

D4D: NSA

| | | | |
|------------------------------|--------------------------------|--------------------------------|------------------------------|
| 1 st Peter Kinder | 2 nd Taylor Hartley | 3 rd Caitlin League | 4 th Amir Ghayour |
| CU Boulder | CU Boulder | CU Boulder | CU Boulder |
| kinderp@colorado.edu | taha6964@colorado.edu | cale4879@colorado.edu | amgh1894@colorado.edu |

1 Introduction

Over the past four months, beginning in January of 2022, our team at CU Boulder has applied our diverse set of skills to address one of the most important issues currently affecting the globe: disinformation. Disinformation, a term many consider to have been coined by Joseph Stalin in 1923 [1], is defined as false information deliberately, and often covertly, spread in order to influence public opinion [2]. However, disinformation has existed in many forms prior to 1923, and in the past decade or so, with the growth of the internet and social media, seems to permeate all aspects of society including politics, business and entertainment. In fact, this deluge of disinformation has correlated to 64% of Americans feeling a great deal of confusion about the basic facts of current events [3]. While some threads of disinformation may not pose much risk to society at large, a sizable amount does. This in turn poses a threat to national security, which has prompted numerous government agencies and the intelligence community to dedicate substantial resources to addressing the problem.

One specific agency that is looking into this problem is the National Security Agency (NSA). Among the many ways that the NSA is seeking to address disinformation, one way it is doing so is by leveraging the academic community. In an annual competition called Hacking for Defense (known as Designing for Defense at CU Boulder), the NSA has presented teams of students with a challenge statement. The particular challenge that our team took on was entitled *Mapping the Disinformation Campaign* and was a challenge to:

Identify and analyze foreign disinformation campaigns in order to inform cybersecurity priorities.

This challenge was relatively broad, making narrowing the scope essential for our team. After months of research and consulting with our sponsors at the NSA, mentors in the class, and more than 75 interviewees with subject expertise, our team decided to approach the challenge in a somewhat literal sense of *mapping* the disinformation campaign. One of the pivotal papers early on in this process was entitled *Narrative Maps: An Algorithmic Approach to Represent and Extract Information Narrative* in which the authors explained an approach to algorithmically map narratives [4]. We believed that narratives were somewhat parallel in nature to campaigns, and that applying this approach to disinformation might yield valuable insight. Further along in the process, our team merged with another team that was tasked with a challenge similar to ours. This team had developed a machine learning algorithm that was able to determine whether individual pieces of information could be classified as organized (dis)information operations with up to a 95% accuracy level. Additionally, during our earlier research, we had come across numerous papers that were able to achieve similar results. In response, our team expanded our focus to cluster mapping as a way to identify distinct campaigns from the individual pieces of information that were tagged as being part of a general (dis)information operation by our partner's algorithm.

In the end, our team developed what we consider as two distinct products, cluster mapping and narrative mapping, to help combat disinformation. Some of the applications of our products would be better suited for the commercial world, but the applications discussed in this paper will focus on value to our sponsor, the NSA. For the remainder of this paper, the following sections will address:

- Cluster mapping of disinformation.
- Narrative mapping of disinformation.
- Discussion of alternative applications not related to our sponsor.
- Reflections on the challenge.

2 Cluster Mapping

The ecosystem of the internet is enormously vast, with approximately 4.66 billion individuals accessing the internet on a daily basis. Furthermore, approximately 4.2 billion individuals use social media daily, which is a relative hotbed for disinformation [5]. Found within all of this activity is the rampant spread of disinformation. So, even if one were able to theoretically detect every single piece of disinformation, sifting through the massive amount of data to extract understanding and know where to divert resources remains a challenge.

The objective of cluster mapping is to make sense out of all of this data in order to provide situational awareness and facilitate informed decision making. To explore this approach, our team used headline data from news sources determined to be sources of conspiracy or fake-news by a company named [Media Bias/Fact Check](#). Besides time and technical constraints, our team felt that approaching the problem this way would serve as a proof of concept and also be easier to understand and relate to for the reader. While we understand that much of the content used in this proof of concept if not disinformation per se, our method can be applied similarly to disinformation flagged by a machine learning algorithm.

2.1 Dataset

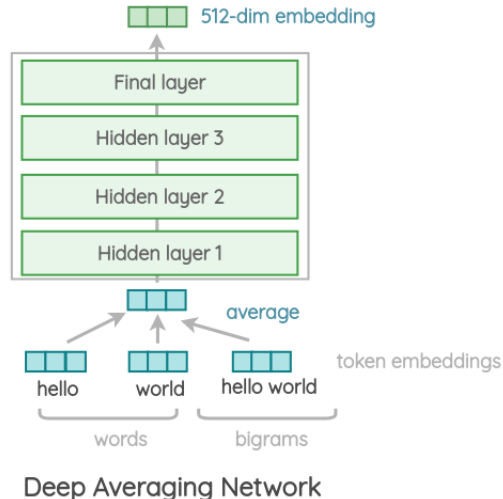
To obtain the headlines required for this proof of concept, our team aggregated data from [The GDELT Project](#). The GDELT Project is a vast trove of data with it's main focus on the world's news media. The data that was obtain ranged from between 01/01/2022 and 04/01/2022. Headlines from organizations determined to be sources of conspiracy or fake-news by Media Bias/Fact Check were the only headlines included. Within the dataset, it was determined that duplicate headlines existed from different sources. These duplicates were kept in the interest of trying to mirror how the same piece of disinformation can be shared by different sources. Additionally, some of the headlines contained text that was not relevant to the actual headline. For example, some headlines included text such as "<en>" and others included source signatures at the end, such as "– Society Child – Sott . net." To address this, our team employed regular expression (regex) to clean the text. An example of a headline before and after the cleaning can be seen below:

Before: Thousands march to honor WWII Nazi collaborator – Society Child – Sott . net

After: Thousands march to honor WWII Nazi collaborator

2.2 Embedding and Encoding

After the data was aggregated and cleaned, it needed to be converted to a numerical form so that clustering could be applied. To do this, our team embedded and encoded the text using an open-source model and python package. The model that we used was the [Universal Sentence Encoder Version 4](#), which was trained with a deep averaging network encoder (subset of deep learning).



The result was a 512 dimension numerical vector that represented each individual sentence or headline. Having a numerical representation of each headline allowed for a computational comparison between the headlines. However, while we chose to utilize the Universal Sentence Encoder Version 4 due to the respectable reputation the model had earned, working with a 512 dimension vector for clustering was not particularly efficient. Due to this, a python package called Uniform Manifold Approximation and Projection (UMAP) was used to reduce the vectors down to two dimensions, which could be interpreted as positional coordinates in an xy-plane. In short, the embeddings would be fitted to the UMAP reducer and then transformed into two dimension (further detailed explanation on this process, please visit this [link](#)). At this point, all of the headline text was converted to two dimensional vectors and ready to be clustered.

2.3 Clustering

There are many different ways to apply clustering, including hierarchical, centroid, distribution, or density based approaches. The approach our team took was largely based on the idea that disinformation, and information in general, is a very broad and constantly evolving space. To account for this, we utilized a clustering algorithm that took a minimum cluster size as a parameter instead of a predefined number of clusters. Based on this minimum cluster size, a number of clusters were “surfaced” on each run of the clustering algorithm. The name of this approach is HDBSCAN and is a rather complex and involved process. For further reading on this approach, please visit this [link](#).

If one were to run the clustering algorithm on a dataset one day, and then run it the next day on a slightly expanded dataset, the clusters would be different due to the relative nature of clustering. So, this part of our solution only established the baseline for cluster mapping and could be repeated on a daily, weekly, or monthly basis to identify new emerging clusters. In our proof of concept, employing this algorithm yielded results that were similar to those seen below.

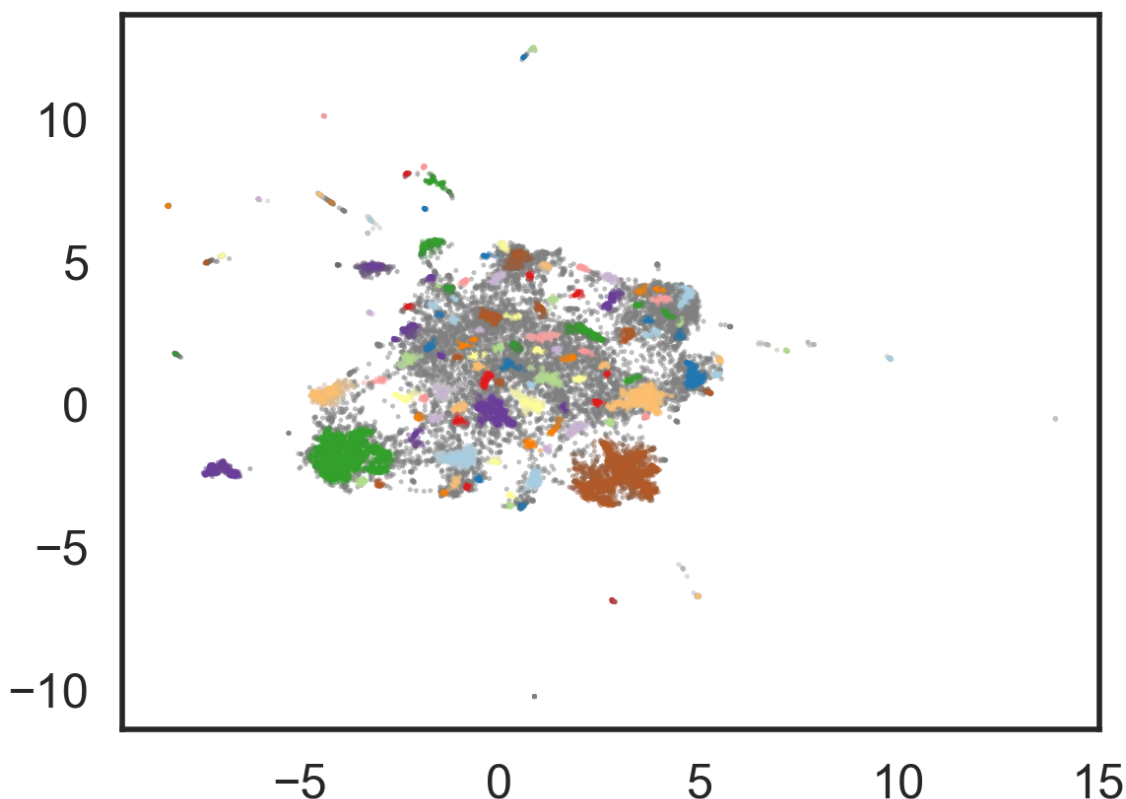


Figure 1: Depicts clusters of two dimensional numerical representations of headlines. The gray points represent headlines that don’t belong to any cluster and are considered noise.

It should be noted that due to the lack of available color palettes, some colors are repeated. But, clearly there are numerous separate clusters that have been identified. For example, the brown cluster near $(3, -2)$ and the green cluster near $(-4, -1)$ are two large individual clusters. The takeaway is that the headlines that correspond to those point and are part of those clusters are similar to each other. This could be roughly interpreted as them belonging to a certain story, narrative, or campaign. However, numerous gray points in the visualization, representing noise, don't belong to a cluster. That noise could contain important information, and needs to be assigned to a cluster for mapping.

To assign this noise to a cluster, the [HDBSCAN soft clustering approach](#) was used. Essentially what this does is instead of simply assigning a point to a cluster or deeming that it is noise, the algorithm assigns a probability vector (summing to 1) to each point. For example, a probability vector would look like such:

$[0.87, 0.03, \dots, 0.01, 0.00]$

The probability represents that likelihood that the point belongs to a particular cluster. Then, by assigning each point to the cluster for which it has the highest probability of belonging, every point is assigned to a cluster. For the above example, the headline would be assigned to the first cluster, for which is has an 87% probability of belonging to. Resultantly, the cluster visualization now is represented as such:

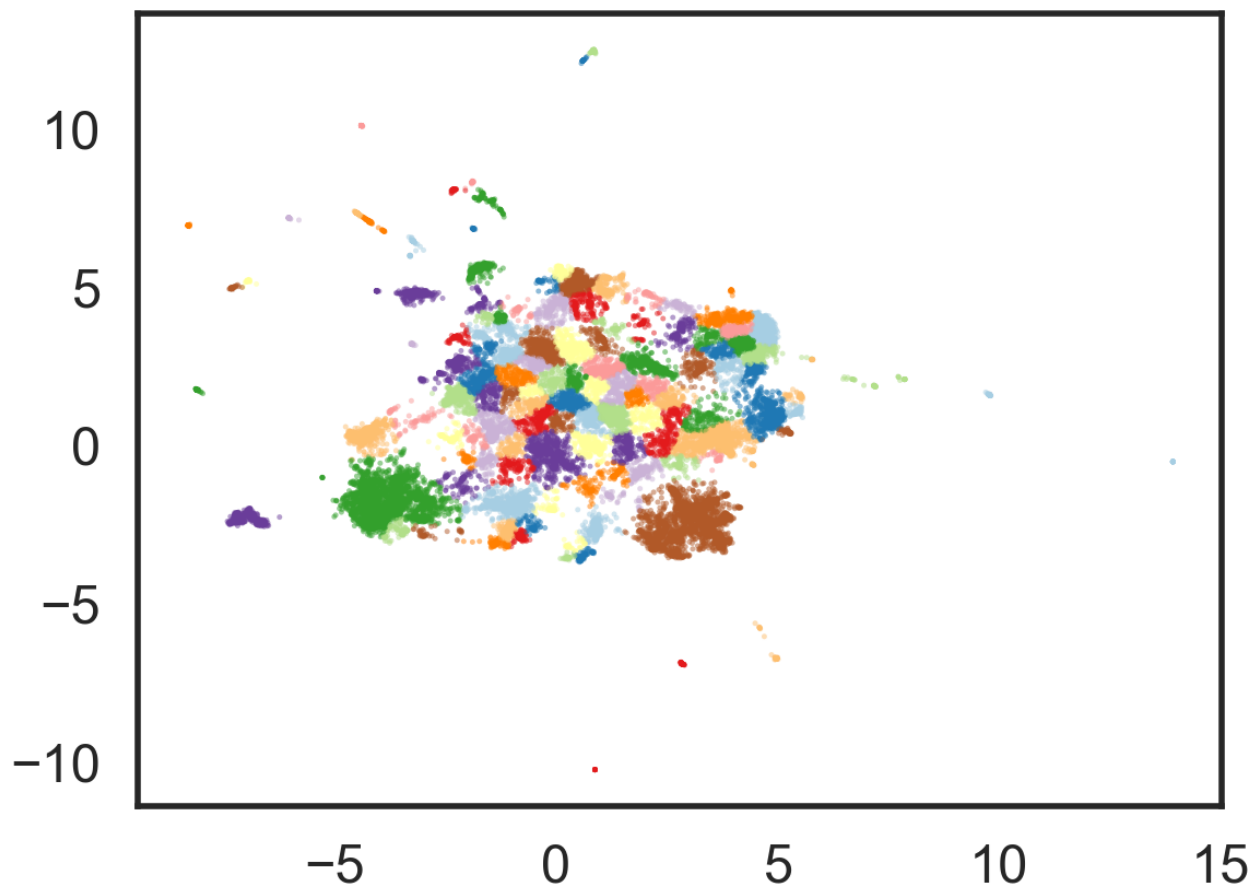


Figure 2: Depicts soft clusters of two dimensional numerical representations of headlines. Each point, representing a headline, is assigned a probability vector representing the probability that it belongs to a certain cluster. The cluster that each headline has the highest probability of being a member is the cluster that the headline is assigned to. Using this process, no headline is considered noise.

At this point, the foundations of the clusters have been established and will be referred to as baseline clusters. These will serve as a point of reference for mapping the identified clusters and extracting information. As noted earlier, this process can be repeated daily, weekly, or monthly, and is primarily focused on identifying emerging clusters to begin to establish situational awareness.

After the baseline clusters has been identified, measuring how the cluster evolves is the next critical step. Key to this process is assigning new headlines to the baseline clusters. To do this, the rough centers of the baseline clusters are determined by calculating the mean of the coordinates of each numerical headline representation within the specific cluster. Then, the distance between the numerical representation of the new headline and all of the rough centers is calculated. The new headline is then assigned to the cluster with the minimum distance between its rough center and the new headline. Then, by assigning that new headline to the cluster for which it is closest to, the analyst could monitor the evolution and other metrics of the cluster. After calculating the rough centers of each cluster, they were plotted on top of the previous visualization in the interest of understanding.

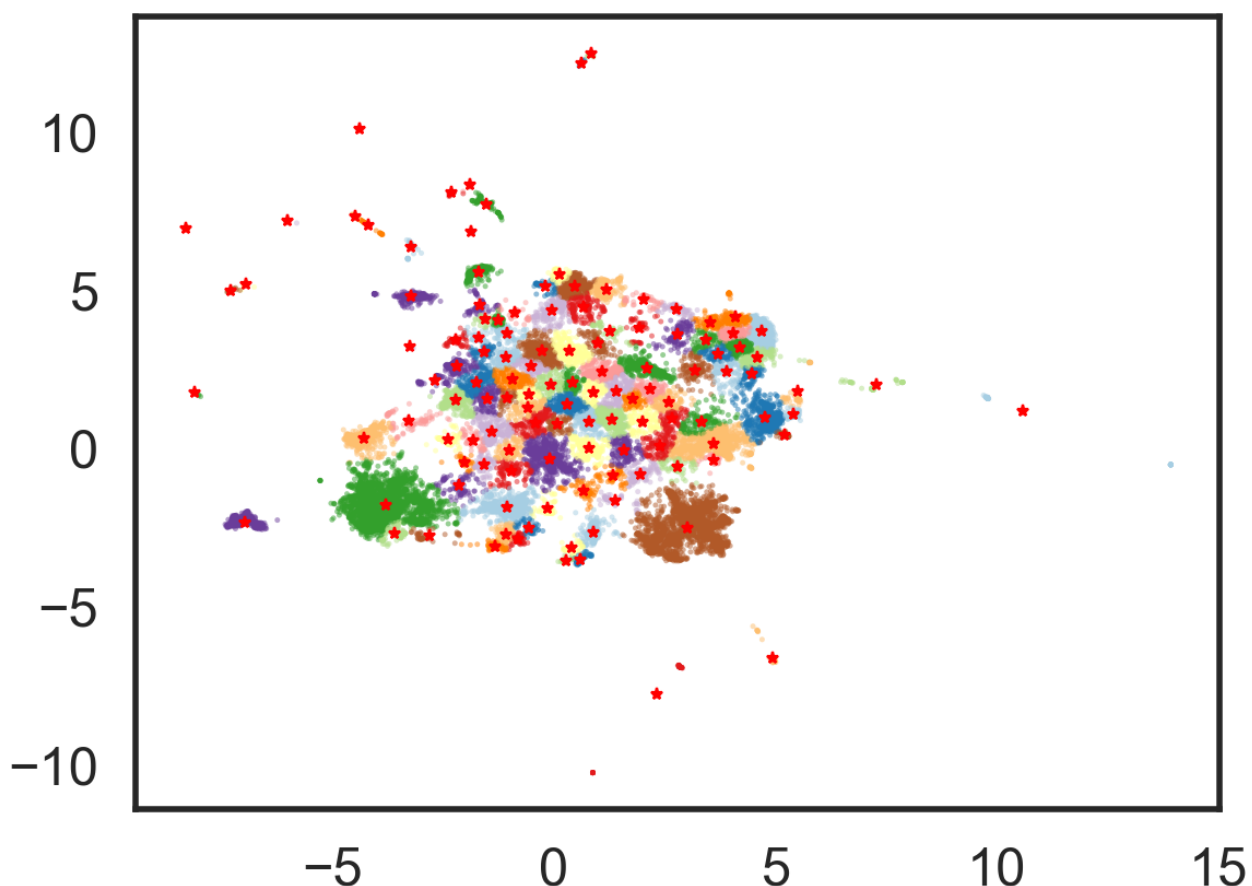


Figure 3: Depicts soft clusters of numerical representations of headlines with the centers of clusters represented by red stars.

Furthermore, the data from the next week following the baseline clusters in the dataset was plotted in blue, with the centers as well, to provide a visualization of how those particular clusters grew in the following week. This process would essentially happen in real-time, but providing the entire next week is easier to visualize and understand. Based on the lack of growth observed in the respective cluster positions, perhaps an analyst could conclude that the brown cluster was less critical to manage, while the green cluster needed urgent attention.

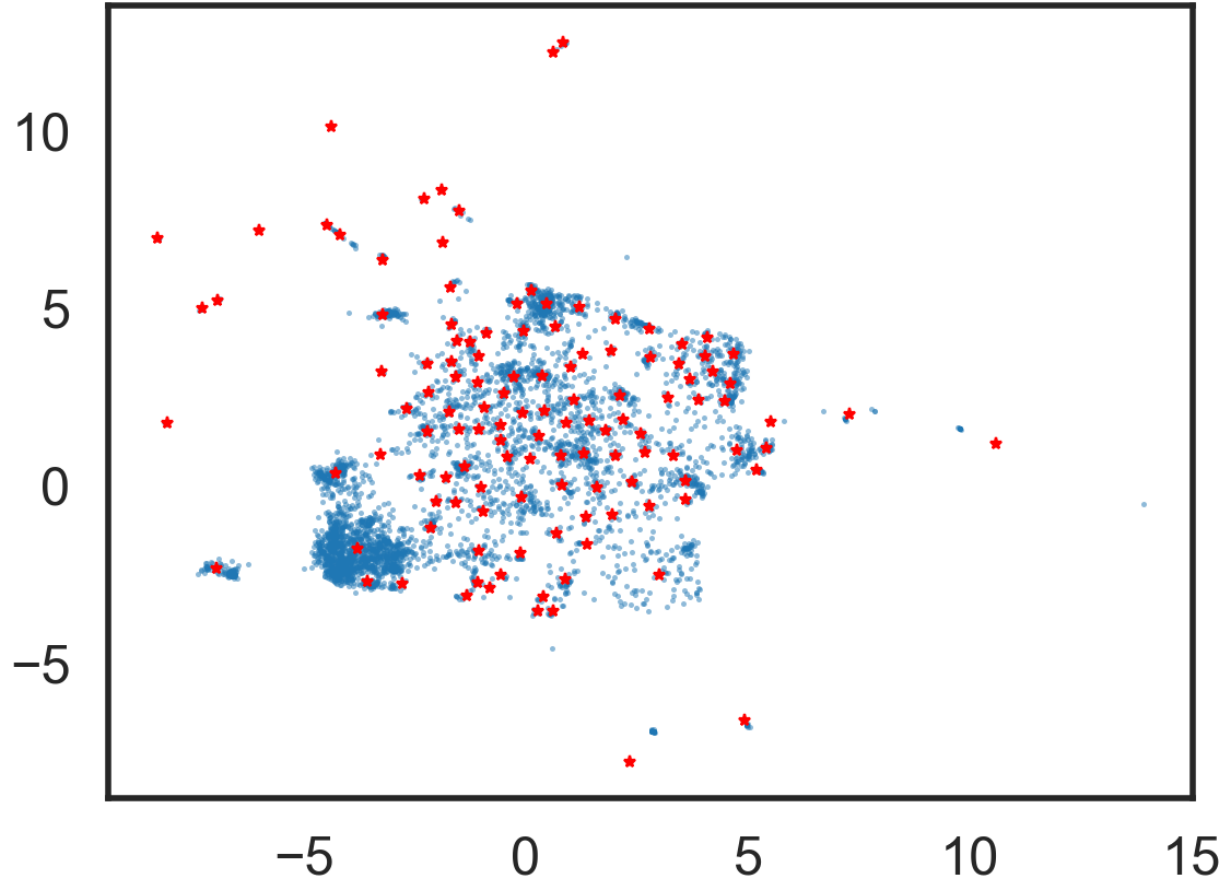


Figure 4: Depicts numerical representations of headlines from the week following the baseline clusters. As can be seen, the cluster where the large green cluster was in Figure 3 did not grow significantly. This could suggest that methods to address this cluster have been successful or that the cluster no longer needs to be monitored. The cluster where the large brown cluster was in Figure 3 did grow significantly. This could suggest that the methods to address the cluster haven't been as successful as hoped or that the cluster needs to continue to be monitored.

2.4 Information Extraction

At this point numerous clusters have been identified, but questions surrounding cluster evolution over time, themes of the clusters, effectiveness of counter strategies, and whether/where more resources should be devoted remain. To address these questions, our team identified a number of ways to extract insight from the clusters. The first general way is by measuring the growth of the clusters. One way to do this is measuring the daily count of disinformation. When identifying a cluster, the daily count of the baseline cluster can be used to help inform whether resources should be devoted to monitoring the cluster. For example, if the cluster is growing quickly this could be cause to monitor the cluster.

After it has been decided to monitor the cluster, the daily count of disinformation could be used to assess whether the steps taken to address it are effective. For example, if the daily count of the new disinformation in the cluster is trending down, this could suggest that the steps taken have been effective. Additionally, if the daily count of new disinformation falls to low enough levels, this could suggest that the cluster no longer needs to be monitored.

Below is an example of the the previously described scenario. The blue bars represent the counts from the baseline cluster, and the orange bars represent the counts from monitoring the cluster. The black arrows are superimposed to highlight the trends an analyst would observe.

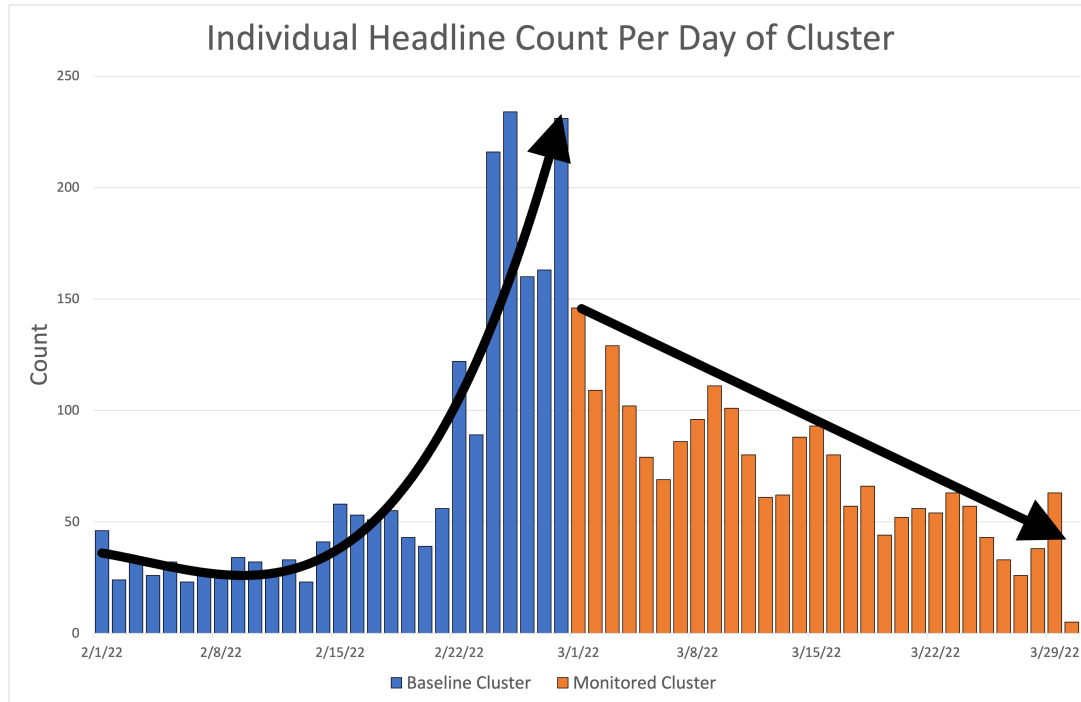


Figure 5: Depicts the daily count of individual headlines in a cluster.

Another way to visualize the trends of a cluster is the week-on-week average growth rate. Below of an example of this using the same data from the chart in Figure 5. As can be seen, the growth rate of the cluster before the baseline cluster was identified on the 28th of February was around 100%, or almost doubling. After it was identified and monitored, the growth rate dropped precipitously leveling out at around under 20%.

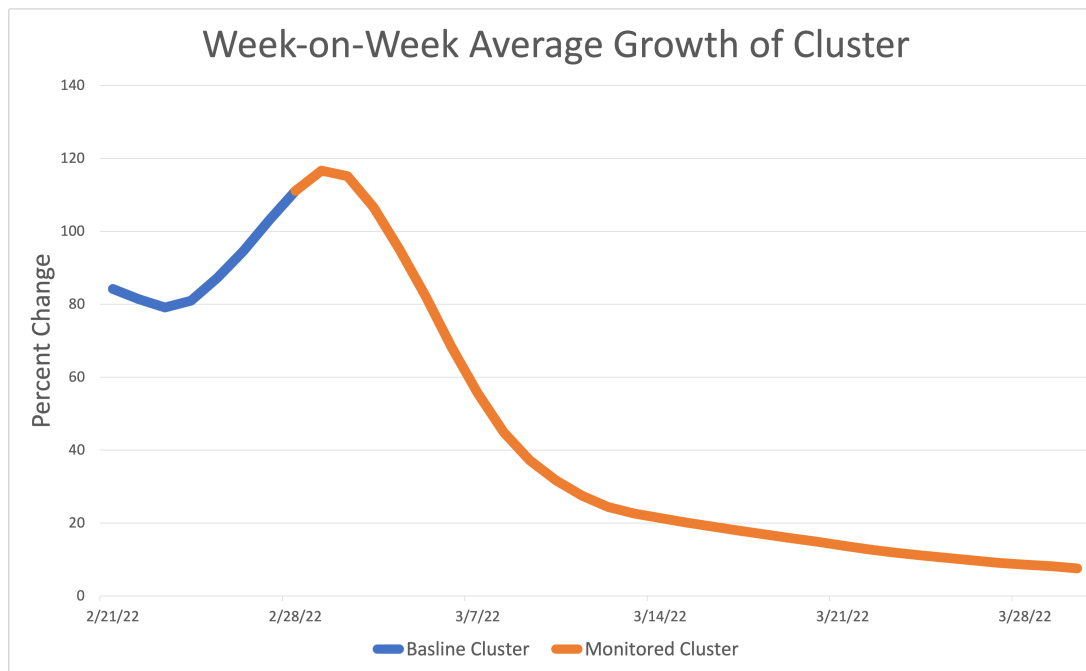


Figure 6: Depicts the average count of the trailing week to the week previous to that.

Another important consideration for these clusters are the sources that are disseminating most of it. By clustering the data together, counts can be made the associated sources. This data could be used in numerous ways, and one example that comes to mind is how RT.com and other state-aligned actors have been delisted or banned in certain countries. In regard to the application of our product, we thought that a word cloud could serve as a quick way to develop a situational awareness of the main sources involved. Example word clouds for two different clusters would look something similar to this:



Figure 7: Depicts two word clouds for two different clusters. The size of the text reflects the prevalence of that specific source. As can be seen, domains such as breitbart.com and theepochtimes.com frequently appear in both clusters. However, some domains are more associated with a specific cluster. For example, canadafreepress.com or snsnews.com frequently appear in one cluster, but not in another.

This process could be applied to specific websites, social media users, or other types of sources as well. This process could also be applied to the actual textual content of the headlines. One way we approached this was by tagging the parts of speech of the headlines. To do so, we used [FlairNLP](#); a python package for natural language processing. After tagging each word for what part of speech it was, we also lemmatized the words using the [NLTK](#) python package. Essentially, lemmatization is putting words into a more generalized form, and was useful for counting and visualizing them in the word clouds. Example word clouds from two different clusters generated from proper nouns would look similar to this:



Figure 8: Depicts two word clouds for two different clusters. The size of the text reflects the prevalence of that specific proper noun. As can be seen, proper nouns like Ukraine, Russia and Putin are prevalent in one cluster, while proper nouns like Freedom Convoy, Canada, and Trucker are prevalent in the other. Perhaps an analyst could quickly conclude that the one of the themes is of greater importance than the other and allocate resources accordingly. Additionally, an analyst could quickly use this to determine whether sub clusters are needed. For example, Taiwan and Xi appear in one cluster and could warrant a more focused attention.

Most readers are aware of the themes of the above clusters, but imagine a disinformation campaign that an analyst is just coming across. The players, themes, and locations are some elements that aren't immediately apparent. Quickly learning this type of information and developing situational awareness can help in the effort to combat the campaign. In another example, here are word clouds of the verbs in the clusters:



Figure 9: Depicts two word clouds for two different clusters. The size of the text reflects the prevalence of that specific verb. As can be seen, verbs like warn and invade are prevalent in one cluster, while verbs like freeze and mandate are prevalent in the other. Perhaps an analyst could quickly conclude one theme is more related to the military and the other is more political and assign those clusters to teams with greater specialization in those realms. Additionally, an analyst could perform a search of the headlines in the cluster that included some of those keywords to gain a more granular understanding of the dominant themes.

Another valuable piece of information that could be extracted from the clusters is the inferred exposure of the cluster. In tandem with the identified sources of the disinformation, information such as average website visits from Google Analytics and metrics like followers on social media can be used to roughly gauge how many people are being exposed to the disinformation cluster. This in turn could be used to determine where to focus the most effort on combating the disinformation.

2.5 Clustering Conclusion

Clustering disinformation can be a valuable tool for developing situational aware used to combat the spread of disinformation. In a perfect world, enough resources would be available to combat all of the existing disinformation out there. But since resources are limited, clustering can help reduce the harm. The next step for this product would be to perform this method on disinformation tagged from machine learning algorithms. We believe that although some adjustments will need to be made, demonstrating this process on headline data suggests that the process will provide similar value for more traditional disinformation.

3 Narrative Mapping

The narrative mapping product is in the same realm as clustering mapping, but instead of focusing on a cluster, focuses on the relation of a particular piece of disinformation to the rest of the ecosystem. The intent for narrative mapping is to develop a greater understanding of disinformation that can help to detect it in a number of way, such as machine learning, fact checking and crowd sourcing. The process involved a lot of the similar methods discussed in the clustering section, so this section will focus on specific aspects related to narrative mapping.

With a similar motivation to the clustering proof of concept, this approach was based off headlines from news outlets. However, the dataset was not limited to sources deemed to be conspiracy or fake-news by Media Bias/Fact Check.

3.1 Dataset

To obtain the headlines required for this proof of concept, our team pulled data from The GDELT Project. The GDELT Project is a vast trove of data with the main focus on the world's news media. The data that was obtain ranged from between 2022-01-01 and 2022-04-01. Within the dataset, it was determined that duplicate headlines existed from different sources. In this particular case, the duplicates would not be helpful in mapping the narrative. So, the first occurrence from each headline was the headline that

was kept. Additionally, the headlines included text that needed to be cleaned in a similar manner to the clustering approach.

After this dataset was constructed, features needed to be extracted from the headline for when the process began. The main reason was for computational efficiency, because while the process could be ran using the entire dataset each time, doing so took a much longer time. So, from each headline, keywords were extracted by utilizing a python library called [YAKE](#). This is essentially an unsupervised machine learning algorithm, and it was determined that only unigrams would be extracted. Then, after the keywords were extracted, the words were stemmed using the [NLTK](#) python package. The stemming of a word is essentially reducing it down to its most fundamental state. This was done because words that are essentially the same are used in different ways, and a computer cannot tell that they are the same if they are spelled differently. All of this allowed a subset of the headlines dataset to be retrieved and computed during the actual process.

3.2 Cosine & Angular Similarity

The dataset would be continuously updated as new headlines are published, and would serve as the foundation for the narrative maps that are generated. When a narrative map originating from a particular headline started to be generated, the keywords would be extracted and stemmed. Then, headlines from the dataset with common keywords would be pulled into a subset. At this point, all of the headlines would be encoded and embedded into numerical vectors as was done for the clustering. However, instead of clustering the data, the cosine and angular similarity between the originating headline and the headlines in the subett would be computed. This was done based on the 512 dimension vector from the Universal Sentence Encoder Version 4 in the interest of greater precision. However, it is helpful to consider two dimensional vectors when thinking about cosine similarity. Essentially, the cosine in the first quadrant of a 2 dimensional space ranges between 0 and 1. In the case of comparing two vectors that represent text, if the cosine is 1, then the two text items are exactly the same. If the cosine is 0, then the two text items are extremely dissimilar. It can be helpful to visualize this, imaging that each arrow (vector) represents a piece of text:

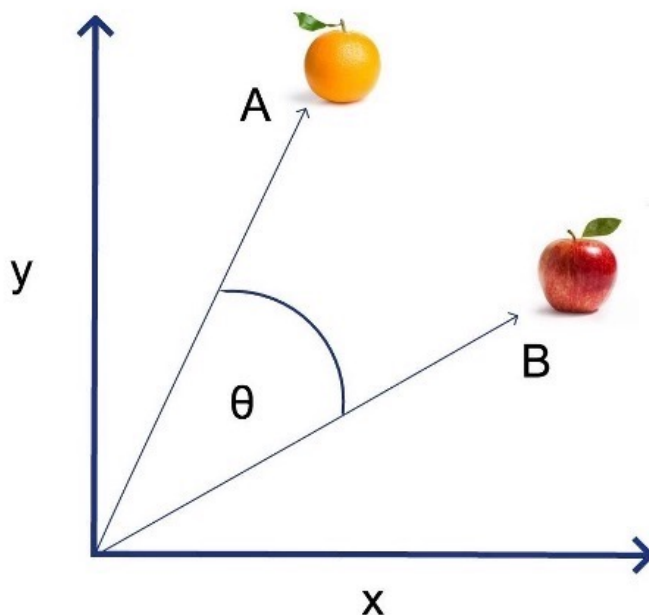


Figure 10: Depicts two dimensional numerical representation of the words “apple” and “orange.” These are similar to each other, but not the same thing. This is represented by the cosine of the angle θ not being 1, but maybe being somewhere around $\frac{\sqrt{3}}{2} \approx 0.87$ by just eyeballing it.

After the cosine similarity scores between the originating headline and headlines in the dataset were computed, the angular similarity was derived. This was done because many of the headlines were similar to the originating headline, and the angular similarity is more effective at distinguishing near parallel vectors from one another. This score would serve as one of the main factors in determining which headlines to include in the narrative map.

3.3 Biases

After the similarity between the actual text was determined, the biases of the sources were introduced. This can be incorporated in different ways, one of which we will specifically talk about in the [Discussion section](#). The main reason for this the anecdotal understanding that sources with different biases tend to present different counter narratives. In this proof of concept, political bias was the bias that was considered.

The main use of this was to compute a cohesive narrative instead of just carrying through the bias of the originating headline. The intent was to see how the originating headline fit into the overall narrative surrounding the originating headline. For example, if the original headline was determined to be from a source with a bias of “center”, then one narrative included in the map would be from “left” or “left-center” sources, while the other narrative would be from “right” or “right-center” sources. While this is not a perfect approach, and our team recognizes that the categorization of biases is likely biased in itself, incorporating this aspect yielded the best results by our judgement after viewing many generated narrative maps.

3.4 Sentiment Analysis

Another factor that we used when computing the narrative maps was sentiment analysis using [FlairNLP](#). The main reason for doing so was to differentiate between headlines that had very close similarity scores. However, the weighting of the sentiment was relatively low in the algorithm due to discovered biases in the sentiment score. For example, we noticed that one headline received an extremely negative sentiment when it included the word transgendered, but an extremely positive sentiment when the word was removed. Overall, the sentiment analysis score was primarily used as a tie breaker in edge cases.

3.5 Time Period

The last consideration was the timing of each headline. This particular aspect of the process can be adjusted as need be, but for the commercial aspect of our product we approached in the following way. It was determined that we did not want to overload an observer with too much information, so nine headlines were to be included in each map. Additionally, headlines are only relevant for so long, so it was determined that a month’s worth of time provided sufficient coherence to the narrative. Then, a single headline in one, two, three, and four week time periods are included to provide coherence. A sample result is as follows:



3.6 Narrative Mapping Conclusion

While we think that the clustering approach provides clearer value in the context of our sponsor’s problem, we are optimistic that narrative mapping can provide additional insight to the detection of disinformation. However, we do think that some of the more apparent value of narrative mapping can be found in the commercial sector and as an out of the box solution for our partner’s sponsor. We will touch on both of these ideas in the Discussion section.

4 Discussion

While the focus of the products has been in the context of our sponsor’s problem, we believe that additional applications exist. Three additional applications that will be laid out include a commercial application of cluster mapping, narrative maps as a covert survey for assessing the effectiveness of information operations, and narrative maps to increase the effectiveness of crowd sourced fact checking.

4.1 Commercial Application

As seen in the proof of concept for cluster mapping, emerging narratives can be identified in the media ecosystem. This ability could be of value to companies that want to manage their reputation. If a cluster emerges that is related to a specific company, that cluster can be tracked. Additionally, with the information extracted from the cluster, a counter campaign can be constructed. The effectiveness of the counter campaign can then be measured by assessing whether the growth rate of the cluster declines. Furthermore, the growth rate can be compared against the actual sales of the company to determine whether the narrative is impacting sales.

4.2 Covert Survey

In reference to the image of the narrative map on page 10, an alternative version of this map could be generated with one narrative composed only of state-sponsored media and the other composed of free media. Users could hover over the individual nodes for a bit more information, click on the nodes to read the actual story, or share the story with their peers. The way that users interact with the map could be measured, that in turn could suggest which narrative is more compelling to the user. This process could be repeated using numerous narrative maps over time to gauge how opinion is changing over time. Aggregating this data for many users could then be used to inform individuals who are conducting information operations whether the content that they are disseminating is effective.

4.3 Crowd Sourcing

One of the ways that the veracity of information is determined is through professional fact checkers. For example, one prominent organization doing this work is [Snopes](#). However, there are a couple downsides to this approach. First, there are a limited number of people in the world considered to be professional fact checkers, and their services are relatively expensive. Additionally, there is a sizable portion of any population that consider professional fact checkers as biased against them. Lastly, a professional fact-checker usually cannot immediately determine whether a piece of information is factual, and there is disagreement among themselves.

In response to some of these issues, the solution of crowd sourcing fact-checking has been investigated. It is based upon the concept of the wisdom of crowds, which is the idea that the aggregate opinion of enough people can match that of experts. One major study on this subject, *Scaling Up Fact-Checking Using the Wisdom of Crowds*, found that at a crowd size of around 10-15 people, the correlation of veracity ratings of a piece of information matched or exceed the correlation among three professional fact checkers [6]. In the study, the members of the crowd made their assessment, on average, in approximately 30 seconds using only a headline and lede. The narrative map could be used in place of the headline and lede to provide additional context without greatly increasing the time required to make an assessment. This could potentially make the process require a lower group size to achieve the same correlation, resulting in cost savings and a greater economic feasibility to scale the process.

5 Reflection

When reflecting back on our project, we feel that we have created the foundations for what could be a very useful tool in combating disinformation, both from the perspective of the government and in a commercial application. However, we wanted to provide some feedback on some challenges we encountered that we feel might be helpful for the next team accepting a similar challenge. The three main themes of the challenges we encountered were the scope of the challenge, the impact of classification, and the political aspect.

5.1 Scope

Perhaps the biggest challenge of the three was the broad scope of disinformation. While our team had a general understanding of this from the onset, the first couple of weeks was largely related to the narrowing of the scope. In our opinion, this was beneficial to our understanding of the problem, but detracted from our ability to develop a product that delivered value. Maybe guidance related to a certain approach, such as machine learning, clustering, or fact-checking would have been helpful. Instead of determining the approach, our team could have investigated a particular approach and whether we could improve upon it. We think that we could have developed a more valuable product if this had been the case.

5.2 Classification

Another large challenge was the impact of the classified nature of certain aspects of the challenge. Disinformation is a very important topic, and discussions with subject matter experts was often difficult due to this. Additionally, obtaining data related to this challenge was rather difficult because disinformation isn't as effective if it is known to be disinformation. Maybe if a parallel problem that wasn't as restricted as the topic of disinformation was presented, some of these issues would have been alleviated. In the end, our team approached this problem using headlines as the parallel to disinformation due to some of these issues.

5.3 Political

Determining what is disinformation can be quite easy in some instances, but extremely nuanced in others. On top of that, disinformation is often blended into misinformation, which is more politically charged than disinformation. This made some of the research more difficult, and slowed down the process.

References

- [1] “Disinformation.” Wikipedia, Wikimedia Foundation, 22 Apr. 2022, <https://en.wikipedia.org/wiki/Disinformation>.
- [2] “Disinformation Definition & Meaning.” Merriam-Webster, Merriam-Webster, <https://www.merriam-webster.com/dictionary/disinformation>.
- [3] Barthel, Michael, et al. “Many Americans Believe Fake News Is Sowing Confusion.” Pew Research Center’s Journalism Project, Pew Research Center, 27 Aug. 2020, <https://www.pewresearch.org/journalism/2016/12/15/many-americans-believe-fake-news-is-sowing-confusion/>.
- [4] Keith, Brian, and Tanushree Mitra. “Narrative Maps: An Algorithmic Approach to Represent and Extract Information Narratives.” ArXiv.org, 26 Oct. 2020, <https://arxiv.org/abs/2009.04508.pdf>.
- [5] Johnson, Joseph. “Internet Users in the World 2021.” Statista, 26 Apr. 2022, <https://www.statista.com/statistics/617136/digital-population-worldwide/>.
- [6] Scaling up Fact-Checking Using the Wisdom of Crowds. <https://www.science.org/doi/10.1126/sciadv.abf4393>.