified ‘validation’ gap that is common in cluster analysis [33].

However, it is important to acknowledge a fundamental issue with using large language models: they are largely black boxes, and in the case of ChatGPT, maintained by a private company. A risk is therefore that implementations and training sets can change over time, affecting the results. We see this when we perform the cluster interpretation using different ChatGPT implementations (see Fig. S7). However, variability is also an issue with human reviewers (see e.g. the outliers in Fig. 6), and there are additional checks that can be used to validate the LLM clustering measures. A quantification of cluster distinctiveness can be performed using all the clustering data, e.g. via keyword analysis, so providing a valuable metric without needing the automated interpretation step of a second large language model. Our analysis also indicates that the top words are a good guide to the nature of the LLM-created cluster. These can be used to perform a human-based interpretation of a cluster as a check against names created by the second LLM. However, to obtain metrics of interpretability and distinctiveness at scale, we have validated that the second LLM

performs well. We hope that this opens up a new approach to analysing short text, allowing new insights to be gained from this increasingly important means of human expression.

3 Methods

3.1 Data

The data used in this paper consisted of 38,639 Twitter user bios from users who used the words ‘trump’ or ‘realDonaldTrump’ between 3rd September and 4th September 2020. The user timelines of these users were collected from Twitter using the Twitter API v2 between September 6th 2020 and September 16th 2020. From these user timelines the username and bio for every user was extracted, however, for this research, only the bio was used.

3.2 Clustering

In this study, clustering models using Python 3 were developed. For data preprocessing, emojis were converted to their CLDR Short Name using the Emoji package [40]. For all models except the Transformer, stopwords were removed and words were lemmatized with NLTK 3.6.5 [41]. LDA and Doc2vec models were built using Gensim 4.1.2 [42]. The Doc2vec method used the parameter values suggested by [12]: 100 dimensions and 75 epochs. The Transformer models employed the all-MiniLM-L6-v2 model [43, 44]. This particular transformer was used as it gives a high performance across different metrics while being fast to run, and has only 384 dimensions when it vectorises text. [44]. A Gaussian mixture model with diagonal covariance was applied to both Doc2vec and Transformer embeddings. Diagonal covariance was chosen as each dimension was orthogonal to each other. The chosen number of clusters and topics (K) was set to 10, a balance between complexity and interpretability.

3.3 Evaluation of model stability

To evaluate the stability of each clustering method, we initially executed each model using 50 distinct random seeds. For every pairwise comparison, we calculated the Adjusted Mutual Information (AMI) between each pair of seeds within the same model. Following this, we computed the overall average AMI and the standard deviation for each model.

3.4 Human Evaluation of Models

Human reviewers (N=39) evaluated clusters/topics generated by LDA, Doc2vec, Transformer, and a randomly created model. Reviewers were recruited from Prolific and were screened to be American and speak English as their first language [45]. From Prolific they were directed to Redcap [46, 47] hosted at the University of Sydney where they saw the data and filled in their responses. For the LDA model, samples included the top 10 high-probability words and 20 bios assigned to each topic. The Doc2vec, Transformer, and random models provided the top 10 frequent words and 20 bios per cluster. Reviewers were shown each cluster from one of the 4 models in a random order, and were asked the following questions:

- “Create a name using less than 10 words to summarize the top 10 words/emojis and the sample bios of the Sample Cluster. If you believe it is not possible to do this with a cluster, write ‘None’.”
- “When you named a Sample Cluster, were you confident that the name summarized the whole cluster?” Reviewers can answer this question by selecting one of the following: Not at all Confident/Not Confident/Neutral/Confident/Very Confident.

For the following questions, reviewers can respond with: Strongly Disagree / Disagree / Neutral / Agree / Strongly Agree

- “‘The situation when the top words/emojis of a cluster fit together in a natural or reasonable way’ Do you agree that the above statement describes the top words/emojis of this Sample Cluster?”
- “‘The situation when the Sample bios of a cluster fit together in a natural or reasonable way’ Do you agree that the above statement describes the sample bios of this Sample Cluster?”
- “‘The Top Words/Emojis provide an accurate summary of the sample bios’ Do you agree that the above statement describes this Sample Cluster?”

See figure S10 and S11 to see how this was presented to the reviewers.

3.5 Ordinal Regression

Survey responses were scaled from 1 (“Strongly Disagree” or “Not at All Confident”) to 5 (“Strongly Agree” or “Very Confident”). The study employed four Bayesian cumulative ordinal regression models, with the dependent variable being the questions the reviewers were asked and the independent variable being the cluster-creating model. All ordinal regression models were created using the BRMS package in R [48]. We used four chains with 5,000 iterations, 2,500 of which were warm-up. This resulted in a total of 10,000 draws from the posterior distribution for each of the four models. The priors for each model were the default from the BRMS package: student’s t-distribution with 3 degrees of freedom, location parameter 0, and scale parameter 2.5. These models generate a latent variable \hat{Y} [34], with thresholds for each response level. A higher \hat{Y} indicates a greater likelihood of “Agree” or “Strongly Agree”. Model performance is evaluated by comparing their coefficient distributions against the intercept (random model), set at 0, to assess how much better they perform.

3.6 Keyword Analysis

In keyword analysis, we identify words that occur more frequently in a specific corpus compared to a standard ‘reference corpus,’ as outlined by Baker (2006) [49]. Initially, we count how often each word appears within individual clusters of the corpus. Next, we determine the frequency of these words across the entire corpus. Using these frequency counts, we calculate the Bayes factor for each word in each cluster. Words with a Bayes factor exceeding 10 are classified as keywords. This threshold indicates that a word is significantly more prevalent in a particular cluster compared to its general frequency in the entire corpus [50].

3.7 Metric Calculations

Coherence

Before calculating coherence, all stopwords and punctuation were removed, and every word was lemmatized using NLTK 3.6.5 [41]. Coherence for every cluster was calculated using both the C_{UMASS} and C_V methods using the Gensim 4.1.2 package in Python [42]. The top 10 words for each cluster were used in the calculation of coherence.

3.7.1 Distance metrics

For all four models, we utilized the embeddings generated from all-mini-lml, calculating the silhouette score using cosine distance for each individual bio within these models (sklearn [51]). Subsequently, we computed the mean silhouette score for the bios in each respective cluster, providing an aggregated measure of their cohesion. Cluster distance was calculated by finding the mean position of each cluster, then getting the pairwise distance between each cluster within a model. The shortest distance for each cluster was then taken as a measure of how well the bios in that group were clustered. Mean standard deviation for each cluster was calculated similarly to cluster distance, but instead of taking the mean of the cluster, it takes the mean standard deviation across all 384 dimensions.

3.7.2 Correlation between reviewer answers and automated metrics

To test how correlated each automated method was with reviewer responses, for each of the four questions we computed the Spearman rank correlation between all 39 reviewers and all 40 clusters. The Spearman rank correlation, a non-parametric measure, assesses the strength and direction of association between two ranked variables. This gives us a distribution of 39 reviewers and their correlation with the automated metrics. Spearman rank correlation was used since the data is ordinal and a rank-based approach gives the most accurate representation of the data [52].

3.8 Cluster Names

By looking at the names our reviewers gave each cluster, as well as looking at the top 10 words, and reviewing a large number of bios, the authors of this paper created names for each cluster for clarity’s sake. If no meaningful name could be created then the cluster would be designated “U” for “Undecided” and a number, the number serves only as a unique identifier. Clusters within the same model that seemed to warrant the same name (e.g. ‘Left leaning’) were given a 1 or a 2 as a unique identifier.

In our study, a cluster’s names, as provided by our reviewers, were compiled into a corpus. We quantified the frequency of each word within these corpora, subsequently identifying the top 5 most frequently occurring words for every cluster. Additionally, to assess the relative prominence of these words, we converted their counts into percentages based on the number of reviewers (i.e. 39). These percentages were then graphically represented.

We replicated this methodology for the ChatGPT version, with a key modification in data collection. Instead of human reviewers, we employed the ChatGPT API [53], inputting a script that instructed ChatGPT to create a name using up to 5 words for a given cluster. This script included the cluster’s top 10 words and a random sample of Twitter bios for context. In cases where ChatGPT deemed a cluster unnameable or just a random amalgamation, it was instructed to return ‘None’. This procedure was repeated 39 times to match the number of human reviewers and plotted in the same way.

3.9 Distinctiveness and naming measure

To understand the distinctiveness of each cluster, we compile all the names given to a particular cluster into a single corpus, by both ChatGPT and our human reviewers. We apply stemming to all words in these corpora to reduce them to their base forms. Next, we use the Jensen-Shannon divergence to measure the linguistic diversity within each cluster [54]. This is given by $JSD(P \parallel Q) = \frac{1}{2}D_{KL}(P \parallel M) + \frac{1}{2}D_{KL}(Q \parallel M)$. Where D_{KL} is the Kullback-Leibler Divergence. For this, we calculate the JSD between each cluster within a model and identify the smallest distance for each one. We then compute the average and standard deviation of these minimum distances. A higher average Jensen-Shannon divergence suggests that the clusters in our model encapsulate a broader range of topics and exhibit more linguistic variety. Additionally, we created an interpretability metric based on consistency of naming, defined as

$$\text{Interpretability} = \frac{\max(\text{Word in corpus})}{\text{Number of reviewers}}. \quad (1)$$

References

- [1] F. Atefeh and W. Khreich, “A Survey of Techniques for Event Detection in Twitter,” *Computational Intelligence*, vol. 31, no. 1, pp. 132–164, 2015.
- [2] K. Ravi and V. Ravi, “A survey on opinion mining and sentiment analysis: Tasks, approaches and applications,” *Knowledge-Based Systems*, vol. 89, pp. 14–46, Nov. 2015.
- [3] S. Pei, L. Muchnik, J. S. Andrade Jr., Z. Zheng, and H. A. Makse, “Searching for superspreaders of information in real-world social media,” *Scientific Reports*, vol. 4, no. 1, p. 5547, Jul. 2014.
- [4] V. N. Gudivada, “Chapter 12 - Natural Language Core Tasks and Applications,” in *Handbook of Statistics*, ser. Computational Analysis and Understanding of Natural Languages: Principles, Methods and Applications, V. N. Gudivada and C. R. Rao, Eds. Elsevier, Jan. 2018, vol. 38, pp. 403–428.
- [5] L. Hong and B. D. Davison, “Empirical study of topic modeling in Twitter,” in *Proceedings of the First Workshop on Social Media Analytics*, ser. SOMA ’10. New York, NY, USA: Association for Computing Machinery, Jul. 2010, pp. 80–88.
- [6] X. Cheng, X. Yan, Y. Lan, and J. Guo, “BTM: Topic Modeling over Short Texts,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 12, pp. 2928–2941, Dec. 2014.
- [7] M. H. Ahmed, S. Tiun, N. Omar, and N. S. Sani, “Short text clustering algorithms, application and challenges: A survey,” *Applied Sciences*, vol. 13, no. 1, p. 342, 2022.
- [8] J. Qiang, Z. Qian, Y. Li, Y. Yuan, and X. Wu, “Short Text Topic Modeling Techniques, Applications, and Performance: A Survey,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 03, pp. 1427–1445, Mar. 2022.
- [9] C. Doogan and W. Buntine, “Topic Model or Topic Twaddle? Re-evaluating Semantic Interpretability Measures,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 3824–3848.
- [10] J. Chang, S. Gerrish, C. Wang, J. Boyd-graber, and D. Blei, “Reading Tea Leaves: How Humans Interpret Topic Models,” in *Advances in Neural Information Processing Systems*, vol. 22. Curran Associates, Inc., 2009.
- [11] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet Allocation,” *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Mar. 2003.
- [12] S. A. Curiskis, B. Drake, T. R. Osborn, and P. J. Kennedy, “An evaluation of document clustering and topic modelling in two online social networks: Twitter and Reddit,” *Information Processing & Management*, vol. 57, no. 2, p. 102034, Mar. 2020.
- [13] M. Z. Rodriguez, C. H. Comin, D. Casanova, O. M. Bruno, D. R. Amancio, L. d. F. Costa, and F. A. Rodrigues, “Clustering algorithms: A comparative approach,” *PLOS ONE*, vol. 14, no. 1, p. e0210236, Jan. 2019.

- [14] Q. V. Le and T. Mikolov, “Distributed Representations of Sentences and Documents,” *arXiv.org*, May 2014.
- [15] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186.
- [16] A. Subakti, H. Murfi, and N. Hariadi, “The performance of BERT as data representation of text clustering,” *Journal of Big Data*, vol. 9, no. 1, p. 15, Feb. 2022.
- [17] OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Baltescu, H. Bao, M. Bavarian, J. Belgum, I. Bello, J. Berdine, G. Bernadett-Shapiro, C. Berner, L. Bogdonoff, O. Boiko, M. Boyd, A.-L. Brakman, G. Brockman, T. Brooks, M. Brundage, K. Button, T. Cai, R. Campbell, A. Cann, B. Carey, C. Carlson, R. Carmichael, B. Chan, C. Chang, F. Chantzis, D. Chen, S. Chen, R. Chen, J. Chen, M. Chen, B. Chess, C. Cho, C. Chu, H. W. Chung, D. Cummings, J. Currier, Y. Dai, C. Decareaux, T. Degry, N. Deutsch, D. Deville, A. Dhar, D. Dohan, S. Dowling, S. Dunning, A. Ecoffet, A. Eleti, T. Eloundou, D. Farhi, L. Fedus, N. Felix, S. P. Fishman, J. Forte, I. Fulford, L. Gao, E. Georges, C. Gibson, V. Goel, T. Gogineni, G. Goh, R. Gontijo-Lopes, J. Gordon, M. Grafstein, S. Gray, R. Greene, J. Gross, S. S. Gu, Y. Guo, C. Hallacy, J. Han, J. Harris, Y. He, M. Heaton, J. Heidecke, C. Hesse, A. Hickey, W. Hickey, P. Hoeschele, B. Houghton, K. Hsu, S. Hu, X. Hu, J. Huizinga, S. Jain, S. Jain, J. Jang, A. Jiang, R. Jiang, H. Jin, D. Jin, S. Jomoto, B. Jonn, H. Jun, T. Kaftan, Ł. Kaiser, A. Kamali, I. Kanitscheider, N. S. Keskar, T. Khan, L. Kilpatrick, J. W. Kim, C. Kim, Y. Kim, H. Kirchner, J. Kiros, M. Knight, D. Kokotajlo, Ł. Kondraciuk, A. Kondrich, A. Konstantinidis, K. Kosic, G. Krueger, V. Kuo, M. Lampe, I. Lan, T. Lee, J. Leike, J. Leung, D. Levy, C. M. Li, R. Lim, M. Lin, S. Lin, M. Litwin, T. Lopez, R. Lowe, P. Lue, A. Makanju, K. Malfacini, S. Manning, T. Markov, Y. Markovski, B. Martin, K. Mayer, A. Mayne, B. McGrew, S. M. McKinney, C. McLeavey, P. McMillan, J. McNeil, D. Medina, A. Mehta, J. Menick, L. Metz, A. Mishchenko, P. Mishkin, V. Monaco, E. Morikawa, D. Mossing, T. Mu, M. Murati, O. Murk, D. Mély, A. Nair, R. Nakano, R. Nayak, A. Neelakantan, R. Ngo, H. Noh, L. Ouyang, C. O’Keefe, J. Pachocki, A. Paino, J. Palermo, A. Pantuliano, G. Parascandolo, J. Parish, E. Parparita, A. Passos, M. Pavlov, A. Peng, A. Perelman, F. d. A. B. Peres, M. Petrov, H. P. d. O. Pinto, Michael, Pokornyy, M. Pokrass, V. Pong, T. Powell, A. Power, B. Power, E. Proehl, R. Puri, A. Radford, J. Rae, A. Ramesh, C. Raymond, F. Real, K. Rimbach, C. Ross, B. Rotsted, H. Roussez, N. Ryder, M. Saltarelli, T. Sanders, S. Santurkar, G. Sastry, H. Schmidt, D. Schnurr, J. Schulman, D. Selsam, K. Sheppard, T. Sherbakov, J. Shieh, S. Shoker, P. Shyam, S. Sidor, E. Sigler, M. Simens, J. Sitkin, K. Slama, I. Sohl, B. Sokolowsky, Y. Song, N. Staudacher, F. P. Such, N. Summers, I. Sutskever, J. Tang, N. Tezak, M. Thompson, P. Tillet, A. Tootoonchian, E. Tseng, P. Tuggle, N. Turley, J. Tworek, J. F. C. Uribe, A. Vallone, A. Vijayvergiya, C. Voss, C. Wainwright, J. J. Wang, A. Wang, B. Wang, J. Ward, J. Wei, C. J. Weinmann, A. Welihinda, P. Welinder, J. Weng, L. Weng, M. Wiethoff, D. Willner, C. Winter, S. Wolrich, H. Wong, L. Workman, S. Wu, J. Wu, M. Wu, K. Xiao, T. Xu, S. Yoo, K. Yu, Q. Yuan, W. Zaremba, R. Zellers, C. Zhang, M. Zhang, S. Zhao, T. Zheng, J. Zhuang, W. Zhuk, and B. Zoph, “GPT-4 Technical Report,” <https://arxiv.org/abs/2303.08774v4>, Mar. 2023.
- [18] A. H. Huang, H. Wang, and Y. Yang, “FinBERT: A Large Language Model for Extracting Information from Financial Text*,” *Contemporary Accounting Research*, vol. 40, no. 2, pp. 806–841, 2023.
- [19] Y. Si, J. Wang, H. Xu, and K. Roberts, “Enhancing clinical concept extraction with contextual embeddings,” *Journal of the American Medical Informatics Association*, vol. 26, no. 11, pp. 1297–1304, Nov. 2019.
- [20] I. Chalkidis and D. Kampas, “Deep learning in law: Early adaptation and legal word embeddings trained on large corpora,” *Artificial Intelligence and Law*, vol. 27, no. 2, pp. 171–198, Jun. 2019.
- [21] S. Altmäe, A. Sola-Leyva, and A. Salumets, “Artificial intelligence in scientific writing: A friend or a foe?” *Reproductive Biomedicine Online*, vol. 47, no. 1, pp. 3–9, Jul. 2023.
- [22] L. Floridi and M. Chiriatti, “GPT-3: Its Nature, Scope, Limits, and Consequences,” *Minds and Machines*, vol. 30, no. 4, pp. 681–694, Dec. 2020.
- [23] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu, Z. Wu, L. Zhao, D. Zhu, X. Li, N. Qiang, D. Shen, T. Liu, and B. Ge, “Summary of ChatGPT-Related research and perspective towards the future of large language models,” *Meta-Radiology*, vol. 1, no. 2, p. 100017, Sep. 2023.
- [24] T. Susnjak, “Applying BERT and ChatGPT for Sentiment Analysis of Lyme Disease in Scientific Literature,” Feb. 2023.
- [25] F. Gilardi, M. Alizadeh, and M. Kubli, “ChatGPT outperforms crowd workers for text-annotation tasks,” *Proceedings of the National Academy of Sciences*, vol. 120, no. 30, p. e2305016120, Jul. 2023.

- [26] H. Jelodar, Y. Wang, C. Yuan, X. Feng, X. Jiang, Y. Li, and L. Zhao, “Latent Dirichlet allocation (LDA) and topic modeling: Models, applications, a survey,” *Multimedia Tools and Applications*, vol. 78, no. 11, pp. 15 169–15 211, Jun. 2019.
- [27] J. Anderson, “The Adaptive Nature of Human Categorization,” *Psychological Review*, vol. 98, no. 3, pp. 409–429, 1991.
- [28] C. A. Doan and R. Vigo, “A comparative investigation of integral- and separable-dimension stimulus-sorting behavior,” *Psychological Research*, vol. 87, no. 6, pp. 1917–1943, Sep. 2023.
- [29] F. Bröker, B. C. Love, and P. Dayan, “When unsupervised training benefits category learning,” *Cognition*, vol. 221, p. 104984, Apr. 2022.
- [30] N. J. MacKinnon and D. R. Heise, *Self, Identity, and Social Institutions*. New York: Palgrave Macmillan US, 2010.
- [31] A. Sternisko, A. Cichocka, and J. J. Van Bavel, “The dark side of social movements: Social identity, non-conformity, and the lure of conspiracy theories,” *Current Opinion in Psychology*, vol. 35, pp. 1–6, 2020.
- [32] C. M. L. Mackay, F. Cristoffanini, J. D. Wright, S. D. Neufeld, H. F. Ogawa, and M. T. Schmitt, “Connection to nature and environmental activism: Politicized environmental identity mediates a relationship between identification with nature and observed environmental activist behaviour,” *Current Research in Ecological and Social Psychology*, vol. 2, p. 100009, Jan. 2021.
- [33] A. Hoyle, P. Goel, D. Peskov, A. Hian-Cheong, J. Boyd-Graber, and P. Resnik, “Is Automated Topic Model Evaluation Broken?: The Incoherence of Coherence,” Oct. 2021.
- [34] P.-C. Bürkner and M. Vuorre, “Ordinal Regression Models in Psychology: A Tutorial,” *Advances in Methods and Practices in Psychological Science*, vol. 2, no. 1, pp. 77–101, Mar. 2019.
- [35] K. R. Clarke, P. J. Somerfield, and R. N. Gorley, “Testing of null hypotheses in exploratory community analyses: Similarity profiles and biota-environment linkage,” *Journal of Experimental Marine Biology and Ecology*, vol. 366, no. 1, pp. 56–69, Nov. 2008.
- [36] Z. C. Lipton, “The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery,” *Queue*, vol. 16, no. 3, pp. 31–57, Jun. 2018.
- [37] R. F. Sear, N. Velásquez, R. Leahy, N. J. Restrepo, S. E. Oud, N. Gabriel, Y. Lupu, and N. F. Johnson, “Quantifying COVID-19 Content in the Online Health Opinion War Using Machine Learning,” *IEEE Access*, vol. 8, pp. 91 886–91 893, 2020.
- [38] S. Syed and M. Spruit, “Full-Text or Abstract? Examining Topic Coherence Scores Using Latent Dirichlet Allocation,” in *2017 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, Oct. 2017, pp. 165–174.
- [39] B. Shade and E. G. Altmann, “Quantifying the Dissimilarity of Texts,” *Information*, vol. 14, no. 5, p. 271, May 2023.
- [40] T. Kim, K. Wurster, and T. Jalilov, “Emoji 2.0.0,” 2022.
- [41] E. Loper and S. Bird, “Nltk: The natural language toolkit,” *arXiv preprint cs/0205028*, 2002.
- [42] R. Rehurek and P. Sojka, “Gensim–python framework for vector space modelling,” *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic*, vol. 3, no. 2, 2011.
- [43] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, “Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 5776–5788, 2020.
- [44] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. [Online]. Available: <https://arxiv.org/abs/1908.10084>
- [45] Prolific, “Prolific,” <https://www.prolific.com>, London, UK, 2023, version: [Current month(s) and year(s) of use].
- [46] P. Harris, R. Taylor, R. Thielke, J. Payne, N. Gonzalez, and J. Conde, “Research electronic data capture (redcap) – a metadata-driven methodology and workflow process for providing translational research informatics support,” *Journal of Biomedical Informatics*, vol. 42, no. 2, pp. 377–381, 2009.
- [47] P. Harris, R. Taylor, B. Minor, V. Elliott, M. Fernandez, L. O’Neal, L. McLeod, G. Delacqua, F. Delacqua, J. Kirby, S. Duda, and R. Consortium, “The redcap consortium: Building an international community of software partners,” *Journal of Biomedical Informatics*, 2019, doi: 10.1016/j.jbi.2019.103208.

- [48] P.-C. Bürkner, “Bayesian item response modeling in R with brms and Stan,” *Journal of Statistical Software*, vol. 100, no. 5, pp. 1–54, 2021.
- [49] P. Baker, *Using corpora in discourse analysis*. A&C Black, 2006.
- [50] A. Wilson, *Embracing Bayes factors for key item analysis in corpus linguistics*, ser. Language Competence and Language Awareness in Europe. Peter Lang, 2013, pp. 3–11.
- [51] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [52] P. Ornstein and J. Lyhagen, “Asymptotic properties of spearman’s rank correlation for variables with finite support,” *PLoS One*, vol. 11, no. 1, p. e0145595, 2016.
- [53] OpenAI, “Gpt-4 and gpt-4 turbo,” <https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>, 2023, accessed: 2023-12-10.
- [54] J. Lu, M. Henchion, and B. Mac Namee, “Diverging divergences: Examining variants of Jensen Shannon divergence for corpus comparison tasks,” in *Proceedings of the Twelfth Language Resources and Evaluation Conference*, N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis, Eds. Marseille, France: European Language Resources Association, May 2020, pp. 6740–6744. [Online]. Available: <https://aclanthology.org/2020.lrec-1.832>

Supplementary Materials

Justin K. Miller, Tristram J. Alexander

May 14, 2024

1 Correlation between metrics

To better understand how each automated metric compares, we calculated the Pearson rank correlation for each pair of metrics based on their scores within each cluster, as presented in the heatmap (see Fig. S1). Notably, silhouette scores showed a correlation exclusively with the mean standard deviation. This correlation is logical given that both metrics aim to measure similar aspects. However, the silhouette scores' low correlation with other metrics is unexpected. In contrast, Keywords and Distance metrics exhibited a significant correlation, aligning with their common goal of assessing cluster distinctiveness. Despite their different approaches, this correlation validates their effectiveness as tools for measuring the uniqueness of clusters.

2 UMAP Visualisation

To effectively visualize the relationships between different clusters in our 384-dimensional large language model dataset, we employed Uniform Manifold Approximation and Projection (UMAP) for dimensionality reduction [1], see Fig.S 3. The choice of UMAP is particularly pertinent due to its exceptional ability to retain the topology of the data, a feature crucial for our objective of illustrating the spatial relationships between various bios. Unlike Principal Component Analysis (PCA), which primarily captures linear relationships and variance, UMAP excels in preserving both the local and global structures within the data. This is especially important given that our dataset exhibits orthogonal dimensions, as evidenced by the correlation plot in Fig. S2. Although PCA might typically be effective in scenarios with orthogonal dimensions, its linear approach could potentially overlook the more complex, non-linear relationships present in our data. Therefore, UMAP's capacity to handle non-linear structures makes it a more suitable choice for our analysis, ensuring a more accurate and insightful two-dimensional representation of the clusters. The results of this can be seen in Fig. S3.

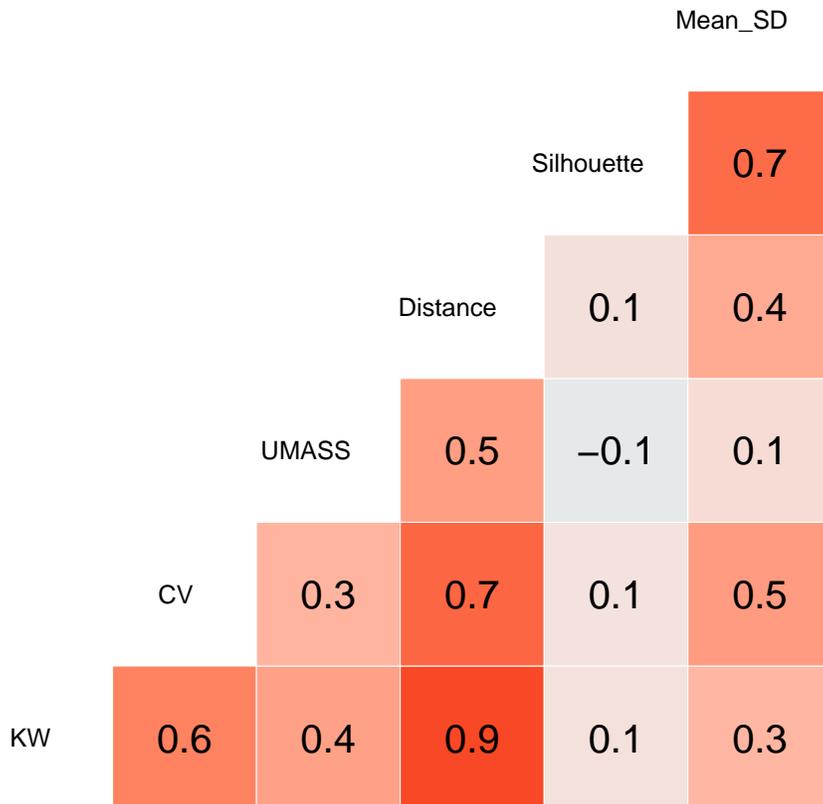


Figure 1: This figure shows the Pearson rank correlation between each of the automated metric scores for each of the 40 clusters.(KW is short for number of Keywords in a cluster)

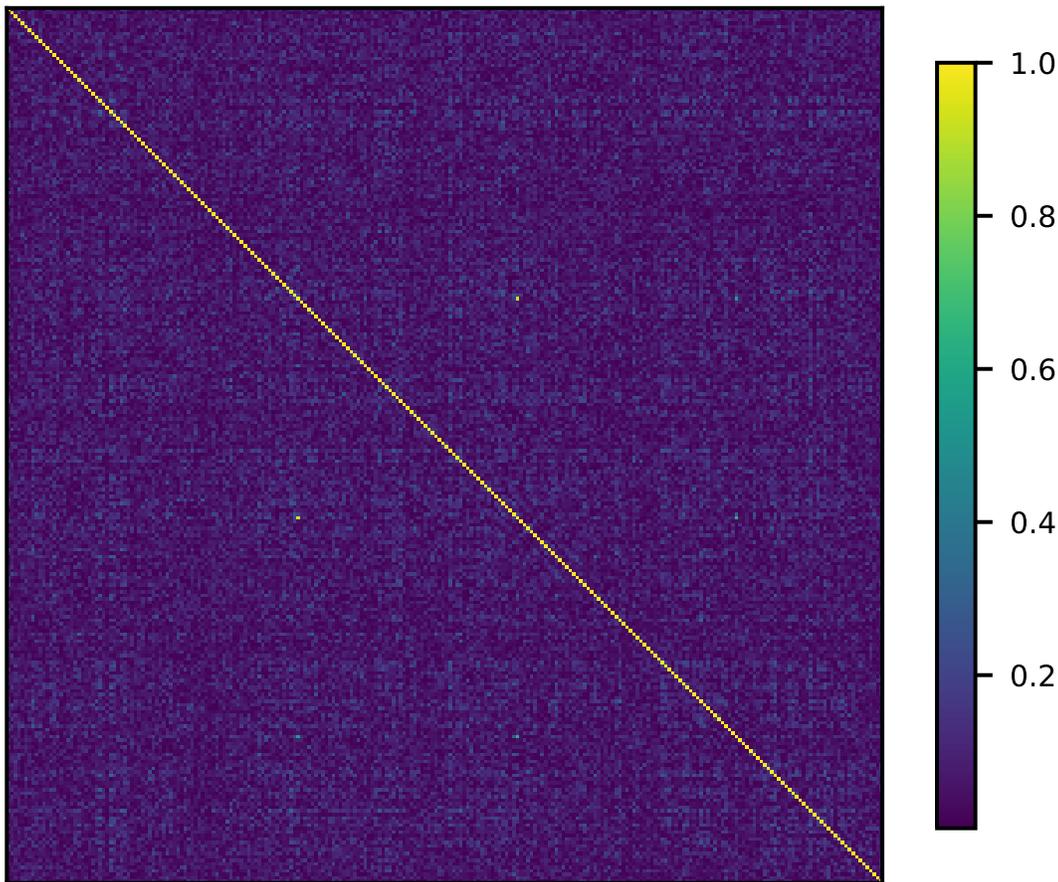


Figure 2: Pearson Correlation between all 384 vectors in the LLM

3 Stability

The three models used in the main paper each incorporate a stochastic component, which means that they can yield entirely different clusters with each run. To assess the stability of our methods, we executed them using various random seeds. For each seed, we calculated the pairwise Adjusted Mutual Information (AMI) between each seed. Our findings revealed that the LDA model had an average AMI of 0.25, with a standard deviation of 0.02. In contrast, the Doc2vec Model exhibited an average AMI of 0.91 and a standard deviation of 0.07, while the LLM showed an average AMI of 0.79 with a standard deviation of 0.09. The distribution of the AMI between each seed can be seen in In figures S4, S5, and S6. LDA appears to be a normal distribution. This seems to indicate that in all 50 seeds, none of them had similarly labelled clusters. This calls into question the reliability

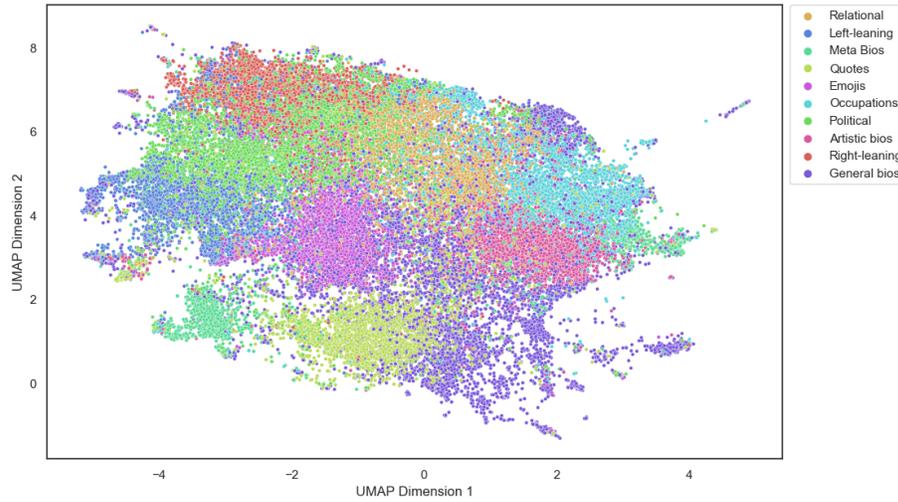


Figure 3: This figure shows all 38k Twitter bios projected onto a 2d plane using UMAP to show each points relation to each other, and which clusters are more similar to each other

of LDA and its implementation in short text. The GMM approaches seem to be quite a lot higher, both of them seem to have around 3 peaks (a small fourth in Doc2vec), potentially indicating that there are three local maxima that the algorithm finds. All three of them seem to have a high AMI indicating they are fairly similar. However, this needs to be further explored to test what tangible differences, if any, there are between these three maxima.

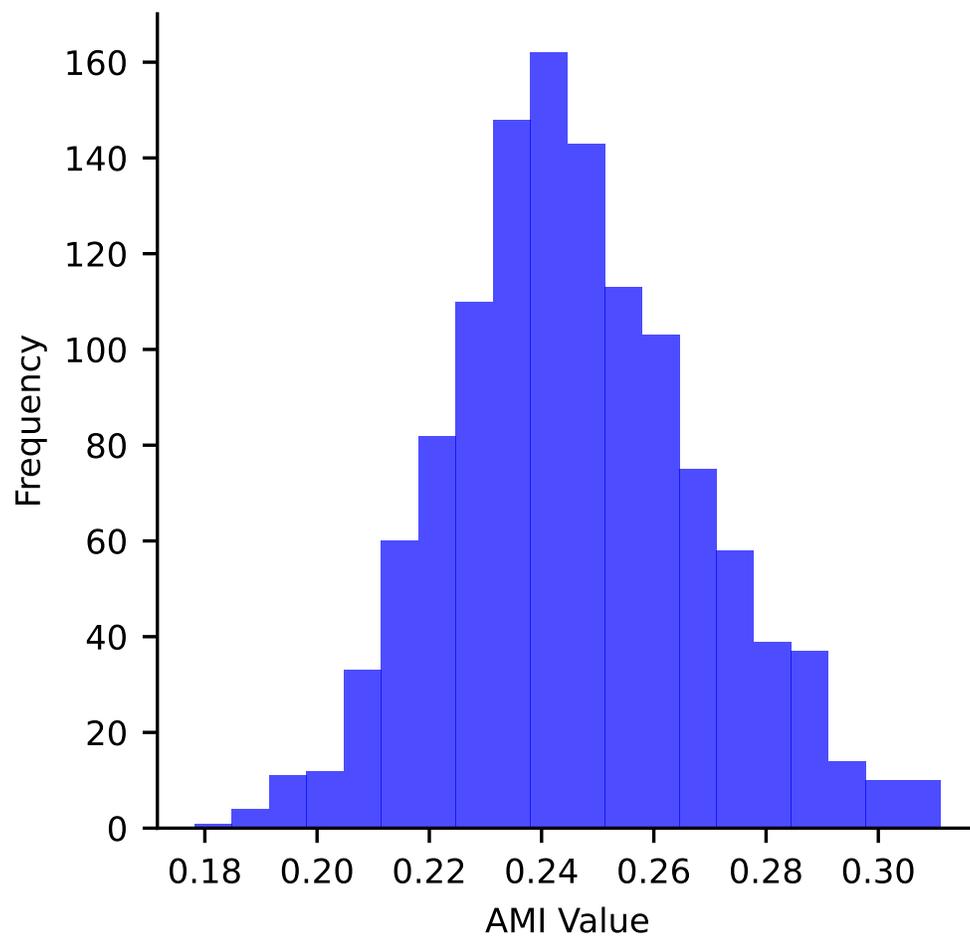


Figure 4: Distribution of pairwise AMI between 50 different seeds of LDA implementation

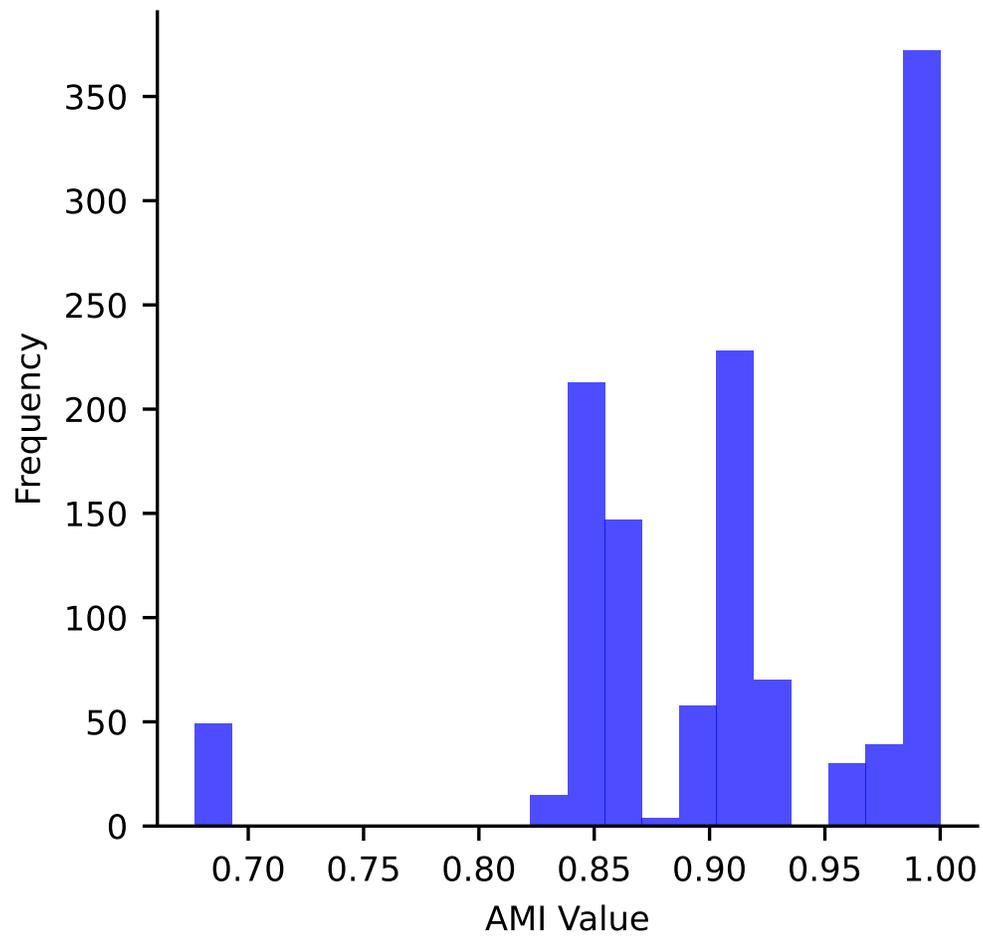


Figure 5: Distribution of pairwise AMI between 50 different seeds of GMM implementation using vectors created by Do2vec

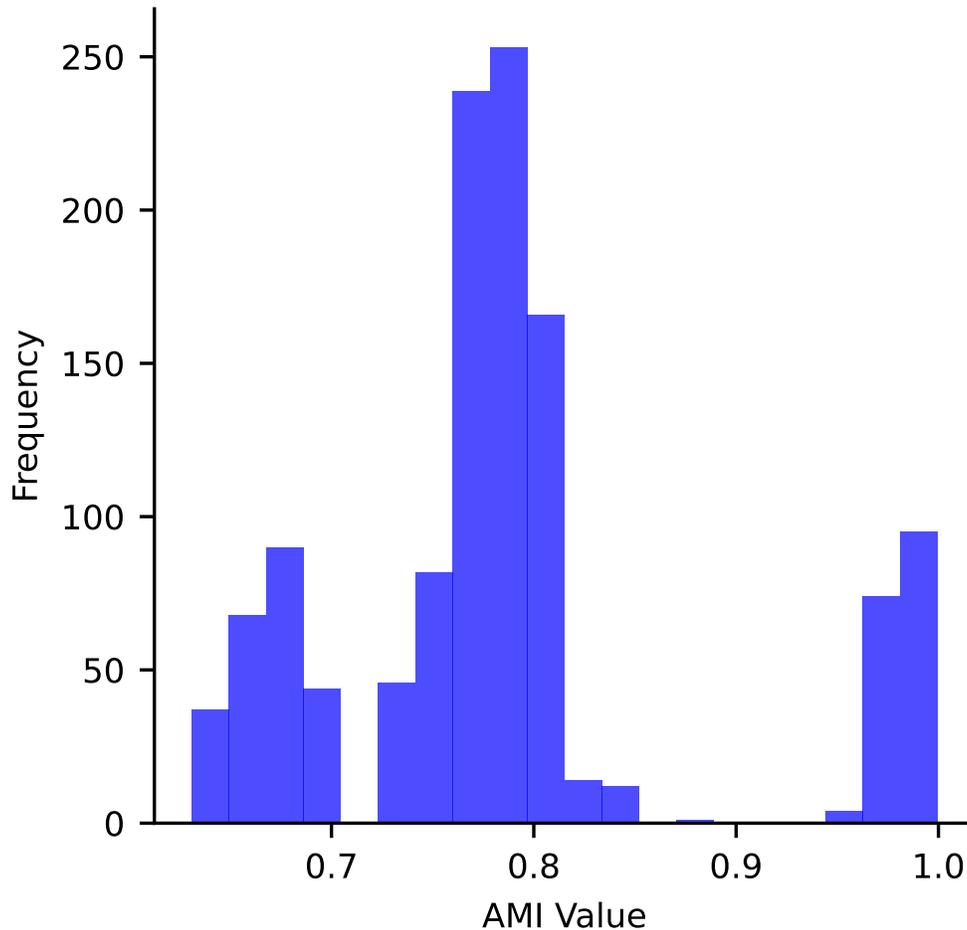


Figure 6: Distribution of pairwise AMI between 50 different seeds of GMM implementation using vectors created by a LLM

4 Other Models of ChatGPT

The use of ChatGPT allows us to explore effects which would be resource-intensive when using human reviewers. To evaluate ChatGPT’s performance with an increased number of bios, we conducted tests using 20, 100, and 500 bios as input, focusing on clusters that represented a range in diversity. This was done to assess the model’s scalability while managing costs, so we limited our tests to two clusters per method. The procedure was consistent with our earlier experiments using the ChatGPT API, with no alterations besides the number of bios, and using GPT-4 Turbo as it can take more bios [2]. After processing the bios through ChatGPT-4 Turbo, we applied the same analysis method from Section 3.8, of the method section in the main paper, to determine

the frequency of the most commonly used word in the generated names. Rather than naming all clusters, we choose a small subset of clusters from across the models. We can see in Fig.S 7 that increasing the number of bios does not appear to lead to a systematic change in the accuracy of the names produced by ChatGPT.

5 Process Pipeline

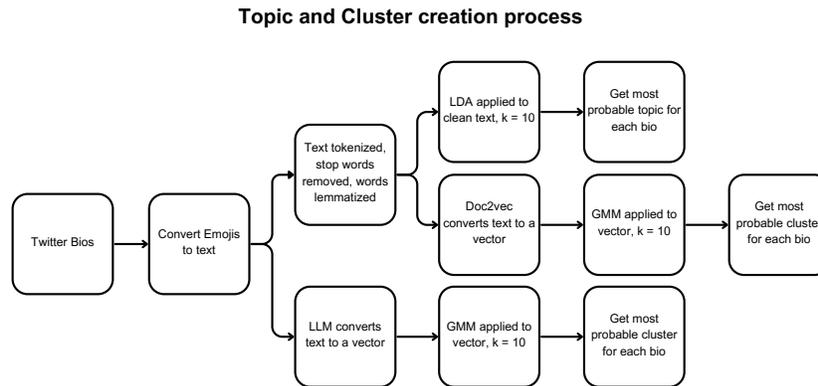


Figure 8: Diagram showing the process to building each model

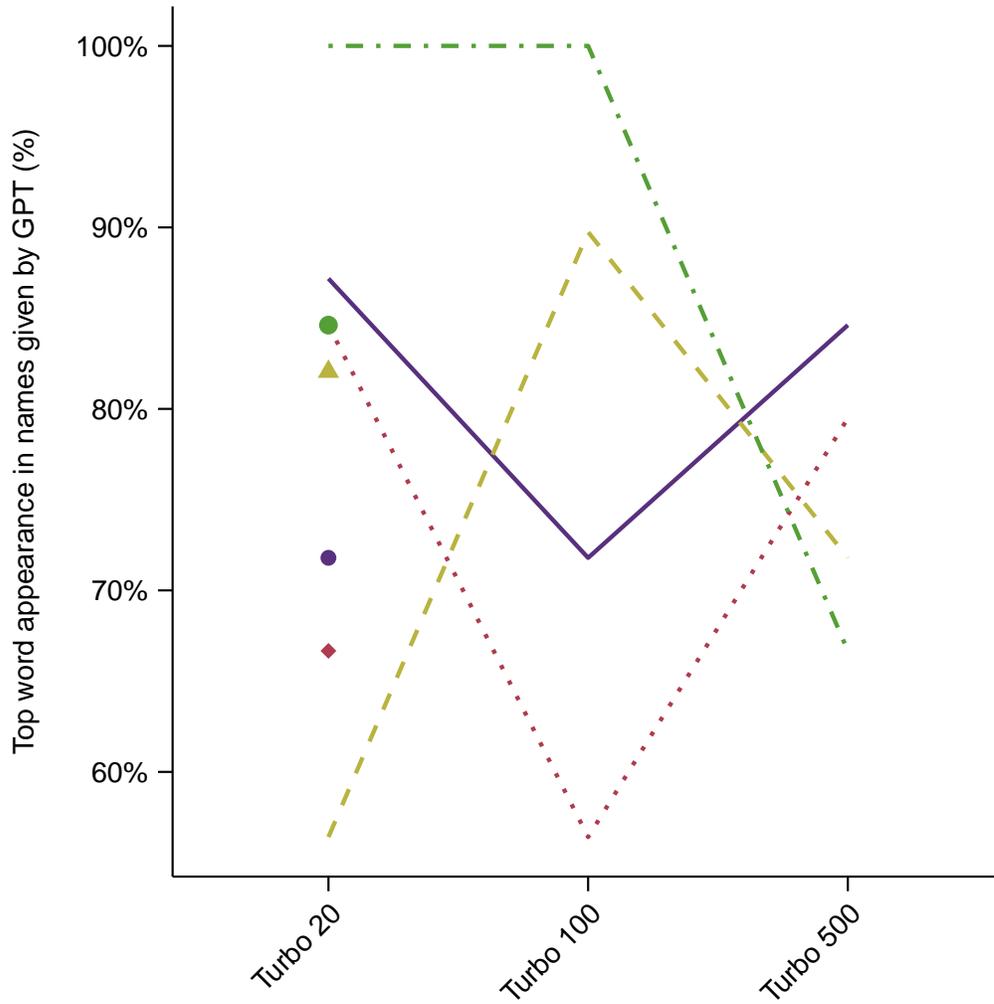


Figure 7: This figure shows the different levels of reliability between gpt-4 being given 20 bios (symbols), and gpt-turbo being given 20, 100, and 500 bios, and being asked to name it.

Process For Review Cluster Creation

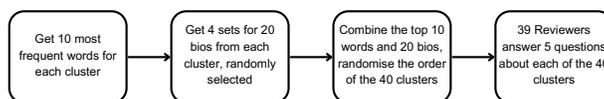


Figure 9: Diagram showing the process by which the cluster samples were created to show to our reviewers

The pipeline for how each cluster was created to show to our reviewers is in Fig. S9 The reviewers are shown each cluster in a random order and are not told which model a cluster came from. For each of the 40 clusters, reviewers must read the top 10 words and the 20 sample bios and then answer the following questions:

- “Create a name using less than 10 words to summarize the top 10 words/emojis and the sample bios of the Sample Cluster. If you believe it is not possible to do this with a cluster, write ‘None’.” This gives a summary of whether the cluster is human-interpretable and a description of the types of bios in the cluster.
- “When you named a Sample Cluster, were you confident that the name summarized the whole cluster?” Reviewers can answer this question by selecting one of the following: Not at all Confident/Not Confident/Neutral/Confident/Very Confident. This question helps understand how easy it was to understand each cluster.

For the following questions, reviewers can respond with: Strongly disagree/ Disagree / Neutral / Agree / Strongly Agree

- “‘The situation when the top words/emojis of a cluster fit together in a natural or reasonable way’ Do you agree that the above statement describes the top words/emojis of this Sample Cluster?” This question tests whether the top 10 words/emojis of a cluster are coherent.
- “‘The situation when the Sample bios of a cluster fit together in a natural or reasonable way’ Do you agree that the above statement describes the sample bios of this Sample Cluster?” This question tests whether the sample bios of each cluster are coherent.

- “The Top Words/Emojis provide an accurate summary of the sample bios’ Do you agree that the above statement describes this Sample Cluster?” This tests if the top words/emojis and sample bios are coherent with each other.

The data given by the reviewers was then ingested and cleaned. The responses to the questions above were converted to a 1-5 scale, where a 5 represents “Strongly Agree” or “Very Confident” and a 1 represents “Strongly Disagree” or “Not at All Confident”. For the confidence in naming question, some reviewers were unable to name the cluster but gave it a score of 4 or 5. In these cases, the data was changed so that every cluster that had “None” for its name was given a score of 1 for the question on confidence in naming.

This study had ethics approval from the University of Sydney. Reviewers were sourced from Prolific and responses were collected through Redcap. We collected 40 reviewers, however, had to exclude one due to them misunderstanding the prompt and copying and pasting bios as the names for each cluster, as well as a non-serious attempt at the other questions. An example of how a reviewer saw each cluster can be seen in figures S10 and S11

6 Inter-Coder Reliability

Inter-coder reliability measures the agreement between reviewers which is important given the subjective nature of their task. If there is low inter-coder agreement, this calls into question the interpretability of the clusters regardless of any other analysis. Since the data is ordinal and reviewers rate every single cluster, the appropriate method to use is the Intraclass Correlation Coefficient (ICC) [3, 4]. The ICC is a number between 0 and 1, with a higher number indicating greater agreement between reviewers. To provide context for the ICC scores, Koo et al. provided a guideline [5]: a score less than 0.5 is poor, a score between 0.5-0.75 is moderate, a score between 0.75-0.9 is good, and a score above 0.9 is excellent.

ICC estimates and their 95% confident intervals were calculated using R version 4.2.2 [6] and the statistical package Psych version 2.2.9 [7] based on an average rating, two-way random effects, absolute agreement, multiple raters/measurements. Table S1 shows the results of the results for each of the questions reviewers were asked.

Cluster 1

Please read the below top words and sample clusters and then answer the questions.

Top Words

Wife Mom Love Mother Husband Father Proud Married Family

Sample Bios

Father, husband, grower of my own facial hair. #AmericalsGreaterThanTrump Love each other, FFS!! #HateIsTheOnlyMinorityNow #BLM

Wife. Mom. LDS. Player of games, singer of songs, crosser of stitches, grower of plants. I mute Covid truthters. Ain't no follow-back girl.

Mom, wife, smartass, agnostic, curious learner, skeptical seeker, music lover, yogi + hiking, biking, skiing, GenX Colorado girl. 🌱 #Biden2020 #Resist #BLM

I am an immigrant. I am a mother. When I came to US seeking better opportunities for my children and myself I had to start my life from scratch. Being a

Married with 2 kids. Business owner of residential remodeling company Dallas Home Renovations.

Mom's Registered Caregiver | Licensed Psychotherapist #LMFT | Certified Personal Trainer & Health Coach | #MentalHealth Contributor | @cronkite_asu #DigDownDeep

Just trying to follow Christ. Wife to a man that only God could have chosen for me, mom of people who astound me daily, book lover, singer, business owner.

Mom of 2 amazing humans ❤️
Curious Georgette 🇺🇸 3rd Twitter Acct.
Marine Corps and Navy Daughterus
#NothingCanStopWhatsComing #HCQ+
#TRUMP2020 #WWG1WGA

I'm not a regular mom, I'm a cool mom.

Believer in fairness for all. Don't care about colour or sexuality. We're all the same underneath. Golden retriever lover 3 over 40 years. Mum of two great men.
....Married 1971..... 2-Sons...1- n-Heaven. 1-Daughter...1-Brother

Figure 10: Screenshot from Redcap of the reviewer being shown the contents of each cluster

<p>1) Create a name using less than 10 words which you think best captures the nature of the cluster, based on what you can infer from the top 10 words/emojis and the sample bios of the Sample Cluster. If you cannot identify a name for the cluster write "None".</p> <p><i>* must provide value</i></p>	<input type="text"/>
<p>2) When you named a Sample Cluster were you confident that the name summarised the whole cluster? If you answered "None" for the naming then please choose "Not at all Confident"</p> <p><i>* must provide value</i></p>	<p><input type="radio"/> Not at all Confident</p> <p><input type="radio"/> Not Confident</p> <p><input type="radio"/> Neutral</p> <p><input type="radio"/> Confident</p> <p><input type="radio"/> Very Confident</p>
<p>3) "The top words/emojis of a cluster fit together in a natural or reasonable way" Do you agree that the above statement describes the top words/emojis of this Sample Cluster?</p> <p><i>* must provide value</i></p>	<p><input type="radio"/> Strongly disagree</p> <p><input type="radio"/> Disagree</p> <p><input type="radio"/> Neutral</p> <p><input type="radio"/> Agree</p> <p><input type="radio"/> Strongly Agree</p>
<p>4) "The Sample bios of a cluster fit together in a natural or reasonable way" Do you agree that the above statement describes the bios of this Sample Cluster?</p> <p><i>* must provide value</i></p>	<p><input type="radio"/> Strongly disagree</p> <p><input type="radio"/> Disagree</p> <p><input type="radio"/> Neutral</p> <p><input type="radio"/> Agree</p> <p><input type="radio"/> Strongly Agree</p>
<p>5) "The Top Words/Emojis provide an accurate summary of the sample bios" Do you agree that the above statement describes this Sample Cluster?</p> <p><i>* must provide value</i></p>	<p><input type="radio"/> Strongly disagree</p> <p><input type="radio"/> Disagree</p> <p><input type="radio"/> Neutral</p> <p><input type="radio"/> Agree</p> <p><input type="radio"/> Strongly Agree</p>
<input type="button" value="Submit"/>	

Figure 11: Screenshot from Redcap of the reviewer being shown the questions and where they input their answers. Same cluster as the previous figure

Metric	Intraclass Correlation	95% Confidence Interval		F Test with True Value 0			
		Lower Bound	Upper Bound	Value	df1	df2	Sig
Confidence in Naming	.900	.0849	.94	13.5	39	1482	9.3×10^{-73}
Coherence of Top Words	.922	.883	.953	14.7	39	1482	1.3×10^{-79}
Coherence of Sample Bios	.911	.866	.946	13.7	39	1482	1.5×10^{-73}
Coherence between Top Words and Sample Bios	.926	.889	.955	15.7	39	1482	7.1×10^{-85}

Table 1: ICC Calculation in R Using Two-way random, average measure (ICC(2,k))

From this it can be seen that the ICC is considered excellent across all four questions, indicating that the reviewers were in relative agreement on their ratings of each cluster. A limitation of this method is that it assumes that the data is continuous, but the reviewer responses are ordinal.

7 Automated Methods

7.1 Automated Metric Results

The results for every cluster’s automated metrics can be found in table S2 Using the automated measures to judge and compare each of the 4 models, every automated metric except for C_{UMASS} agrees with the ranking given by human reviewers: LLM, LDA, Doc2vec, and Random. However, when looking at the ranking of the individual clusters there is a much greater difference. Overall it seems that all the automated methods fail when looking at individual clusters, but in a small sample of 4 models can rank them effectively.

7.1.1 Coherence

Both coherence methods found that LDA 1 was the most coherent, this was in contrast to the reviewers who were not able to name the cluster and gave on average quite low scores across the 4 questions. Doc2vec had comparatively good C_{UMASS} scores. In terms of C_V , it scored generally lower, with its highest score being for Doc2vec 4 which was one of the few that reviewers were able to consistently name. The LLM clusters that were political in nature (0,2, and 9) scored quite well on the coherence metrics, however, other clusters that were rated quite high by reviewers were given low scores. On the whole, it does not appear that either metric of coherence matches up with human interpretability.

7.1.2 Cluster Distance

The results from all models show low silhouette scores, with the majority of the scores close to zero. The reason for this is because of the high standard deviation in each dimension from the embeddings generated for the LLM and used to generate the silhouette score for all four models. Fig. S12 shows that with synthetic data, the clusters generated by a GMM are able to match the silhouette score of the true clusters. As standard deviation increases it can be seen that the silhouette score decreases even amongst the true label. So the low silhouette scores are not a reflection of bad clustering but rather high standard deviation in the data.

7.1.3 Keyword Analysis

Keyword analysis showed vast differences between the three models in terms of lexical content. On average, the transformer model created the clusters with the most keywords and highest average keyness scores. In more practical terms, this implies that the transformer model created the clusters that were most distinct from the original corpus of bios in terms of word frequency. In contrast, Doc2Vec was by far the worst-performing model, including several clusters with no keywords. This implies that the distribution of words in said clusters is near identical to the original corpus. Therefore, to make clusters that are lexically distinct from one another, the Transformer model performs the best.

Silhouette scores as a function of standard deviation: 384 dimensions

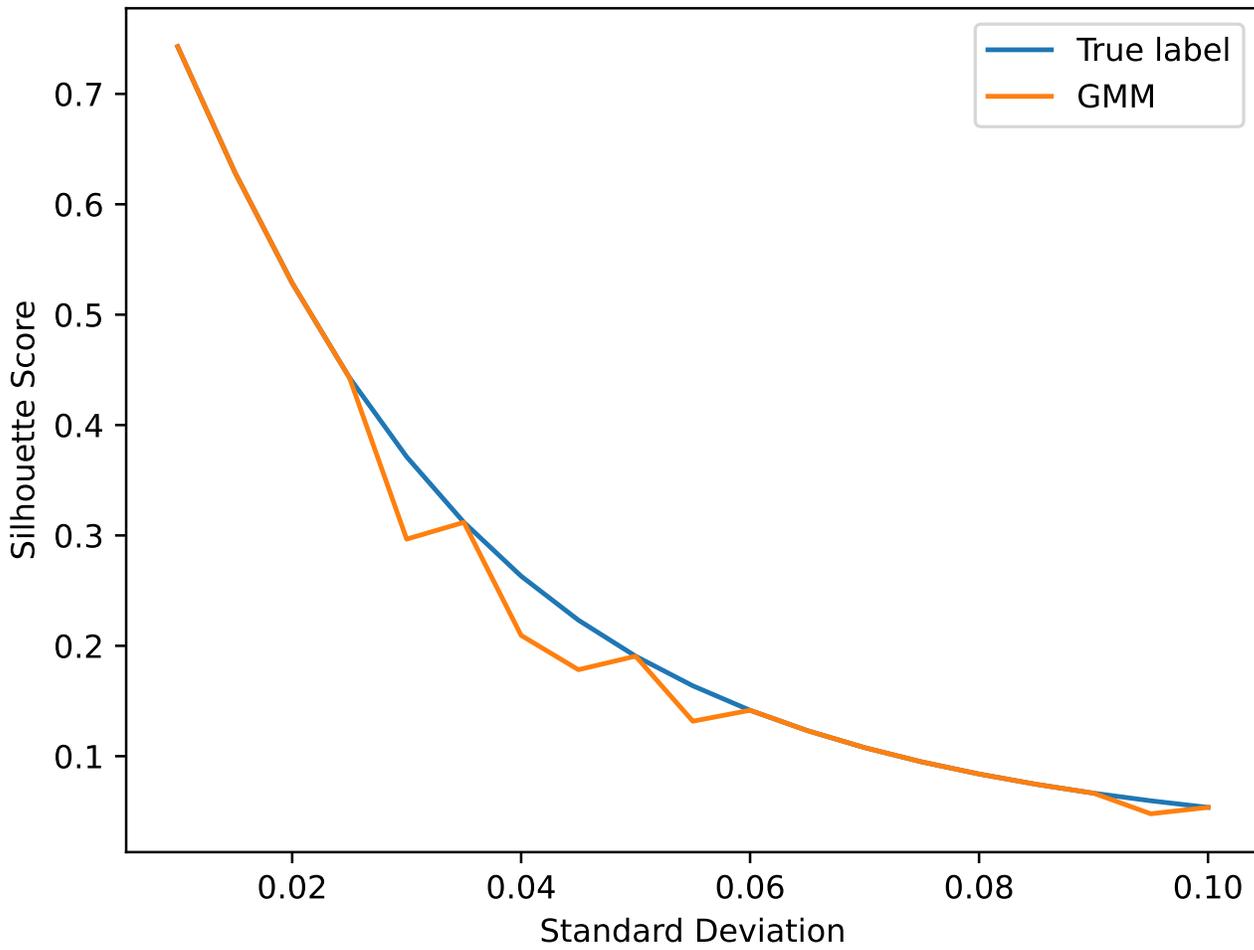


Figure 12: Silhouette Score as Standard deviation increases on Synthetic data

Cluster id	Name	GPT Name	Keywords	CV	UMASS	Distance	Silhouette	Mean SD
Doc2vec 0	U1	Political Views and Personal Attitudes	0	0.26	-3.5	0.012	-0.033	0.048
Doc2vec 1	U2	Minimal Information MAGA Supporters	32	0.2	-2.1	0.053	0.004	0.047
Doc2vec 2	U3	Diverse Personal Interests and Beliefs	5	0.29	-4	0.012	-0.035	0.047
Doc2vec 3	U4	Varied Interests and Political Opinions	62	0.33	-2.6	0.022	-0.024	0.047
Doc2vec 4	Political_1	Mixed Political Sentiments and Personal Interests	8	0.45	-3.3	0.002	-0.021	0.046
Doc2vec 5	U5	Mixed Political Views	0	0.26	-3.5	0.005	-0.026	0.047
Doc2vec 6	Political_2	Mixed Political Sentiments	0	0.34	-3.6	0.003	-0.023	0.046
Doc2vec 7	U6	Diverse Political Interests	0	0.27	-8.3	0.004	-0.010	0.046
Doc2vec 8	U7	Personal Interests and Political Affiliations	0	0.44	-8.8	0.004	0.011	0.045
Doc2vec 9	U8	Political Perspectives and Personal Interests	0	0.42	-4.4	0.002	-0.022	0.046
LDA 0	Left leaning_1	Activists and Biden Supporters	38	0.35	-2.7	0.086	-0.015	0.045
LDA 1	U9	Personal Perspectives and Interests	42	0.69	-4	0.082	-0.050	0.047
LDA 2	U10	Diverse Interests and Professions	39	0.25	-2.9	0.037	-0.060	0.048
LDA 3	Right Leaning	Patriotic Trump Supporters	39	0.62	-2.7	0.115	0.107	0.042
LDA 4	Relationships	Family and Occupation Focus	58	0.36	-2.8	0.076	-0.005	0.045
LDA 5	U11	Opinion Expressers	54	0.42	-3.3	0.049	-0.047	0.048
LDA 6	U12	Veterans and Varied Interests	32	0.31	-2.3	0.043	-0.059	0.047
LDA 7	U13	Varied Interests and Political Engagement	51	0.31	-2.7	0.052	-0.039	0.046
LDA 8	U14	American Sociopolitical Views and Beliefs	55	0.35	-2.8	0.061	-0.056	0.046
LDA 9	Left leaning_2	Life, Politics, and Personal Interests	60	0.35	-2.6	0.037	-0.062	0.047
LLM 0	Left Leaning	Democratic Supporters and Resisters	131	0.54	-2.6	0.222	0.028	0.041
LLM 1	Artistic Bios	Creative Professionals and Enthusiasts	230	0.4	-3.2	0.246	0.033	0.044
LLM 2	Political	Political Affiliation and National Pride	155	0.56	-2.8	0.246	-0.029	0.043
LLM 3	Occupations	Professionals and Students with Various Interests	264	0.29	-4.1	0.183	0.008	0.045
LLM 4	Relational	Family-oriented Patriotic Individuals	129	0.38	-2.9	0.304	0.074	0.042
LLM 5	General	Eclectic Personal Expressions	191	0.37	-7.7	0.222	-0.070	0.048
LLM 6	Twitter	Twitter-focused Political Commentators	133	0.45	-8.1	0.162	0.035	0.043
LLM 7	Quotes	Life Philosophies and Personal Beliefs	270	0.43	-7.7	0.334	-0.019	0.046
LLM 8	Right Leaning	Patriotic Trump Supporters	224	0.7	-1.9	0.291	0.160	0.038
LLM 9	Emojis	Polarized Political Users	216	0.48	-2.6	0.162	0.099	0.041
Random 0	U15	Political Interests and Personal Identities	1	0.05	-5.1	0.001	-0.002	0.046
Random 1	U16	Diverse Political and Personal Interests	1	0.08	-4.1	0.001	-0.003	0.046
Random 2	U17	Polarized Political Views	0	0.15	-5	0.001	-0.002	0.046
Random 3	U18	Political Standpoints and Diverse Interests	0	0.09	-4.6	0.001	-0.003	0.046
Random 4	U19	Mixed MAGA Supporters and Resisters	0	0.21	-4.2	0.001	-0.001	0.046
Random 5	U20	Mixed Political Views and Personal Interests	2	0.13	-4.3	0.001	-0.005	0.047
Random 6	U21	Diverse Personal Interests and Politics	1	0.18	-4.2	0.001	-0.002	0.046
Random 7	U22	American Politics and Personal Interests	0	0.25	-5	0.001	-0.005	0.047
Random 8	U23	Mixed Political Sentiments	0	0.24	-4.2	0.001	-0.003	0.046
Random 9	U24	Political Views and Personal Interests	0	0.14	-4.7	0.001	-0.001	0.046

Table 2: Table showing every clusters score for each automated metric, including how it was named by one example use of Chatgpt

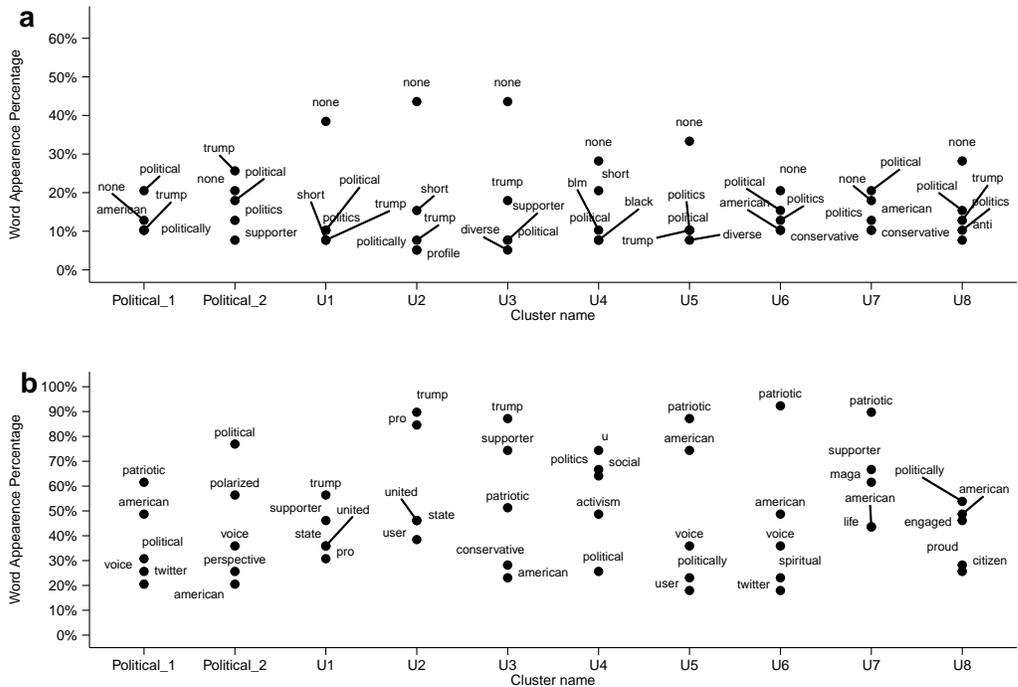


Figure 14: The top five words by fraction of appearance used by (a) reviewers and (b) ChatGPT to name the clusters created by Doc2vec. Along the x-axis are the names given to each cluster by the authors of this paper.

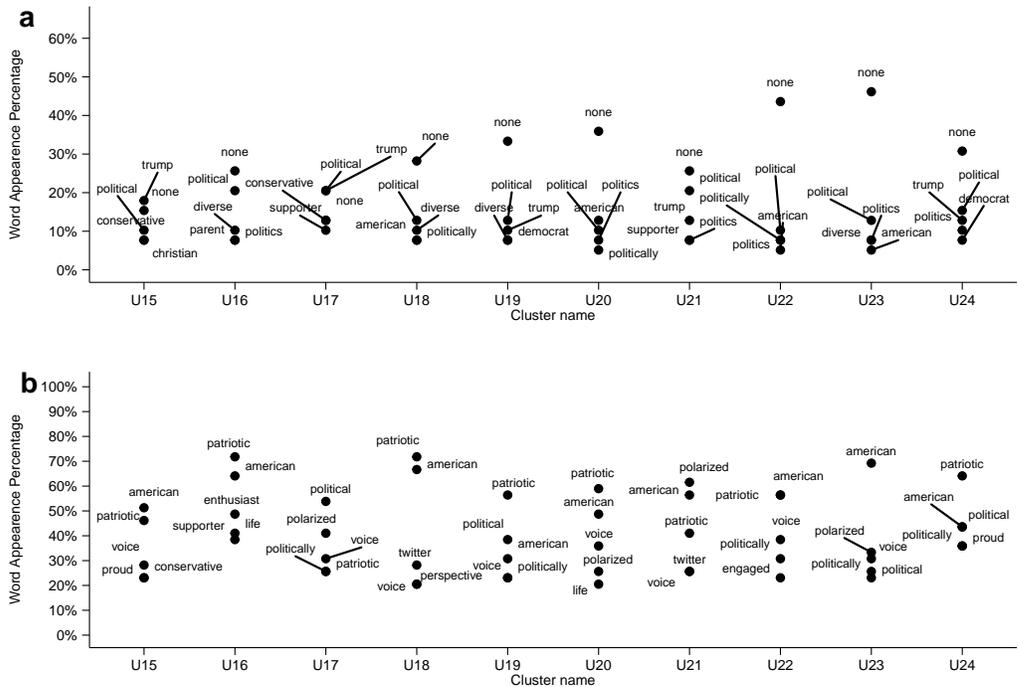


Figure 15: The top five words by fraction of appearance used by (a) reviewers and (b) ChatGPT to name the clusters created by the Random Model. Along the x-axis are the names given to each cluster by the authors of this paper.

Cluster ID	Cluster Name	Top 10 Frequent Words
LDA 0	Left leaning 1	Resist, Blm, Http, Blacklivesmatter, Matter, Life, Black, Bidenharris2020, Fbr, Theresistance
LDA 1	Undecided	New, Side, Party, Local, National, flag-mexico, Analyst, Troll, Old, General
LDA 2	U9	sparkle, Live, Sheher, Fuck, Im, Director, Life, One, 20, World
LDA 3	Right Leaning	flag-united-states, Maga, Trump, Patriot, Kag, Love, Trump2020, God, Wwg1Wga, red-heart
LDA 4	Relationships	Mom, Father, Husband, Wife, Retired, Proud, Mother, Lover, Fan, Dad
LDA 5	U11	Opinion, Like, Tweet, Endorsement, Hehim, Im, La, De, Http, View
LDA 6	U12	Army, Veteran, Vet, Retired, Follow, U, Twitter, Fan, Proud, Dont
LDA 7	U13	water-wave, Love, Politics, News, Music, Lover, Book, Dog, World, Resist
LDA 8	U14	Love, Trump, Truth, U, Country, Dont, God, America, Believe, Im
LDA 9	Left leaning 2	Vote, Im, Blue, Human, Right, Life, Fan, Living, Take, Time
Doc2vec 0	U1	flag-united-states, Maga, water-wave, Trump, Just, Trump2020, 2020, Kag, Wwg1Wga, Life
Doc2vec 1	U2	flag-united-states, Maga, Just, Https, Love, Wwg1Wga, Trump, Ig, Trump2020, Kag
Doc2vec 2	U3	flag-united-states, Maga, Trump, Love, Just, Trump2020, Patriot, Https, Conservative, God
Doc2vec 3	U4	flag-united-states, La, Blacklivesmatter, Black, Lives, Matter, Politics, Maga, Husband, News
Doc2vec 4	Political 1	flag-united-states, Trump, Love, Maga, red-heart, water-wave, Resist, God, Life, Proud
Doc2vec 5	U5	flag-united-states, Maga, Trump, Love, water-wave, Conservative, Proud, Life, Resist, Https
Doc2vec 6	Political 2	flag-united-states, Trump, Love, Maga, water-wave, God, Resist, Conservative, red-heart, Kag
Doc2vec 7	U6	flag-united-states, Trump, Maga, Love, water-wave, Proud, Resist, Mom, Life, God
Doc2vec 8	U7	flag-united-states, Trump, Love, Maga, water-wave, red-heart, God, Proud, Life, Mom
Doc2vec 9	U8	flag-united-states, Trump, Maga, Love, Resist, Proud, Conservative, God, Kag, Mom
LLM 0	Left Leaning	water-wave Resist, Blm, Bidenharris2020, Blacklivesmatter, Trump, Vote, Fbr, Resistance, flag-united-states
LLM 1	Artistic Bios	Writer, Lover, Music, Artist, Love, Https, Fan, Art, Author, Life
LLM 2	Political	Trump, Conservative, Love, flag-united-states, America, Proud, American, Country, Democrat, Liberal
LLM 3	Occupations	Retired, Business, Fan, Student, Https, Science, Veteran, Politics, Engineer, Father
LLM 4	Relational	Wife, Mom, Love, Mother, Husband, flag-united-states, Father, Proud, Married, Family
LLM 5	General	Just, Https, Fan, Love, Don, Like, Life, Time, flag-united-states, La
LLM 6	Twitter	Twitter, News, Tweets, Politics, Https, Follow, Tweet, Trump, Don, Account
LLM 7	Quotes	Life, Love, God, Truth, World, People, Don, Just, Good, Matter
LLM 8	Right Leaning	flag-united-states, Maga, Trump, Kag, Trump2020, Conservative, Wwg1Wga, God, Patriot, Love
LLM 9	Emojis	flag-united-states, red-heart, water-wave, blue-heart, Love, Maga, sparkle, Trump, star, God

Table 3: Top Words for all models

References

- [1] L. McInnes, J. Healy, and J. Melville, “Umap: Uniform manifold approximation and projection for dimension reduction,” *arXiv preprint arXiv:1802.03426*, 2018.
- [2] OpenAI, “New models and developer products announced at devday,” <https://openai.com/blog/new-models-and-developer-products-announced-at-devday>, November 2023, accessed: [2023-12-10].
- [3] N. Gisev, J. S. Bell, and T. F. Chen, “Interrater agreement and interrater reliability: key concepts, approaches, and applications,” *Research in Social and Administrative Pharmacy*, vol. 9, no. 3, pp. 330–338, 2013.
- [4] K. O. McGraw and S. P. Wong, “Forming inferences about some intraclass correlation coefficients.” *Psychological methods*, vol. 1, no. 1, p. 30, 1996.
- [5] T. K. Koo and M. Y. Li, “A guideline of selecting and reporting intraclass correlation coefficients for reliability research,” *Journal of chiropractic medicine*, vol. 15, no. 2, pp. 155–163, 2016.
- [6] R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2022. [Online]. Available: <https://www.R-project.org/>
- [7] W. Revelle, *psych: Procedures for Psychological, Psychometric, and Personality Research*, Northwestern University, Evanston, Illinois, 2022, r package version 2.2.9. [Online]. Available: <https://CRAN.R-project.org/package=psych>