

Probabilistic Model of Narratives Over Topical Trends in Social Media: A Discrete Time Model

Toktam A. Oghaz¹, Ece C. Mutlu^{2,*}, Jasser Jasser^{1,*}, Niloofar Yousefi², Ivan Garibay^{1,2}
 {toktam@cs,ece.mutlu@,jasser.jasser@,niloofar.yousefi@,igaribay@}.ucf.edu

¹Department of Computer Science,

²Department of Industrial Engineering,
 University of Central Florida

ABSTRACT

Online social media platforms are turning into the prime source of news and narratives about worldwide events. However, a systematic summarization-based narrative extraction that can facilitate communicating the main underlying events is lacking. To address this issue, we propose a novel event-based narrative summary extraction framework. Our proposed framework is designed as a probabilistic topic model, with categorical time distribution, followed by extractive text summarization. Our topic model identifies topics' recurrence over time with a varying time resolution. This framework not only captures the topic distributions from the data, but also approximates the user activity fluctuations over time. Furthermore, we define significance-dispersity trade-off (SDT) as a comparison measure to identify the topic with the highest lifetime attractiveness in a timestamped corpus. We evaluate our model on a large corpus of Twitter data, including more than one million tweets in the domain of the disinformation campaigns conducted against the White Helmets of Syria. Our results indicate that the proposed framework is effective in identifying topical trends, as well as extracting narrative summaries from text corpus with timestamped data.

CCS CONCEPTS

• **Computing methodologies** → *Information extraction; Topic modeling*; • **Information systems** → **Data mining**.

KEYWORDS

Topic Modeling; Graphical Models; Narrative¹ Extraction; Topic Detection and Tracking; Extractive Text Summarization

1 INTRODUCTION

Social media and microblogging platforms, such as Twitter and Facebook, are becoming the primary sources of real-time content regarding ongoing socio-political events, such as United States Presidential Election in 2016 [11], and natural and man-made emergencies, such as COVID-19 pandemic in 2020 [9]. However, without the appropriate tools, the massive textual data from these platforms makes it extremely challenging to obtain relevant information on significant events, distinguish between high-quality and unreliable content [17], or identify the opinions within a polarized domain [13].

The challenges mentioned above have been studied from different aspects related to topic detection and tracking within the field

of Natural Language Processing (NLP). Researchers have developed automatic document summarization tools and techniques, which intend to provide concise and fluent summaries over a large corpus of textual data [20]. Preserving the key information in the summary and producing summaries that are comparable to human-created narratives are the primary goals of the extractive and abstractive approaches for automatic text summarization [2]. News websites are a prime example of such techniques, where automatic text summarization algorithms are applied to generate news headlines and titles from the news content [31]. The shortage of labeled data for text analysis has encouraged researchers to develop novel unsupervised algorithms that consider co-occurrence of words in documents as well as emerging new techniques such as exploiting an additional source of information similar to Wikipedia knowledge-based topic models [37, 38]. Additionally, unsupervised learning enables training general-purpose systems that can be used for a variety of tasks and applications as strong classifiers [7]. In this regard, statistical models of co-occurrence such as Latent Dirichlet Allocation (LDA) [6], discover the relevant structure and co-occurrence dependencies of words within a collection of documents to capture the distribution of topic latent variable from the data. Although an abundant timestamped textual data, particularly from social media platforms and news reports are available for analysis, the changes in the distribution of data over time have been neglected in most of the topic mining algorithms proposed in the literature [35]. For instance, time-series analysis on datasets over the events relative to 2012 US presidential election suggests that modeling topics and extracting summaries without considering the text-time relationship lead to missing the rise and fall of topics over time, the changes in terms of correlations, and the emergence of new topics [12].

Although continuous-time topic models such as [35] have been proposed in the literature, topical models with continuous-time distribution cannot model many modes in time, which leads to deficiency in modeling the fluctuations. Additionally, continuous-time models suffer from instability problems in the case of analyzing a multimodal dataset that is sparse in time.

In this paper, we propose a probabilistic model of topics over time with categorical time distribution to detect topical recurrence, designed as an LDA-based generative model. To achieve probabilistic modeling of narratives over topical trends, we incorporate the components of narratives including named-entities and temporal-causal coherence between events into our topical model. We believe that what differentiates a narrative model² from topic analysis and summarization approaches is the ability to extract relevant sequences

¹In this study, we have used the terms narrative and story interchangeably.

*Equal contribution.

²In this paper, we refer to event-based topic modeling over topical trends as narrative modeling.

of text relative to the corresponding series of events associated with the same topic over time. Accordingly, our proposed narrative framework integrates unsupervised topic mining with extractive text summarization for narrative identification and summary extraction. We compare the identified narratives by our model with the topics identified by Latent Dirichlet Allocation (LDA) [6] and Topics over Time (TOT) [35]. This comparison includes presenting numerical results and analysis for a large corpus of more than one million tweets in the domain of disinformation campaigns conducted against the White Helmets of Syria. The collected dataset contains tweets spanning 13 months within the years 2018 and 2019. Our results provide evidence that our proposed method is effective in identifying topical trends within a timestamped data. Furthermore, we define a novel metric called significance-dispersity trade-off (SDT) in order to compare and identify topics with higher lifetime attractiveness in timestamped data. Finally, we demonstrate that our proposed model discovers time localized topics over events that approximates the distribution of user activities on social media platforms.

The remaining of this paper is organized as follows: First, an overview of the related works is provided in Section 2. In Section 3, we provide a detailed explanation of our proposed method followed by the experimental setup and results. Finally, in Section 5 we conclude the paper and discuss future directions.

2 BACKGROUND AND RELATED WORK

In this section, we first provide a background on narrative analysis and how literature has investigated stories in social media. Then, we present an overview of topic modeling and text summarization.

2.1 Narrative analysis

Narratives can be found in all day-to-day activities. The fields of research on narrative analysis include narrative representation, coherence and structure of narratives, and the strategies, aim, and functionality of storytelling [22]. From a computational perspective, narratives may relate to topic mining, text summarization, machine translation [33], and graph visualization. The later can be achieved via using directed acyclic graphs (DAGs) to demonstrate relationships over the network of entities [15]. Narrative summaries can be constructed from an ordered chain of individual events with causality relationships amongst events, appeared within a specific topic [18]. The narrative sequence may report fluctuations over time relative to the underlying events. Additionally, the story-like interpretation of the text is a must to imply a narrative [25].

Since social media have been admitted as a component of today's society, many studies have investigated narratives in social media content [14, 25, 34]. These Narratives contain small autobiographies that have been developed in personal profiles and cover trivial everyday life events. Other types of narratives appearing in social media platforms consist of breaking news and long stories of past events [25]. Some types of narratives, such as breaking news, result in the emergence of other narratives related to the predictions or projections of events in near future [14]. These literature view social media conversation cascades as stories that are co-constructed by the tellers and their audience, and are circulating amongst the public within and across social media platforms.

Moreover, the events have been considered as the causes of online user activity that can be identified via activity fluctuations over time [3, 25]. Developing appropriate tools for social media narrative analysis can facilitate communicating the main ideas regarding the events in large data.

2.2 Topic Mining and Text Summarization

As social media activities generate abundant timestamped multi-modal data, many studies such as [8] have presented algorithms to discover the topics and develop descriptive summaries over social media events. probabilistic models to discover word patterns that reflect the underlying topics in a set of document collections [1]. The most commonly used approach to topic modeling is Latent Dirichlet Allocation (LDA) [19]. LDA is a generative probabilistic model with a hierarchical Bayesian network structure that can be used for a variety of applications with discrete data, including text corpora. Using LDA for topic mining, a document is a bag-of-words that has a mixture of latent topics [6]. Many advanced topic modeling approaches have been derived from LDA, including Hierarchical Topic Models [15, 16] that learn and organize the topics into a hierarchy to address a super-sub topic relationship. This approach is well-suited for analyzing social media and news stories that contain rich data over a series of real-world events [30]. Topic models over time with continuous-time distribution [5] and dynamic topic models [35] intend to capture the rise and falls of topics within a time range. However, continuous-time topic models, such as beta or normal time distribution, cannot model many modes in time. Furthermore, the smooth time distribution over topics does not allow recognizing distinct topical events in the timestamped dataset, where topical events reflect the event-based topic activity fluctuations over time.

Topic modeling and summarization of social media data is challenging as a result of certain restrictions, such as the maximum number of characters allowed on the Twitter platform. As short-text or microblogs have low word co-occurrence and contextual information, models designed for short-text topic analysis and summarization may obtain context information with short-text aggregation to enrich the relevant context before further analysis [27].

Document summarization techniques are generally categorized into abstractive and generative text summarization models. Herein, we consider extractive text summarization methods. Several algorithms for extractive text summarization have been proposed in the literature that assign a salient score to sentences [10]. To summarize a text corpus with short text, [29] presents an automatic summarization algorithm with topic clustering, cluster ranking and assigning scores to the intermediate features, and sentence extraction. Some other approaches, particularly for the Twitter data include aggregating tweets by hashtags or conversation cascades [27, 32], and obtaining summaries for a targeted event of interest as one or a set of tweets that are representative of the topics [8].

Additionally, neural network-based summarization models [23, 28], commonly with an encoder-decoder architecture, leverage attention mechanism for contextual information among sentences or ROUGE evaluation metric to identify discriminative features for sentence ranking and summarization. However, these architectures require labeled datasets and might not apply to short-text.

Text summarization with compression using neural networks is proposed by [36] which applies joint extraction and syntactic compression to rank compressed summaries with a neural network. Our focus in the present work is on probabilistic topic modeling and extractive text summarization to provide descriptive narratives for the underlying events that occurred over a period of time.

3 METHODOLOGY

In this section we explain our narrative framework. The framework comprises of 2 steps: I. Narrative modeling based on topic identification over time and II. extractive summarization from the identified narratives. To discover the narratives over topical events, first, we use our discrete-time generative narrative model as an unsupervised learning algorithm to learn distribution of textual contents from daily conversation cascades. Then, we extract narrative summaries over topical events from sentences in the time categories. This is achieved by sampling from the identified distribution of narratives and perform sentence ranking. Narrative modeling and summarization steps are explained below in separate subsections.

3.1 Narrative Modeling

To model narratives, we design our topic model such that the discovered topics present a series of timely ordered topical events. Accordingly, the topical events deliver a narrative covering distinct events over the same topic. In this regard, we present Narratives Over Categorical time (NOC), a novel probabilistic topic model that discovers topics based on both word co-occurrence and temporal information to present a narrative of events. According to the topic-time relationship explained above, we refer to the topics or narratives, topical events as events, and the extracted timely ordered sentences of documents with high probability of belonging to each event as the extracted narrative summary. To fully comply with the definition of narrative, we assume a causality relation between the conversation cascades in social media. However, we do not investigate the causality relation across the conversation cascades or named-entities.

The differences between our Narrative model with dynamic topic models [5], topic models with continuous time distribution [35], and hierarchical topic models [16, 26] include: not filtering the data for an specific event, imposing sharp transition for topic-time changes with time slicing, discovering topical events without scalability and sparsity issues, allowing multimodal topic distribution in time as a result of categorical time distribution, and selecting an appropriate slicing size such that distinct topical events be recognizable. Additionally, categorical time distribution enables discovering topical events with varying time resolution, for instance, weekly, biweekly, and monthly.

Time discretization brings the question of selecting the appropriate slicing size or the number of categories that depends on the characteristics of the dataset under study. On the contrary, topical models with continuous time distribution cannot model many modes in time. Additionally, continuous time models such as [35] suffer from instability problem if the dataset is multimodal and sparse in time. Furthermore, categorical time enables discovering topic recurrence which results in identifying topical events related to distinct narrative activities, which is of our interest in

Table 1: Symbols and definitions used in this paper

Variable Descriptions	Symbol
Number of topics	T
Number of documents	D
Number of unique words	V
Number of word tokens in document d	N_d
Multinomial distribution of topics for document d	θ_d
Multinomial distribution of words for topic z	ϕ_z
Categorical distribution of time for topic z	ψ_z
Topic of the i th token in document d	z_{di}
i th word token in document d	w_{di}
Timestamp for i th word token in document d	t_{di}
Time category for timestamp associated with a token	b_k
Entropy of topic z	H_z
j th sentence of document d	s_{dj}

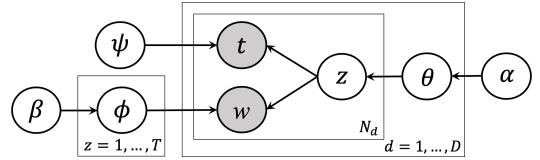


Figure 1: The graphical model for NOC with Gibbs sampling.

this paper. Narrative activities in social media refer to the amount of textual content that is circulating in online platforms over time, corresponding to a specific topic.

The generative process in NOC, models timestamps and words per documents using Gibbs sampling which is a Markov Chain Monte Carlo (MCMC) algorithm. The graphical model of NOC is illustrated in Figure 1. As can be seen from the graphical model, the posterior distribution of topics is dependent on both text and time modalities. This generative procedure can be described as follows:

- I. For each topic z , draw T multinomials ϕ_z from a Dirichlet prior β ;
- II. For each document d , draw a multinomial θ_d from a Dirichlet prior α ;
- III. For each word w_{di} in d :
 - (a) draw a topic z_{di} from multinomial θ_d ;
 - (b) draw a word w_{di} from multinomial $\phi_{z_{di}}$;
 - (c) draw a timestamp t_{di} from categorical $\psi_{z_{di}}$;

where the timestamps t_{di} for words w_{di} in each document d are identical. The list of symbols and their descriptions can be found in table 1. The model parameterization is as below:

$$\begin{aligned}
 \theta_d | \alpha &\sim \text{Discrete}(\alpha) \\
 \phi_z | \beta &\sim \text{Discrete}(\beta) \\
 z_{di} | \theta_d &\sim \text{Multinomial}(\theta_d) \\
 w_{di} | \phi_{z_{di}} &\sim \text{Multinomial}(\phi_{z_{di}}) \\
 t_{di} | \psi_{z_{di}} &\sim \text{Categorical}(\psi_{z_{di}})
 \end{aligned} \tag{1}$$

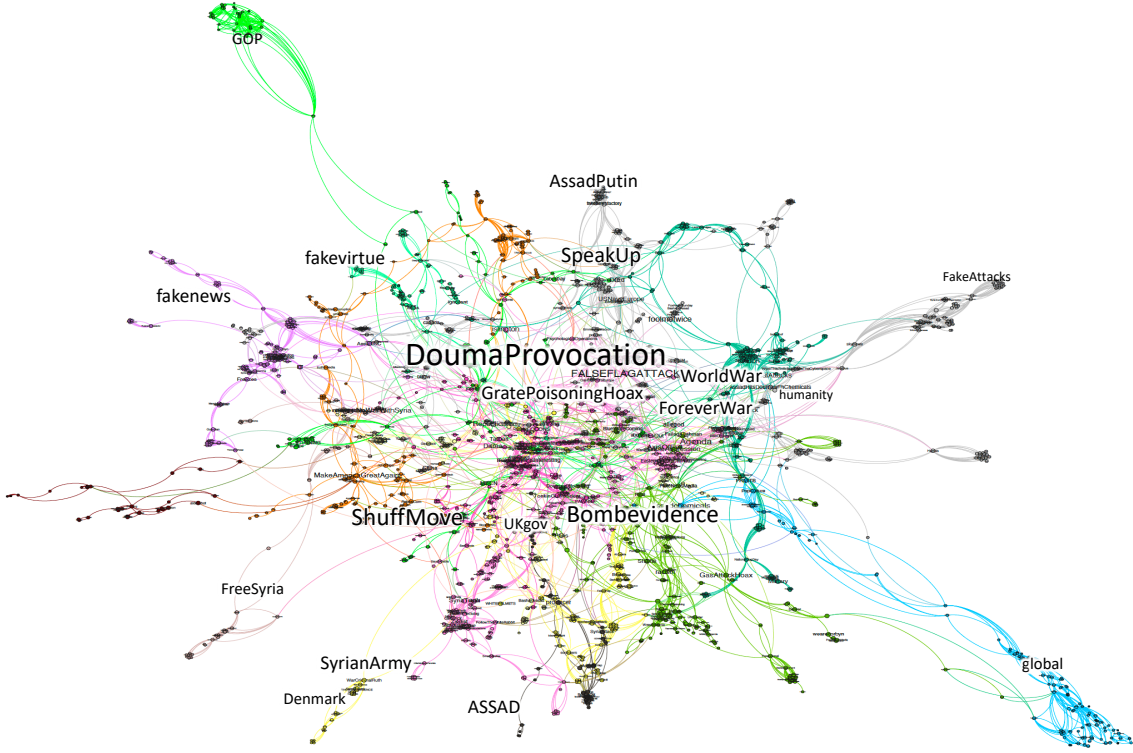


Figure 2: The hashtag co-occurrence graph for twitter dataset on the domain of White Helmets of Syria for a period of 13 month, from 2018 to 2019. This graph represents a down-sampled version of the hashtag co-occurrence in this data for the sake of visualization.

In this model, Gibbs sampling provides an approximate inference instead of exact inference. To calculate the probability of topic assignment to word w_{di} , we first need to calculate the joint probability of the dataset as $\mathbb{P}(z_{di}, w_{di}, t_{di} | w_{-di}, t_{-di}, z_{-di}, \alpha, \beta, \psi)$ and use chain rule to derive the probability of $\mathbb{P}(z_{di} | w, t, z_{-di}, \alpha, \beta, \psi)$ as below, where $-di$ subscript refers to all tokens except w_{di} :

$$\begin{aligned} \mathbb{P}(z_{di} | w, t, z_{-di}, \alpha, \beta, \psi) &\propto (m_{dz_{di}} + \alpha_{z_{di}} - 1) \\ &\times \frac{n_{z_{di}w_{di}} + \beta_{w_{di}} - 1}{\sum_{v=1}^V (n_{z_{di}v} + \beta_v) - 1} p(t_{z_{di}} \in b_k) \end{aligned} \quad (2)$$

where n_{zv} refers to the number of words v assigned to topic z , m_{dz} refers to the number of word tokens in document d that are assigned to topic z , and b_k represents the k th time slice. The details on the Gibbs sampling derivation can be found in the Appendix section. After each iteration of Gibbs sampling, we update the probability of $p(t_{z_{di}} \in b_k)$ as follows:

$$p(t_{z_{di}} \in b_k) = \frac{1}{K} \sum_{k=1}^K \mathbb{I}(t_{z_{di}} \in b_k) \quad (3)$$

where $\mathbb{I}(\cdot)$ is equal to 1 when $t_{z_{di}} \in b_k$, and 0 otherwise.

In this paper, we report results with bi-weekly categorical time resolution. To determine the values for hyper-parameters α and β and to investigate the sensitivity of the model to these values, we repeated our experiment with symmetric Dirichlet distributions using values $\alpha \in [0.1, 0.5, 1]$, $\beta \in [0.01, 0.1, 0.5, 0.8, 1]$. We observed

that the model did not show significant sensitivity to the values of these hyper-parameters. Thus, we fix $\alpha = 1$ and $\beta = 0.5$, both as symmetric Dirichlet distributions. We initialize the hyperparameter ψ in 2 ways for comparison: I. random initialization (model referred as NOC_R); and II. based on the probability of user activity per time category, illustrated in Figure 3c, (model referred as NOC_A).

To estimate the number of topics for our experiments, we first visualize the tweets' hashtag co-occurrence graph. We measure the graph modularity to examine the structure of the communities in this graph. We observe the highest modularity score of 0.41 using modularity resolution equal to 0.85. Figure 2 illustrates a downsample version of this graph, where each color represents a modularity class. The edges of the graph are weighted according to the number of hashtags' co-occurrence in the document collection. Our modularity analysis suggests that few distinct hashtag communities exist. Additionally, the dataset under study contains tweets associated with a single domain. As a result, we assume the number of topics to be relatively low. To choose an appropriate number of topics, we repeated LDA with the number of topics as $T \in [4, \dots, 20]$ with increments of size 1. We evaluated the c_v coherence of topics identified by LDA and observed the highest coherence score for $T = 5$ and $T = 5$, respectively. Thus, we report our experimental results using these values.

3.2 Narrative Summary Extraction

We employ the discovered probabilities of topics over documents, θ , probabilities of words per topic, ϕ , and probabilities of topics per time category, ψ to perform sentence ranking. This ranking allows extracting the sentences with the higher scores of belonging to each topic. This is achieved via performing weighted sampling on the collection of documents based on the probabilities of topics per time category ψ and draw D documents from θ . The weighted sampling leads to drawing more documents from the time categories b_k with a higher ψ as this time slices contain more documents related to the topic z . Each document contains a sequence of sentences $(s_1, s_2, \dots, s_j) \in d$ from the aggregated conversation cascades per day. Information on the aggregation of conversation cascades and document preparation can be found in section 4.2.

Since the social media narrative activity over a topic evolves from the circulation of identical or similar textual content in the platform, the content involves significant similarity. For instance, the Twitter conversation cascades include replies, quotes, and comments, where replies and quotes duplicate the textual content. Therefore, we applied Jaro-Winkler distance over the timely ordered sentences and dismissed the sentences with similarity above 70%, while keeping the longest sentence. After removing redundant text as described earlier, we calculate the probability of each sentence s_j by measuring the sum of the probabilities of topics for words $w_{di} \in s_j$. Then, we select the sentences with the highest accumulative probability of words w per topic z . Summary coherence was induced as suggested in [4] by ordering the extracted sentences according to their timestamps such that the oldest sentences appear first. Table 4 in the Appendix section contains the extracted narrative summaries for 5 topics for a sample run.

4 EXPERIMENTS AND RESULTS

As mentioned earlier, the discovered topics by NOC present a series of timely ordered topical events. Thus, the topical events deliver a narrative covering distinct social media events over the same topic. Figure 3 demonstrates the generated narrative distributions with NOC, where the hyperparameter ψ was randomly initialized (referred to as NOC_R). This figure represents that the identified narratives by our model are distinct from each other and the collapsed distribution of all narratives approximates the distribution of social media user activity over time.

The identified narratives can be evaluated using effective evaluation metrics for topic models. Accordingly, we calculate pointwise mutual information [24] to measure the coherence of a topic z as follows:

$$Coh_z = \frac{2}{K(K-1)} \sum_{j < k \leq K} \log \frac{p(w_j, w_k)}{p(w_j)p(w_k)} \quad (4)$$

where K is the number of most probable words for each narrative, $p(w_j)$ and $p(w_k)$ refer to the probabilities of occurrence for words w_j and w_k , and $p(w_j, w_k)$ represents the probability of co-occurrence for the two words in the collection of documents.

We compare our model with LDA and TOT[35], where TOT is a probabilistic topic model over time with Beta distribution for time. Table 2 displays the average coherence score measured across the discovered topics by LDA, TOT, and NOC. For NOC, we investigate

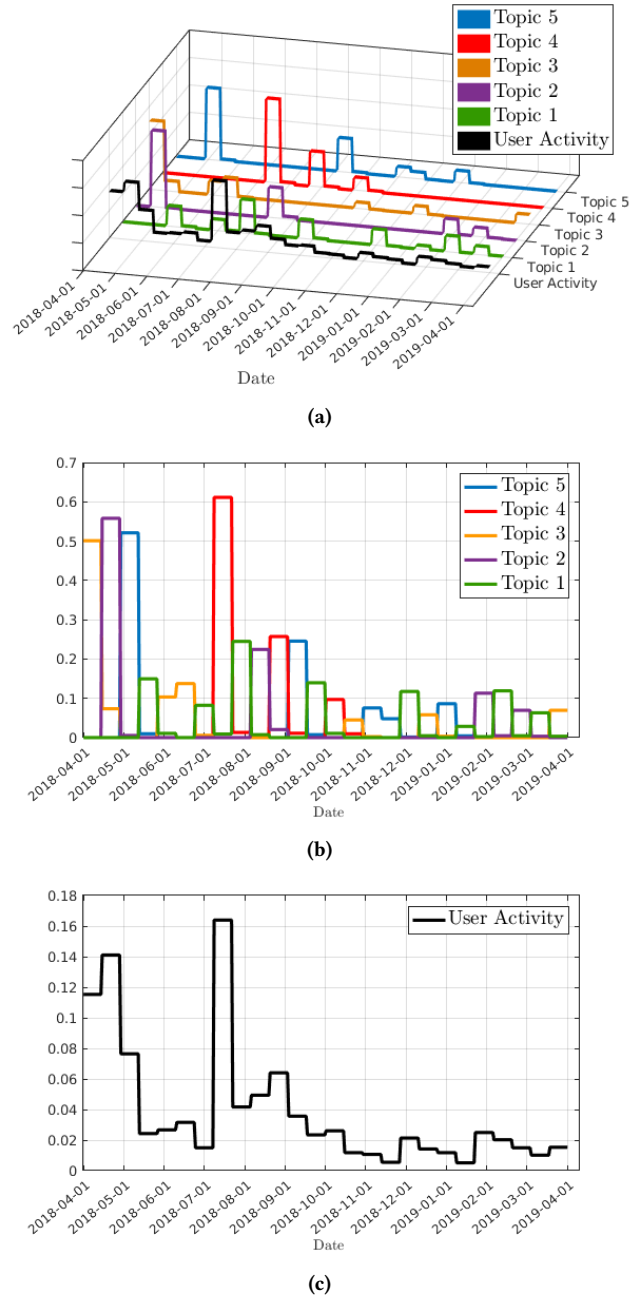


Figure 3: The distribution of extracted topics and user activity over time: a) The distribution of user activity over time is depicted by black, followed by the 5 distribution of the extracted topics, using NOC_R ; b) The collapsed distribution of the 5 extracted topics; c) The distribution of user activity in time. The results suggest that the distribution of extracted topics approximate the distribution of user activity over the same time period.

Table 2: The comparison of coherence scores for 4 models:

<i>Model</i>	<i>LDA</i>	<i>TOT</i> ³	<i>NOC_R</i>	<i>NOC_A</i>
<i>T = 5</i>	5.980	6.36	7.95	8.23
<i>T = 6</i>	5.546	5.99	7.75	7.98

initializing the parameter ψ with random and user activity-based initialization, referred as NOC_R and NOC_A , respectively. We consider $K = 500$ most probable words from each topic. This comparison suggests that the narratives identified by NOC are more coherent than the identified topics by LDA , with an improvement in coherence of about 35%. The observed improvement comparing with TOT was about 27%. Additionally, initializing the hyperparameter ψ in NOC using the distribution of user activity improves the narrative coherence by about 3%.

4.1 Proposed Evaluation Metric

The topic attractiveness to social media users can be investigated as a measure of the length of conversation cascades, the number of initiated textual content, and the number of unique users performing an activity relative to the underlying topic. The user activity fluctuations for timestamped data may contain activity bursts that are illustrative of significant events. Similarly, the generation and propagation of textual content within an online platform can illustrate the narrative activity relative to the events over time, where a burst represents a significant narrative activity. Additionally, the recurrence of a topic can be considered as an attractiveness measure for the associated topic.

In this regard, we propose the significance-dispersity trade-off (SDT) metric to compare the identified narratives against each-other. SDT measures the lifetime attractiveness of the identified narratives based on the distribution of narratives over topical events. The proposed metric quantifies the significance of the narrative activities and recurrence of a topic via employing the Shannon entropy for the discovered narrative distributions. The intuition behind the SDT score is that the value of the entropy is maximum when the probability distribution is uniform. On the contrary, this value is minimum if the distribution is delta function. This is visualized in Figure 4 in the Appendix section.

We define dispersity of a categorical time topic distribution as a measure of the dispersion of the time categories. Based on this definition, SDT score of topic z can be obtained as:

$$SDT_z = H^\gamma (H_{max} - H)^{1-\gamma}, \quad (5)$$

where H is the Shannon entropy for the categorical distribution of time for topic z :

$$H_z = - \sum_k^K p_z \log_2 p_z, \quad (6)$$

$$H_{max} = \log_2(K),$$

and K refers to the number of time slices in the distribution. We assume that social media topics with high lifetime attractiveness

³We used the implementation available on https://github.com/ahmaurya/topics_over_time.

Table 3: The comparison of SDT scores for 5 narratives:

<i>Narrative</i>	<i>T₁</i>	<i>T₂</i>	<i>T₃</i>	<i>T₄</i>	<i>T₅</i>
$\gamma = 0$	1.59	2.90	2.39	3.21	2.75
$\gamma = 0.4$	2.08	2.40	2.36	2.36	2.40
$\gamma = 0.7$	2.54	2.08	2.33	1.87	2.16
$\gamma = 1$	3.11	1.80	2.31	1.49	1.95
User Activity	353,280	317,686	244,674	247,895	175,343

are significant and recurrent. However the probability distribution imposes a trade-off on the two. The parameter α provides a weighted geometric mean of H and $H_{max} - H$ that enables promoting either significance or recurrence, dependent on the application under study. A larger value of parameter α promotes dispersity for SDT score, and a smaller amount of this parameter promotes mode significant. The bounds for the SDT score are:

$$SDT_i = \begin{cases} 0 & \text{if } H = 0 \text{ \& } \gamma! = 0 \\ 0 & \text{if } H = H_{max} \text{ \& } \gamma! = 1 \\ \gamma^\gamma (1 - \gamma)^{1-\gamma} H_{max} & \text{if } H = \gamma H_{max} \text{ \& } 0 < \gamma < 1 \end{cases} \quad (7)$$

where $H = 0$ occurs when the distribution under study is uniform, and $H = H_{max}$ relates to delta distribution. Since the time categorical distribution of our narrative model allows many modes in time, recurrent narratives can be identified. Additionally, the narrative activity fluctuations can be modeled using categorical time distribution in topic analysis. Table 3 provides a comparison for the SDT scores measured for the 5 identified narratives, using varying values of α . The illustration of the distribution of the extracted narratives can be seen in Figure 3a. We can clearly see in this figure that narratives 1 and 3 have the highest dispersity. On the contrary, narratives 4 and 2 have the highest significance. We compare SDT_i for narrative i with the number of user activity associated with narrative z . The results suggest that SDT score can be used to identify the narrative with higher lifetime attractiveness in a timestamped dataset. In our experiments, this is achieved for topic 1 when the value of γ is greater than or equal to 0.7. As it can be seen, this topic is associated with the highest user activity count, reported in the same table.

4.2 Dataset Description and Pre-processing

To analyze topical events and provide narratives, we investigate the Twitter dataset on the domain of White Helmets of Syria over a period of 13 month from April 2018 to April 2019. This dataset was provided to us by Leidos Inc¹ as part of the Computational Simulation of Online Social Behavior (SocialSim)² program initiated by the Defense Advanced Research Projects Agency (DARPA). We analyze more than 1,052,000 tweets from April 2018 to April 2019.

To prepare the model inputs, we filter the tweets from non-English text. Then, we clean up the data by removing usernames, short URLs, as well as emoticons. Additionally, we remove the stopwords, perform Part of Speech (POS) tagging and Named

¹<https://www.leidos.com/>

²<https://www.darpa.mil/program/computational-simulation-of-online-social-behavior>

Entity Recognition (NER) on each tweet using Stanford Named Entity Recognizer³ model. Using the NER tool, we extract persons, locations and organizations and removed all pseudo-documents that do not contain named entities similar to [21]. Furthermore, We remove the tweets shorter than 3 words.

As the Twitter maintains a maximum allowed character limit of 280 characters, collected tweets lack context information and have very low word co-occurrence. We tackle the challenge of topic modeling on short-text tweets and to include plentiful context information by preparing pseudo-documents for our model inputs via aggregating daily root, parent, and reply/quote/retweet comments in each conversation cascade. We maintain the order of the conversation according to the timestamps associated with each tweet. This text aggregation method results in preparing pseudo-documents rich of context and related words with a daily time resolution. We use the pre-processing phase output as the model input pseudo-documents, referred as documents in this paper.

5 CONCLUSION AND FUTURE DIRECTIONS

In this paper, we addressed the problem of narrative modeling and narrative summary extraction for social media content. We presented a narrative framework consisting of I. Narratives over topic Categories (NOC), a probabilistic topic model with categorical time distribution; and II. extractive text summarization. The proposed narrative framework identifies narrative activities associated with social media events. Identifying topics' recurrence and significance over time categories with our model allowed us to propose significance-dispersity trade-off (SDT) metric. SDT can be employed as a comparison measure to identify the topic with the highest lifetime attractiveness in a timestamped corpus. Results on real-world timestamped data suggest that the narrative framework is effective in identifying distinct and coherent topics from the data. Additionally, the results illustrate that the identified narrative distributions approximate the user activity fluctuations over time. moreover, informative, and concise narrative summaries for timestamped data are produced. Further improvement of the narrative framework can be achieved via incorporating the causality relation cross the social media conversation cascades and social media events into account. Other future directions include identifying topical hierarchies and extract summaries associated with each hierarchy.

ACKNOWLEDGMENTS

This work was supported by the Defense Advanced Research Projects Agency (DARPA) under grant number FA8650-18-C-7823. The views and opinions expressed in this article are the authors' own and should not be construed as official or as reflecting the views of the University of Central Florida, DARPA, or the U.S. Department of Defense.

REFERENCES

[1] Rubayyi Alghamdi and Khalid Alfalqi. 2015. A survey of topic modeling in text mining. *Int. J. Adv. Comput. Sci. Appl.(IJACSA)* 6, 1 (2015).
 [2] Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D Trippe, Juan B Gutierrez, and Krys Kochut. 2017. Text summarization techniques: a brief survey. *arXiv preprint arXiv:1707.02268* (2017).

[3] Jeffery Ansah, Lin Liu, Wei Kang, Jixue Liu, and Jiuyong Li. 2020. Leveraging burst in twitter network communities for event detection. *World Wide Web* (2020), 1–26.
 [4] Regina Barzilay, Noemie Elhadad, and Kathleen R McKeown. 2001. Sentence ordering in multidocument summarization. In *Proceedings of the first international conference on Human language technology research*. Association for Computational Linguistics, 1–7.
 [5] David M Blei and John D Lafferty. 2006. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*. 113–120.
 [6] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research* 3, Jan (2003), 993–1022.
 [7] Sophie Burkhardt and Stefan Kramer. 2019. A Survey of Multi-Label Topic Models. *ACM SIGKDD Explorations Newsletter* 21, 2 (2019), 61–79.
 [8] Freddy Chong Tat Chua and Sitaram Asur. 2013. Automatic summarization of events from social media. In *Seventh international AAAI conference on weblogs and social media*.
 [9] Matteo Cinelli, Walter Quattrociocchi, Alessandro Galeazzi, Carlo Michele Valensise, Emanuele Brugnoli, Ana Lucia Schmid, Paola Zola, Fabiana Zollo, and Antonio Scala. 2020. The covid-19 social media infodemic. *arXiv preprint arXiv:2003.05004* (2020).
 [10] Soumi Dutta, Vibhash Chandra, Kanav Mehra, Sujata Ghatak, Asit Kumar Das, and Saptarshi Ghosh. 2019. Summarizing microblogs during emergency events: A comparison of extractive summarization algorithms. In *Emerging Technologies in Data Mining and Information Security*. Springer, 859–872.
 [11] Gunn Enli. 2017. Twitter as arena for the authentic outsider: exploring the social media campaigns of Trump and Clinton in the 2016 US presidential election. *European journal of communication* 32, 1 (2017), 50–61.
 [12] Samah Gad, Waqas Javed, Sohaib Ghani, Niklas Elmqvist, Tom Ewing, Keith N Hampton, and Naren Ramakrishnan. 2015. ThemeDelta: Dynamic segmentations over temporal topic models. *IEEE transactions on visualization and computer graphics* 21, 5 (2015), 672–685.
 [13] Ivan Garibay, Alexander V Mantzaris, Amirarsalan Rajabi, and Cameron E Taylor. 2019. Polarization in social media assists influencers to become more influential: analysis and two inoculation strategies. *Scientific Reports* 9, 1 (2019), 1–9.
 [14] Alexandra Georgakopoulou. 2017. 17 Small Stories Research: A Narrative Paradigm for the Analysis of Social Media. *The Sage Handbook of social media research methods* (2017).
 [15] Goran Glavaš, Jan Šnajder, Parisa Kordjamshidi, and Marie-Francine Moens. 2014. HiEve: A corpus for extracting event hierarchies from news stories. In *Proceedings of 9th language resources and evaluation conference*. ELRA, 3678–3683.
 [16] Thomas L Griffiths, Michael I Jordan, Joshua B Tenenbaum, and David M Blei. 2004. Hierarchical topic models and the nested chinese restaurant process. In *Advances in neural information processing systems*. 17–24.
 [17] Naeemul Hassan, Gensheng Zhang, Fatma Arslan, Josue Caraballo, Damian Jimenez, Siddhant Gawsane, Shohedul Hasan, Minumol Joseph, Aaditya Kulkarni, Anil Kumar Nayak, et al. 2017. ClaimBuster: the first-ever end-to-end fact-checking system. *Proceedings of the VLDB Endowment* 10, 12 (2017), 1945–1948.
 [18] Bram Jans, Steven Bethard, Ivan Vulić, and Marie-Francine Moens. 2012. Skip n-grams and ranking functions for predicting script events. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 336–344.
 [19] Hamed Jelodar, Yongli Wang, Chi Yuan, Xia Feng, Xiahui Jiang, Yanchao Li, and Liang Zhao. 2019. Latent Dirichlet Allocation (LDA) and Topic modeling: models, applications, a survey. *Multimedia Tools and Applications* 78, 11 (2019), 15169–15211.
 [20] Amrah Maryam and Rashid Ali. 2018. Twitter Based Event Summarization. In *International Conference on Advances in Computing and Data Sciences*. Springer, 165–174.
 [21] Andrew J McMinin and Joemon M Jose. 2015. Real-time entity-based event detection for twitter. In *International conference of the cross-language evaluation forum for european languages*. Springer, 65–77.
 [22] Elliot G Mishler. 1995. Models of narrative analysis: A typology. *Journal of narrative and life history* 5, 2 (1995), 87–123.
 [23] Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Ranking sentences for extractive summarization with reinforcement learning. *arXiv preprint arXiv:1802.08636* (2018).
 [24] David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. 2010. Automatic evaluation of topic coherence. In *Human language technologies: The 2010 annual conference of the North American chapter of the association for computational linguistics*. Association for Computational Linguistics, 100–108.
 [25] Ruth Page. 2013. Seriality and storytelling in social media. *StoryWorlds: A Journal of Narrative Studies* 5 (2013), 31–54.
 [26] Jay Pujara and Peter Skomoroch. 2012. Large-scale hierarchical topic models. In *NIPS Workshop on Big Learning*, Vol. 128.
 [27] Xiaojun Quan, Chunyu Kit, Yong Ge, and Sinno Jialin Pan. 2015. Short and sparse text topic modeling via self-aggregation. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.

³<https://nlp.stanford.edu/software/CRF-NER.html>

- [28] Pengjie Ren, Zhumin Chen, Zhaochun Ren, Furu Wei, Jun Ma, and Maarten de Rijke. 2017. Leveraging contextual sentence relations for extractive summarization using a neural attention model. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 95–104.
- [29] Lidan Shou, Zhenhua Wang, Ke Chen, and Gang Chen. 2013. Sumblr: continuous summarization of evolving tweet streams. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*. 533–542.
- [30] PK Srijith, Mark Hepple, Kalina Bontcheva, and Daniel Preotiu-Pietro. 2017. Sub-story detection in Twitter with hierarchical Dirichlet processes. *Information Processing & Management* 53, 4 (2017), 989–1003.
- [31] Jiwei Tan, Xiaojun Wan, and Jianguo Xiao. 2017. From Neural Sentence Summarization to Headline Generation: A Coarse-to-Fine Approach. In *IJCAL*. 4109–4115.
- [32] Johnny Torres, Carmen Vaca, Luis Terán, and Cristina L. Abad. 2020. Seq2Seq models for recommending short text conversations. *Expert Systems with Applications* 150 (2020), 113270.
- [33] Josep Valls Vargas. 2017. *Narrative information extraction with non-linear natural language processing pipelines*. Drexel University.
- [34] Kristin Veel. 2018. Make data sing: The automation of storytelling. *Big Data & Society* 5, 1 (2018), 2053951718756686.
- [35] Xuerui Wang and Andrew McCallum. 2006. Topics over time: a non-Markov continuous-time model of topical trends. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. 424–433.
- [36] Jiacheng Xu and Greg Durrett. 2019. Neural extractive text summarization with syntactic compression. *arXiv preprint arXiv:1902.00863* (2019).
- [37] Kang Xu, Guilin Qi, Junheng Huang, and Tianxing Wu. 2017. Incorporating Wikipedia concepts and categories as prior knowledge into topic models. *Intelligent Data Analysis* 21, 2 (2017), 443–461.
- [38] Liang Yao, Yin Zhang, Baogang Wei, Lei Li, Fei Wu, Peng Zhang, and Yali Bian. 2016. Concept over time: the combination of probabilistic topic model with wikipedia knowledge. *Expert Systems with Applications* 60 (2016), 27–38.

APPENDIX

A GIBBS SAMPLING DERIVATION FOR THE DISCRETE-TIME NARRATIVE MODEL

Starting with the joint distribution $\mathbb{P}(\mathbf{w}, \mathbf{t}, \mathbf{z} | \alpha, \beta, \psi)$, we can use conjugate priors to simplify the equations as below:

$$\begin{aligned}
\mathbb{P}(\mathbf{w}, \mathbf{t}, \mathbf{z} | \alpha, \beta, \psi) &= \mathbb{P}(\mathbf{w} | \mathbf{z}, \beta) p(\mathbf{t} | \psi, \mathbf{z}) \mathbb{P}(\mathbf{z} | \alpha) \\
&= \int \prod_{d=1}^D \prod_{i=1}^{N_d} \mathbb{P}(w_{di} | \phi_{z_{di}}) \prod_{z=1}^T p(\phi_z | \beta) d\Phi \prod_{d=1}^D \prod_{i=1}^{N_d} p(t_{di} | \psi_{z_{di}}) \\
&\times \int \prod_{d=1}^D \left(\prod_{i=1}^{N_d} \mathbb{P}(z_{di} | \theta_d) p(\theta_d | \alpha) \right) d\Theta \\
&= \int \prod_{z=1}^T \prod_{v=1}^V \phi_{zv}^{n_{zv}} \prod_{z=1}^T \left(\frac{\Gamma(\sum_{v=1}^V \beta_v)}{\prod_{v=1}^V \Gamma(\beta_v)} \prod_{v=1}^V \phi_{zv}^{\beta_v - 1} \right) d\Phi \\
&\times \int \prod_{d=1}^D \prod_{z=1}^T \theta_{dz}^{m_{dz}} \prod_{d=1}^D \left(\frac{\Gamma(\sum_{z=1}^T \alpha_z)}{\prod_{z=1}^T \Gamma(\alpha_z)} \prod_{z=1}^T \theta_{dz}^{\alpha_z - 1} \right) d\Theta \quad (8) \\
&\times \prod_{d=1}^D \prod_{i=1}^{N_d} p(t_{di} | \psi_{z_{di}}) \\
&= \left(\frac{\Gamma(\sum_{v=1}^V \beta_v)}{\prod_{v=1}^V \Gamma(\beta_v)} \right)^T \left(\frac{\Gamma(\sum_{z=1}^T \alpha_z)}{\prod_{z=1}^T \Gamma(\alpha_z)} \right)^D \prod_{d=1}^D \prod_{i=1}^{N_d} p(t_{di} | \psi_{z_{di}}) \\
&\times \prod_{z=1}^T \frac{\prod_{v=1}^V \Gamma(n_{zv} + \beta_v)}{\gamma(\sum_{v=1}^V (n_{zv} + \beta_v))} \prod_{d=1}^D \frac{\prod_{z=1}^T \Gamma(m_{dz} + \alpha_z)}{\gamma(\sum_{z=1}^T (m_{dz} + \alpha_z))},
\end{aligned}$$

where \mathbb{P} and p refer to the probability mass function (PMF) and probability density function (PDF), respectively. The conditional probability $\mathbb{P}(z_{di} | \mathbf{w}, \mathbf{t}, z_{-di}, \alpha, \beta, \psi)$ can be found using the chain rule as:

$$\begin{aligned}
\mathbb{P}(z_{di} | \mathbf{w}, \mathbf{t}, z_{-di}, \alpha, \beta, \psi) &= \frac{\mathbb{P}(z_{di}, w_{di}, t_{di} | w_{-di}, t_{-di}, z_{-di}, \alpha, \beta, \psi)}{\mathbb{P}(w_{di}, t_{di} | w_{-di}, t_{-di}, z_{-di}, \alpha, \beta, \psi)} \\
&\propto \frac{\mathbb{P}(\mathbf{w}, \mathbf{t}, \mathbf{z} | \alpha, \beta, \psi)}{\mathbb{P}(w_{-di}, t_{-di}, z_{-di} | \alpha, \beta, \psi)} \\
&\propto \frac{n_{z_{di} w_{di}} + \beta_{w_{di}} - 1}{\sum_{v=1}^V (n_{z_{di} v} + \beta_v) - 1} (m_{dz_{di}} + \alpha_{z_{di}} - 1) p(t_{di} | \psi_{z_{di}}) \\
&\propto (m_{dz_{di}} + \alpha_{z_{di}} - 1) \frac{n_{z_{di} w_{di}} + \beta_{w_{di}} - 1}{\sum_{v=1}^V (n_{z_{di} v} + \beta_v) - 1} p(t_{z_{di}} \in b_k) \quad (9)
\end{aligned}$$

The probability of $p(t_{di} \in b_k)$ can be measured as follows:

$$p(t_{z_{di}} \in b_k) = \frac{1}{K} \sum_{k=1}^K \mathbb{I}(t_{z_{di}} \in b_k), \quad (10)$$

where $\mathbb{I}(\cdot)$ is equal to 1 when $t_{z_{di}} \in b_k$, and 0 otherwise.

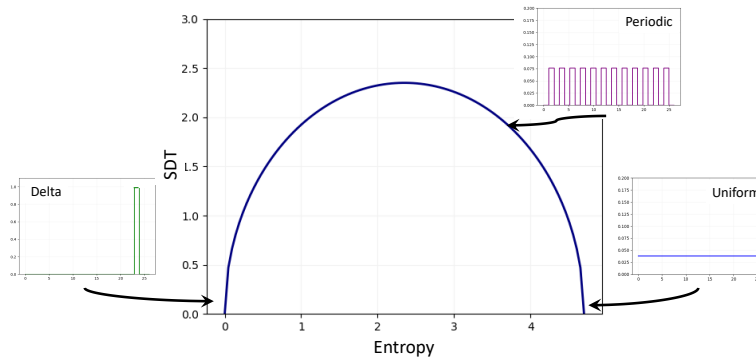


Figure 4: A visualization of SDT for different entropy values using $\gamma = 0.5$. The SDT values for delta, uniform, and periodic distributions are marked on the graph.

Table 4: Representative keywords and narrative summaries per topic.

Topics	Keywords	Summary
Topic 1	Terrorist, Idlib, Civilian, Child, City, Attack, Aleppo, Rescue, Weapon, Killed	WhiteHelmets Syria News: One child was injured in the north of Aleppo. Their aim is to save lives in war zones inside Syria. Has credibly substantiated 336 uses of ChemWeapons in Syria 98% of attacks by Assadallies. These are the WHITE HELMETS or Syria Civil Defense as our US Dept of State calls them!! Russian airstrikes killed two men and one baby in DMZ areas RussianWarCrimes.
Topic 2	Chemical, Attack, Douma, Video, Idlib, Staged, Boy, War, Child, Witness	Remember first they said the video including the pics of the chlorine cylinder was fake. Whitehelmets One America News Pearson Sharp Visits Hospital in Douma Where White Helmets Filmed Chemical Attack Hoax Multiple Eyewitness Doctors Say No Chemical Attack Took Place Syria. This is the video evidence of the airstrike on Zardana an Idlib town controlled by Very expensive camera on the helmet of the WhiteHelmets rescuer. White Helmets making films of chemical attacks with children in Idlib.
Topic 3	Chemical, Attack, Douma, Terrorist, Fake, Child, Propaganda, Video, Russian, Russia	From the fabrication of the plays of the chemist and coverage of the crimes of terrorism to the public cooperation with the Israeli army the white helmets. They are holding children! Another chemical attack is imminent its all they've got left! 4 dead including two children and more than 50 wounded mostly women and children. Love the White Helmets propaganda almost as untruthful as the BBC.
Topic 4	Israel, Terrorist, Idlib, Chemical, Attack, Life, Rescue, Russian, People, Al Qaeda	WHITE HELMETS ARE PREPARING CHEMICAL ATTACK ON CITIZENS AGAIN! Those are basically just members of Al Qaeda Al Nusra right? The Al Qaeda smear is deliberate propaganda. Its war crime only If US intervenes in Kashmir Kashmir will be liberated like Raqqan with a dozen US bases having Thaad missiles aimed at China and with AlQaeda WhiteHelmets taking out children's organs of Kashmiris.
Topic 5	Funding, Freeze, Trump, Terrorist, Group, Chemical, Attack, Idlib, Civilian, News	Trumps USA has built a rationale for its public that it will need to support rebels in holding on to a large chunk of Syria. I wonder how it is possible that criminal associations such as WhiteHelmets and the Syrian Human Rights Observatory can make the world go round as they want by influencing the policies of world leaders. U.S. freezes funding for Syrias White Helmets. White helmets are terrorists. Former Head of Royal Navy Lord West on BBC White Helmets Aren't Neutral They're On The Side Of The Terrorists.

The summaries provided here are the results for a sample run of the proposed narrative framework and do not reflect authors' personal opinions.