

ECOLE POLYTECHNIQUE FEDERALE DE LAUSANNE
SCHOOL OF COMPUTER AND COMMUNICATION SCIENCES



Master's Thesis in Data Science

**Analysing Climate Change Information Campaigns on Online Social
Networks**

Carried out in the Anyscale Learning For All (ALFA) group
at MIT, Cambridge, Massachusetts
Under the supervision of Dr. Erik Hemberg (MIT) & Jérôme Baudry (EPFL)

Done by
Thomas Benchetrit

Acknowledgements

I am deeply grateful to the whole ALFA team for their warm welcome and immediate integration. Special thanks go to Ethan Garza for his captivating conversations and thoughtfulness, and Sam Laney for allowing me to soak in American culture through humor and baseball. I would also like to thank Steve Moskal for his benevolence, sense of humor, and support during my thesis work. Further, I appreciate Thomas Galligani, Aruna Sankaranarayanan, Shashank Srikant, Michael Wang, and Amanda Abrams for enhancing my stay in many meaningful and unforgettable ways. My past six months with this exceptional team have been a rewarding experience, filled with joy, awe, and sadness upon my departure back to France.

I express my profound gratitude to Una-May O'Reilly for the invaluable opportunity given to me to undertake academic research at ALFA. Her expertise, guidance, and sharpness of mind have greatly contributed to my work. I would also like to deeply thank Erik Hemberg for his supervision. Being under his guidance has been enlightening as I learned how to make puns as much as I learned how to do research. He strikes the right balance between enthusiasm and demands and helped me gain a deeper understanding of the role of a researcher, leading me to question myself to become a better researcher. Finally, I now know that I do not have any problems in my life, just a lot of research questions.

I would also like to express my gratitude to Amber Wang and Rishika Bansal for their collaborative efforts. It was delightful to mentor them, their contributions greatly improved the quality of the work presented in this thesis. Lastly, thanks go to Ozlem Garibay, Ivan Garibay, Chathura Jayalath, Jasser Jasser, Bruce Miller, Alex Baekey, and the entire MIPs ICE team at UCF in Orlando. Their warm support in Orlando made me feel welcome and went beyond the warmth of even Florida in mid-June.

Contents

| | |
|--|-----------|
| Abstract | 8 |
| 1 Introduction | 9 |
| 1.1 Research questions | 11 |
| 1.2 Contributions | 12 |
| 1.3 Roadmap | 12 |
| 2 Background | 14 |
| 2.1 Information Campaigns | 14 |
| 2.2 Online Social Network | 15 |
| 2.2.1 Social Network vs. Information Network | 15 |
| 2.2.2 Directed Graphs as a Social/Information Networks model | 17 |
| 2.3 Large Language Models | 17 |
| 2.4 Transfer Entropy | 19 |
| 3 Related Work | 21 |
| 3.1 Discussion space modeling | 21 |
| 3.1.1 Polarization | 22 |
| 3.1.2 Credibility | 23 |
| 3.2 Large Language Models Bias in content generation | 23 |
| 3.3 Information Flow | 24 |
| 4 Methods | 26 |
| 4.1 Datasets and Properties | 26 |
| 4.1.1 Information Credibility | 26 |
| 4.1.2 Datasets | 27 |
| 4.2 Article representation in the discussion space | 29 |
| 4.3 Content generation | 31 |

| | | |
|----------|--|-----------|
| 4.4 | Explicit measures of Social network graphs | 32 |
| 4.5 | Transfer Entropy Graphs | 33 |
| 4.5.1 | Time series creation | 33 |
| 4.5.2 | Transfer Entropy Graph creation | 35 |
| 5 | Experiments & Results | 38 |
| 5.1 | Content generation | 39 |
| 5.1.1 | Persona modeling | 39 |
| 5.1.2 | Article alteration | 42 |
| 5.2 | Comparing Transfer Entropy graph and Explicit Influence Networks . . . | 42 |
| 5.2.1 | Transfer Entropy as an explicit influence predictor | 42 |
| 5.2.2 | Differences between captured actors | 43 |
| 5.3 | Describing information campaign through captured information flow in TEG | 46 |
| 5.4 | Sensitivity Analysis | 47 |
| 6 | Discussion | 50 |
| 6.1 | Content | 50 |
| 6.1.1 | Describe climate change discussion | 51 |
| 6.1.2 | Content Bias in Language Learning Models | 51 |
| 6.1.3 | Comparing Large Language Models capabilities | 52 |
| 6.2 | Actors | 53 |
| 6.2.1 | Explicit Information Networks and Transfer Entropy Graphs are two different but complementary approaches to study information flow | 53 |
| 6.2.2 | Unveiling the actors in the information flows | 54 |
| 6.3 | Behavior | 55 |
| 6.3.1 | Describing Information Campaigns using Transfer Entropy Graphs | 55 |
| 6.3.2 | TEG method | 56 |
| 6.3.3 | Detecting Emerging disinformation campaigns | 57 |
| 7 | Conclusion | 59 |
| | Appendix | 69 |
| | Reference Articles | 69 |
| | Example of low credibility reference article | 69 |
| | Content Generation | 71 |
| | Article alteration prompts | 72 |

| | |
|-----------------------------------|----|
| Transfer Entropy | 76 |
| Graph topology measures | 76 |
| CCSN for other datasets | 76 |

List of Figures

| | | |
|-----|---|----|
| 2.1 | Representation of a sample of Twitter user graph, [AllThingsGraphed.com(2014)] | 16 |
| 2.2 | Representation of an information network. Red nodes correspond to news sources, and blue to users. [Rabii et al.(2022)] | 17 |
| 4.1 | Timeline of the datasets used | 27 |
| 4.2 | Retweet network of one @guardian tweet | 32 |
| 4.3 | Action time series creation for TEG [Vallez(2023)] | 35 |
| 5.1 | Polarization and credibility of the articles generated using GPT 3.5 | 40 |
| 5.2 | Polarization and credibility of the articles generated using GPT 4 | 41 |
| 5.3 | Comparison of polarization/credibility of the articles made by humans, and those generated by GPT 3.5 and GPT4 | 41 |
| 5.4 | Roc curve for each dataset | 43 |
| 5.5 | We compare the proportion of common actors (y-axis) in the list of the top n actors (x-axis) based on outdegree. For example, if we consider the lists of the 20 most important actors in TEG and EIN for Figure 5.5a only 20% of the actors were captured by both TEG and EIN. | 44 |
| 5.6 | Follower and Tweet count (y-axis) for the actors present in Table.5.5 (outdegree) on the COP27 datasets for EIN and TEG | 45 |
| 5.7 | CCSN for the Skripal dataset | 47 |
| 5.8 | CCSN for the COP27 dataset | 47 |
| 5.9 | Sensitivity analysis for Skripal dataset per T and U action classification threshold. The X-axis is time and the y-axis is CCSN, the different lines indicate different thresholds for T and U. | 48 |
| 7.1 | CCSN for the COP26 dataset | 76 |
| 7.2 | CCSN for the Navalny dataset | 77 |

List of Tables

| | | |
|-----|---|----|
| 3.1 | Related work to polarization and credibility metric | 23 |
| 3.2 | Related work to bias measurment in Large Language Models | 24 |
| 3.3 | Related work to information flow | 25 |
| 4.1 | Datasets used, Size corresponds to the number of tweets inside the dataset | 27 |
| 5.1 | Parameter setups for the study using large language models. | 38 |
| 5.2 | Hyper-parameters used for the generation of <code>TEG.time_window</code> corresponds to the time window used for the generation of CCSN time series, and <code>sampling_freq</code> corresponds to the frequency at which the CCSN measure was computed. Δt , <code>time_window</code> and <code>sampling_freq</code> are given in hours. | 39 |
| 5.3 | Mean polarization and credibility of the generated articles for each prompt and each model. The category corresponds to the prompts used. The "human" category corresponds to the credibility and polarization of the human-made articles. | 40 |
| 5.4 | Average credibility and polarization difference between the articles generated using article alteration prompts and the original articles. Stars (*) indicate statistical significance with a threshold p-value of 0.05 | 42 |
| 5.5 | Table of the most important actors in TEG and EIN information networks for different datasets, based on outdegree measure | 44 |
| 5.6 | Mean Follower and Tweet Counts for Skripal Dataset | 46 |

Abstract

Climate change information impacts public understanding and impacts actions to address climate change. We study information campaigns in climate change discussions on Online Social Networks. We use an information theory based metric, Transfer Entropy, to unveil the actors involved, and describe their behavior. We also highlight Large Language Models capabilities and limitations in generating climate-related content. We describe the behavior of information campaigns through simultaneous activity captured by Transfer Entropy. Moreover, we show that Transfer Entropy helps unveil actors involved in such campaigns. Finally, we observe that with the rise of LLMs as content generators, producing content as credible as human-made articles could still be challenging, and that the generation capabilities are biased across the climate change discussion.

Chapter 1

Introduction

Online social networks (OSNs) are widespread sources of information. An increasing number of people rely almost exclusively on these platforms for getting news[Matsa et al.(2018), Shearer and Gottfried(2017)]. However, in such a vast information ecosystem where fact-checking and moderation are difficult, disinformation, misinformation, and fake news thrive[Moravec et al.(2018), Vrontis et al.(2021)].

An information campaign is a strategic, organized, and targeted initiative intended to disseminate content to a targeted audience[Weiss and Tschirhart(1994)]. Misinformation[per(2020)] can be defined as “a claim that contradicts or distorts common understandings of verifiable facts”, whereas disinformation can be defined as “the subset of misinformation that is deliberately propagated”, to serve a particular narrative for example. On the other hand, fake news, could be defined as “deliberately misleading articles designed to mimic the look of actual articles from established news organizations”, and can be used as a tool for disinformation. Several frameworks [François(2020), Smith et al.(2021)] aim to describe information campaigns, identifying several key factors of such campaigns such as actors, their actions and behavior, and the narrative spread.

The implications of such campaigns on society are broad, ranging from loss of trust in government, the scientific community, or medicine to polarization and division within and across communities[Gupta et al.(2023), Engel et al.(2022)]. Climate change stands as a prominent area where disinformation has been noticeably disseminated[King(2023)]. Several actors, such as companies have played a role in spreading climate change misinformation for several decades[Franta(2021)], interfering with public understanding of the issue. One impact of this disinformation is the mitigation of the severity of climate

change in public perception, potentially delaying actions[Spartz et al.(2017)].

It is challenging and complex to study the dynamics of information campaigns in OSNs. OSN platforms undergo constant changes, including the social network topology and content spread. Moreover, limited access to the OSNs application programming interfaces (APIs) restricts researchers from obtaining information about user interactions and content dissemination. Furthermore, the scale of OSNs increases the complexity, with hundreds of millions of active users and content constantly generated every second. Previous research methods have struggled to effectively handle the volume of data[Ghani et al.(2019)] and often fail to capture the full extent and complexity of information flow within these dynamic information ecosystems.

The investigation of content used for information campaigns presents its own challenges. Though there have been multiple attempts at characterizing narratives about climate change online [Al-Rawi et al.(2021)], there is still a need for a structured and automated approach to evaluation. Modeling the discourse on climate change within a discussion space would enable us to identify types of information campaigns on climate change. A discussion space, in the scope of this thesis, is a N-dimensional space in which one content can be represented as a set of N values. Polarization and credibility offer two aspects in structuring discussions. Polarization categorizes the stance of any given viewpoint within the issue’s spectrum, from denial to urgency, which discloses insight into the nature of the dialogue and its factors such as agreement, divisiveness, and neutrality. On the other hand, credibility assesses the reliability of the content. This evaluation of credibility helps distinguish between reliable and unqualified content, hence helping determine the impact and accuracy of the content. Therefore, through the design of these two axes, one can identify the sentiment or position in the debate and evaluate the weight of its argument based on credibility. This two-dimensional modeling of the discussion space enables comprehension of the online discussion on climate change. However, to assess polarization, most works focus on the network structure[Waller and Anderson(2021), Yang et al.(2017)] to attribute polarization to content. Concerning credibility, most models tend to capture event-specific information from data, resulting in poor generalization on unseen events [Shu and Liu(2019)].

A further point of concern arises with the use of Large Language Models (LLMs) as potential methods of content creation. There are apprehensions about these models being utilized to propagate misinformation [Tang et al.(2023), Stiff and Johansson(2022),

Goldstein et al.(2023)]. While LLMs are indeed capable of producing high-quality content [Xu et al.(2023)], their ability to generate climate-related content specifically is yet to be studied.

1.1 Research questions

in this thesis we aim to answer the following questions:

- Are LLM able to generate credible content centered around climate change?
- How is this content be characterized in the discussion space?
- Is there a bias in LLM-generated content regarding climate change?
- Can we describe the behavior of a narrative on an online social network without having complete information on the underlying network?
- Does Transfer Entropy capture the same information flow as more explicit measures?
- Given its information-theoretic foundations, can Transfer Entropy be used as a complementary method to unveil actors and describe information campaigns?
- How can Transfer Entropy Graphs be used as an initial detector of a potential narrative being spread on OSNs?

To address the challenges of describing and understanding information campaigns in online social networks, we study LLM capabilities to generate climate-related content and design methods to assess their polarization and credibility. We also investigate an information theory-based method, Transfer Entropy (TE). TE captures simultaneous actions over time and can be used to identify and analyze information flows in OSNs with Transfer Entropy Graphs (TEG) [Garibay et al.(2022), JASSER(2023), Gunaratne and et al.(2022)]. This implies that TE can identify talkative actors in OSNs, which can be different from actors with many followers. Thus, TE can be a complimentary measure of information flow by overcoming both the limitations of moreexplicit measures and partial access to OSN platform data.

1.2 Contributions

- Description an automated method to assess polarization and credibility of a content around climate change (Section.4.2). It provides a two-dimensional space to analyze climate-related content.
- Assessment LLM capabilities of climate-related content generation. (Section.6.1.3) We show that it GPT4, under certain conditions, can generate credible content on par with humans.
- Characterization the bias LLMs have on the generation of such content (Section.6.1.2). We assess that LLM are biased toward generating pro-climate change content.
- An analysis of how TEGs capture different sets of actors compared to more explicit influence Networks (Section.6.2.1).
- Examples of TEGs for different information events. (Section.6.3.1)
- A TEG analytic, Count Cut Sum Normalized (CCSN), for detecting information campaigns. The CCSN measure is observed to be able to capture coordinated information campaigns (Section.5.3)

1.3 Roadmap

We first provide an overview of concepts and related works that are used in the thesis (Section.2 and Section.3). Then we describe design the discussion space, by providing a method to measure polarization and credibility (Section.4.2). Then we create specific prompts that are given to LLMs to generate content (Section.4.3). Upon examining the polarization and credibility of the resultant content, we uncover a bias in the LLM-produced material (Section.6.1.2). We then design explicit Influence Networks (Section.4.4) and Transfer Entropy Graphs (Section.2.4) as two methods of modeling information flow. We apply these tools to known information campaigns such as the Navalny case and Skripal assassination attempts to evaluate their properties, and then apply them to climate change-related discussions around COP26 and COP27 to describe information campaigns around those events (Section.6.3.1). With this empirical analysis, we draw out the most prevalent actors (Section.6.2.2) at play within each campaign.

We also study the behavior of these information campaigns (Section.6.3.1). Lastly, we perform a sensitivity analysis of the actions employed in the derivation of Transfer Entropy Graphs, to assess their robustness (Section.5.4).

Chapter 2

Background

This chapter provides an overview of information campaigns and how they operate, focusing on the aspect of disinformation campaigns in Section 2.1. We also introduce the concept of Online Social Networks (OSNs), specifically focusing on Twitter, and distinguish between social networks and information networks within OSNs in Section 2.2. This chapter then provides an introduction to Large Language Models, discussing their structure, training, and use, and more specifically the concept of prompting. In section 2.4, the mathematical concept of Transfer Entropy is introduced and explained.

2.1 Information Campaigns

Information campaigns are purposeful and coordinated actions to spread information, aiming to shape the public opinion, manipulate perceptions, and influence behavior [per(2020)]. They involve systematically disseminating content through diverse communication channels, such as social media or news outlets.

These campaigns are driven by specific agendas and motives, often executed by organizations or state actors with strategic objectives [Ramsay and Robertshaw(2019)]. Disinformation campaigns are a subset of information campaigns, relying on distorted or fabricated information. False information is strategically designed to exploit cognitive biases, societal divisions, and technological vulnerabilities for maximum impact and reach. Disinformation [Bennett and Livingston(2018)] also "disrupts the information flows from government officials through the mainstream press to the public". It therefore serves as a way of imposing another narrative through other information channels. Methods utilized

in disinformation campaigns vary but commonly include creating and amplifying deceptive narratives, using misleading visuals or audiovisual content, impersonating trusted sources, and manipulating online algorithms and networks [Muhammed and Mathew(2022)].

The investigation of information campaigns can encompass its detection, determining which narratives have been propagated, identifying the actors involved, or examining the objective sought by the actors engaged in the campaign. The examination of information campaigns therefore depends upon background knowledge of the information environment. From a structural perspective, it requires to have prior expertise on the potential actors and the way they interact with one another. From a content perspective, it requires identifying which narrative is being spread.

2.2 Online Social Network

An Online Social Network (OSN) forms from a digital platform that enables individuals or entities to connect, interact, and share information through the internet. These networks have become increasingly prevalent, facilitating communication and information dissemination on a global scale. Among the diverse array of platforms, Twitter¹ has emerged as a prominent example due to its user features and widespread user base. Twitter allows users to post short messages, known as tweets. Users can follow other accounts and view a chronological feed of tweets from the accounts they follow. The platform also incorporates features like hashtags, which enable the categorization and discovery of tweets on specific topics, as well as retweeting, which allows users to share others' tweets with their followers. Figure.2.1 represents an example of a Twitter user graph showing the following and content-tagging relationships.

2.2.1 Social Network vs. Information Network

It is important to distinguish between a social network and an information network within the context of OSNs. A social network primarily emphasizes interpersonal connections and social interactions among individuals. It serves as a medium for users to maintain and develop relationships, share personal experiences, and engage in conversations with their peers. On the other hand, an information network emphasizes the dissemination and consumption of information rather than social interactions. While social interactions can

¹<https://twitter.com>

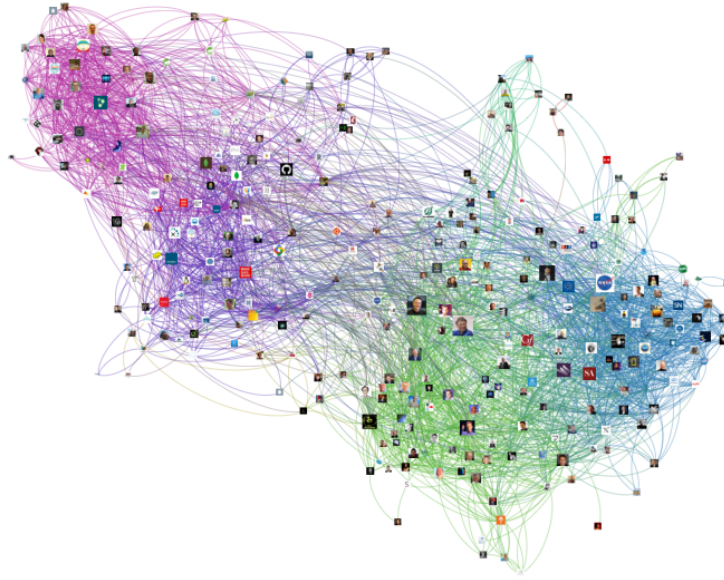


Figure 2.1: Representation of a sample of Twitter user graph,
[AllThingsGraphed.com(2014)]

still occur within an information network, the main focus is on sharing and accessing news, updates, and other informational content. Information networks are accessed by users to follow accounts of news organizations, journalists, subject matter experts, or influencers, i.e. actors that disseminate and spread information. The structure of an information network is shaped by the flow of information rather than interpersonal relationships. Hence an information network tends to have more of a centralized structure, with a small number of nodes being followed by a large portion of the network, as illustrated in Figure.2.2

Twitter supports both social and information networks [Myers et al.(2014)]. Users on Twitter engage in social interactions through tweets, replies, and direct messages, fostering interpersonal connections. Simultaneously, Twitter serves as a tool for the dissemination and consumption of news, current events, and other informational content, with users able to follow accounts that cater to their informational needs and preferences. The interplay between social interactions and information dissemination can impact the diffusion and impact of information campaigns.

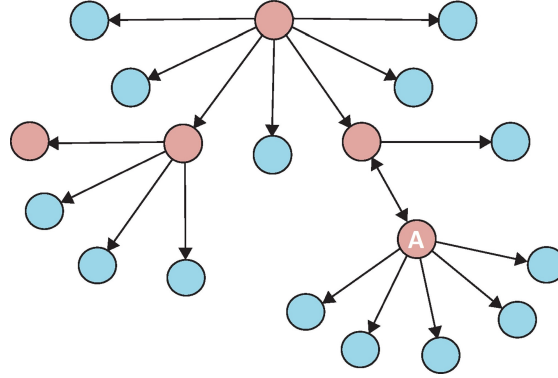


Figure 2.2: Representation of an information network. Red nodes correspond to news sources, and blue to users. [Rabii et al.(2022)]

2.2.2 Directed Graphs as a Social/Information Networks model

A directed graph, denoted as $G = (V, E)$, can serve as a mathematical representation of a social/information network. In this context, vertices (V) represent users/actors, and directed edges (E) capture asymmetric connections like "follower/following", or information flow. The directed edges, denoted as (u, v) , symbolize the relationships between users/actors. These edges indicate that user u follows or is connected to user v . Notably, the direction of the edges signifies the asymmetry in social/information networks, as it does not imply a reciprocal relationship. To understand the potential influence and popularity of users/actors within the network, we examine the out-degree ($d_{\text{out}}(v)$) of a vertex v . The out-degree represents the number of edges originating from user v . Similarly, in-degree ($d_{\text{in}}(v)$) is defined as the number of edges pointing to user v . For example, in Figure.2.2, for node A: $d_{\text{out}}(A) = 6$, $d_{\text{in}}(A) = 1$. Higher out-degree values suggest that a user has a connection with many others, potentially indicating influence.

2.3 Large Language Models

A large language model (LLM) is a type of artificial intelligence model that is characterized by its ability to process and generate human-like text. Examples include OpenAI's GPT-3 and GPT-4[OpenAI(2023)], Meta's LLaMA[Touvron et al.(2023)], Google's PaLM[Chowdhery et al.(2022)] or open-source Bloom[and(2023)]. A large language model is based on a probability distribution over a sequence of tokens in the target language. It is usually designed to optimize the likelihood of the next token at each step in the

sequence, given the previous tokens. A token refers to typically the smallest unit of meaning that the system can understand and process. This could correspond to a word, a subword, or even a single character depending on the granularity of the model. For example, in the English language, in the sentence "The cat sat on the mat", each word ("The", "cat", "sat", "on", "the", "mat") is a token. Let's denote with w the tokens in the sentence, with $w_{<t}$ the sequence of tokens preceding t and w_t the t -th token in the sequence. The model is trained to maximize the log probability of each token in the sequence, given the previous tokens.

$$L(\theta) = \frac{1}{N} \sum_{t=1}^N \log p_{\theta}(w_t | w_{<t}) \quad (2.1)$$

Here, N is the total number of tokens, θ are the parameters of the model and p_{θ} is the conditional distribution parameterized by θ .

In practice, this optimization is carried out using gradient-based methods, specifically maximizing this log-likelihood over a large corpus of sentences.

In Transformer-based large language models such as GPT-3 (Generative Pretrained Transformer 3), the distribution p_{θ} is modeled using a self-attention mechanism that weighs the importance of previous tokens for predicting the next token.

Upon prompting input text, these models predict subsequent words or tokens, learning the syntax and semantics of human language, along with any inaccuracies or biases present in the data they're trained on. Beyond simulating languages, LLMs can have a wide range of applications such as memory aids, drafting emails, coding software, translating languages, answering queries, and document summarization. They might also serve to facilitate natural language processing tasks in industries like healthcare, finance, and entertainment, enabling the creation of AI personal assistants, chatbots, and other interactive tools.

Prompting

Prompting in the context of large language models refers to the process of feeding the model with a sequence of tokens (also known as the prompt) to guide the model's subsequent text generation.

Because LLMs are conditioned on the input tokens to generate the next probable token in a sequence, through a selected input prompt, we can steer the direction of the

model's output. The prompt can be as simple as a few words, a sentence, or a long extract from a document.

Given a prompt $w_{<t^*}$ where t^* is the length of the prompt, the model will then generate a continuation sequence $w_{t^*}, w_{t^*+1}, \dots, w_T$ by repeatedly sampling from the distribution $p(w_t|w_{<t}, \theta)$, where θ are the parameters of the model:

$$\hat{w}_t = \arg \max_{w_t} p_\theta(w_t|w_{<t}), \quad \text{for } t = t^*, t^* + 1, \dots, T \quad (2.2)$$

This computation is usually implemented using a beam search or greedy decoding.

For instance, a prompt like "Translate the following English text to French:" guides the model to activate the learned translation function, while giving a command like "Summarize the following text:" leads it to execute a summarization task.

2.4 Transfer Entropy

Introduced by Schreiber in 2003 [Schreiber(2000)], Transfer Entropy is formally defined between two random discrete variables X and Y with values in sample space Ω_X and Ω_Y respectively as :

$$T_{X \rightarrow Y} = \sum_{y_{n+1}, y_n^{(k)}, x_n^{(l)}} p(y_{n+1}, y_n^{(k)}, x_n^{(l)}) \log \frac{p(y_{n+1}|y_n^{(k)}, x_n^{(l)})}{p(y_{n+1}, y_n^{(k)})} \quad (2.3)$$

where y_n and x_n represent respectively the state of variable Y and X at time n and $y_n^{(k)} = (y_n, \dots, y_{n-k+1})$ and $x_n^{(l)} = (x_n, \dots, x_{n-l+1})$ are short-hand for k -dimensional and l -dimensional delay embedding vectors of state. The sum runs on all possible combinations of states that variables X and Y can assume at adjacent times. This means that the sum runs on all possible vectors in the space $\Omega_Y \times (\Omega_Y)^k \times (\Omega_X)^l = (\Omega_Y)^{k+1} \times (\Omega_X)^l$ such that $(y_{n+1}, y_n^{(k)}, x_n^{(l)}) \in (\Omega_Y)^{k+1} \times (\Omega_X)^l$.

In practice, we chose $k = l = 1$ for computational reasons, which means a 1-dimensional delay. This approximation reduces general Transfer Entropy to

$$T_{X \rightarrow Y} = \sum_{y_{n+1}, y_n, x_n} p(y_{n+1}, y_n, x_n) \log \frac{p(y_{n+1}|y_n, x_n)}{p(y_{n+1}, y_n)} \quad (2.4)$$

in this thesis, the variables X and Y represent the action count and take values in $\Omega_X = \Omega_Y = \mathbb{N}$. To simplify the sample space and bound the computations, we follow an

approach used in computational neuroscience[Panzeri and et al.(2007)] and binarize our variables as shown in Eq. 2.5.

$$X = \begin{cases} 1 & \text{if } C(\Delta t_i) > 0 \\ 0 & \text{else} \end{cases} \quad (2.5)$$

with $C(\Delta t_i)$ being the total number of actions in the considered time bin Δt_i

Therefore, instead of looking at the number of tweets posted at each time bin Δt_i , we check if *at least one tweet was posted* in the bin Δt_i . Using such a binarization, Eq.2.3 reduces to:

$$T_{X \rightarrow Y} = \sum_{(y_{n+1}, y_n, x_n) \in \{0,1\}^3} p(y_{n+1}, y_n, x_n) \log \frac{p(y_{n+1}|y_n, x_n)}{p(y_{n+1}, y_n)} \quad (2.6)$$

in which the triplet (y_{n+1}, y_n, x_n) can take values in at most $2^3 = 8$ possible states. To evaluate the probabilities in Eq.2.6, we empirically count the frequency of each state to estimate its probability for each of the 8 possible states and deduce the probabilities from it. We can only observe the time series for one actor once, in contrast to repeating a stochastic experiment multiple times and using Monte-Carlo methods to estimate probabilities of the underlying time series.

Chapter 3

Related Work

In this chapter, a review of related work is presented. In Section 3.1, we delve into both how previous have tried to assess polarization and credibility of content. Section 3.2 focuses on several studies addressing biases in LLMs and their mitigation efforts, indicating a lack of research done in the specific context of biases in generating climate change content. Lastly, Section 3.3 discusses the use of transfer entropy for modeling information and influence dynamics in online social networks. The related works are summarized in Tables 3.1,3.2,3.3. The current literature, to the best of our knowledge, lacks specific research on assessing content credibility and detecting polarization in the context of climate change discussion based on content embeddings. Also, there is a lack of analysis on bias within large language models (LLMs) specifically for content generation regarding climate change. Finally, the use of Transfer Entropy to study information climate change in this context is still missing characterization. Our research addresses these gaps by developing methods for credibility assessment and polarization detection based on content, as well as analyzing LLM biases within the generation of climate change content. We also provide a study on the use of TE to study information campaigns related to climate change.

3.1 Discussion space modeling

In this section, we present the other studies carried around representation discussion space around credibility and polarization, summarized in Tab.3.1.

3.1.1 Polarization

To the best of our knowledge, there has been limited effort put into developing methods to automatically detect polarization based on content regarding climate change. Instead, most research focuses on analyzing the polarization of users in online social networks, based on network structure. In [Waller and Anderson(2021)], the authors developed a neural-embedding methodology to measure the social organization of online communities on Reddit. They analyzed 5.1 billion comments made in 10,000 communities over 14 years and found that the 2016 US presidential election prompted significant polarization event on Reddit, with right-wing activity driving the increase in polarization.

On the other hand, [Yang et al.(2017)] focused on quantifying content polarization on Twitter during the 2016 US presidential election from a content perspective. They proposed a new approach to measure the polarization of content semantics using word embedding representation and clustering metrics. Their evaluation framework validated the proposed measurement and explored the extent of content polarization across time, geography, and different types of users.

Both papers contribute to understanding polarization on online platforms, but they tackle it from different angles. While [Waller and Anderson(2021)] focuses on the social organization and dynamics of online communities, [Yang et al.(2017)] examines the polarization of content semantics. These studies provide insights into the polarization phenomenon on different platforms (Reddit and Twitter) and during the same political event (2016 US presidential election).

In [Ding et al.(2023)], the authors investigate the relationship between broadcast news media language and social media discourse. They analyze closed captions from CNN and Fox News, along with corresponding discourse from Twitter, and compare how topically important keywords were discussed between the two channels, demonstrating a sharp increase in polarization, especially after 2016. Even though this approach measures polarization based on the content of the captions, it relies on sampled keywords that may not be representative of the captions.

Our approach builds upon these previous works by specifically aiming to detect polarization in online discussions based solely on content, by relying on embeddings that are believed to better capture the meaning of the studied news articles. Also, we study polarization in the context of climate change, which is not covered by the previous studies.

3.1.2 Credibility

In existing literature, numerous studies evaluated the credibility of online content using various approaches. The research by [Savolainen(2023)] focuses on assessing the credibility of misinformation about COVID-19 vaccines on Reddit, with findings indicating the role of author’s reputation, expertise, and integrity in the assessment process. Similarly, [Eastin(2001)] examines the influence of source expertise and content knowledge on perceptions of online health information credibility. It suggests that there’s a tendency to rate the majority of online health information as relatively credible regardless of source expertise or familiarity with content.

The research by [Gupta et al.(2015)] introduces a method for assessing the credibility of tweets in real time during emergencies, using assorted features related to user behavior and tweet content, topic, and diffusion. Additionally, [Yang and Menczer(2023)] investigates whether large language models, like ChatGPT, can assess the credibility of news outlets by demonstrating a correlation between ratings provided by the model and human experts.

While these studies provide insights into the assessment of online content credibility based on a variety of factors, none of them have specifically concentrated on assessing credibility based on content embedding and its similarity with reference articles, and in the context of climate change discussion.

| Metric | Method | Input | Domain | Paper |
|--------------|------------------------------|-------------------------|--------------------|-----------------------------|
| Polarization | Network Embedding | Graph | US Elections 2016 | [Waller and Anderson(2021)] |
| Polarization | Word Embedding | Tweets | US Elections 2016 | [Yang et al.(2017)] |
| Polarization | Closed Captions’ Keywords | News channels’ captions | News and Twitter | [Ding et al.(2023)] |
| Credibility | Author Reputation Evaluation | Reddit Covid-19 posts | Health/Reddit | [Savolainen(2023)] |
| Credibility | Source/Content Evaluation | Online health info | Online health info | [Eastin(2001)] |
| Credibility | Real-time Tweets Evaluation | Emergency tweets | Emergency/Twitter | [Gupta et al.(2015)] |
| Credibility | Large Language Model | News outlet ratings | News credibility | [Yang and Menczer(2023)] |

Table 3.1: Related work to polarization and credibility metric

3.2 Large Language Models Bias in content generation

The related literature on bias in the outputs of large language models is extensive, investigating the causes, implications, and mitigation strategies of these biases. First, [Ferrara(2023)] provides an overview of the challenges and risks related to biases in large language models like ChatGPT, including their origins, ethical concerns, and potential approaches to identify, quantify, and mitigate them. This work highlights the importance

of a multi-disciplinary approach to develop more equitable, transparent, and responsible AI systems.

Other work [Kirk et al.(2021)] investigates the biases existing in popular language models, like GPT-2, when used 'out-of-the-box' for downstream tasks. In their study, they show that GPT-2 reflects the skewed gender and ethnicity distribution found in US Labor Bureau data. On the other hand, [Abid et al.(2021)] explores the largely overlooked area of religious bias in large language models, outlining bias toward Muslim people in GPT-3. It shows that biases are not limited only to aspects such as race and gender, but also extend to religious matters.

Beyond the English-speaking context, [Névéol et al.(2022)] focuses on extending a dataset for measuring social bias to French, pointing out that bias studies require a multilingual approach to account for specific societal and cultural peculiarities of different countries.

Finally, [Dev et al.(2022)] presents a practical framework to measure biases in NLP technologies, offering a comprehensive understanding of both intrinsic and extrinsic bias measures in different application contexts.

All the included works emphasize the complexities of understanding, quantifying, and mitigating biases in language models. However, none of them examines these biases in the context of content generation about climate change. The research we present in our paper fills this gap by specifically focusing on climate change related content generated by these models.

| Model | Domain | Paper |
|--------------|-------------------------|-----------------------|
| ChatGPT | Overview | [Ferrara(2023)] |
| GPT 2 | Gender & Ethnicity | [Kirk et al.(2021)] |
| GPT-3 | Religious Bias | [Abid et al.(2021)] |
| N/A | Bias in French Language | [Névéol et al.(2022)] |
| N/A | Measuring Bias | [Dev et al.(2022)] |

Table 3.2: Related work to bias measurment in Large Language Models

3.3 Information Flow

A recent study proposed the use of transfer entropy for modeling information and influence dynamics in online social networks (OSNs) without relying on explicit network structures provided by individual platforms [Senevirathna et al.(2021)]. This approach allows for

analyzing network dynamics agnostically, regardless of the specific platform and without the need for explicit relationships from the platform.

Traditional methods for studying OSNs depend on explicit causal knowledge, which may not always be available. To address this, researchers have explored coupling inference methods to assess the strength of interactions between users. Granger causality, a well-known method, measures the ability to predict future values of one-time series using prior values of another[Granger(1969)]. Transfer entropy, a measure derived from information theory, quantifies the directed transfer of information between evolving systems and is equivalent to Granger causality for Gaussian variables[Schreiber(2000)]. Other coupling inference methods based on distance measures in high-dimensional spaces, such as cross-mapping and joint distance distribution (JDD), have also emerged[Amigó and Hirata(2018), Sugihara and et al.(2012)].

The use of transfer entropy has gained popularity for inferring social influence in OSNs. For instance, transfer entropy has been applied to analyze user activity on Twitter, specifically focusing on tweets containing URLs, to establish influence relationships between individuals within the network[Ver Steeg and Galstyan(2012)]. Similar studies have been conducted using data from Tencent Weibo, a Chinese social media platform, demonstrating the effectiveness of transfer entropy in quantifying influence dynamics[He and et al.(2013)]. Additionally, transfer entropy has been employed to measure the impact of one user’s tweet content on another, proving to be a reliable predictor of user mentions[Ver Steeg and Galstyan(2013)].

| Method | Input | Output | Paper |
|--------------------|---------------------------|----------------------------------|---|
| Transfer entropy | OSN data | Measures of information transfer | Senevirathna et al (2021) |
| Granger causality | Time series values | Predictive ability | Granger (1969) |
| Cross mapping, JDD | Lorenz systems simulation | Interaction strength measures | Amigo et al (2018), Sugihara et al (2012) |
| Transfer Entropy | User activity (Twitter) | Influence relationships | Versteeg et al (2012) |
| Transfer entropy | Tencent Weibo data | Influence dynamics | He et al (2013) |
| Transfer entropy | User tweet content | User mention predictor | Versteeg et al (2013) |

Table 3.3: Related work to information flow

Chapter 4

Methods

We present in this chapter several datasets focusing on recent climate change events and international political incidents 4.1. We develop a model to represent articles in a discussion space as a measure on two primary axes: credibility and polarization 4.2. We experiment with two methods of content generation, persona modeling and article alteration, to explore the ability of language learning models to generate polarized and credible content 4.3. Additionally, we describe explicit measures of influence in social network graphs 4.4 and present the process of creating transfer entropy graphs 4.5.

4.1 Datasets and Properties

4.1.1 Information Credibility

NewsGuard¹ is a company that aims to provide an assessment of news sources' credibility. Their journalists reviewed and rated a wide range of news sources, which makes up around 95% of online engagement. The reviewing system uses 9 criteria to assign a "trust score" between 0 and 100 to each news source. These ratings were done by humans. NewsGuard considers news websites scoring above 60 as trustworthy (T) and untrustworthy (U) otherwise. The dataset includes ratings from NewsGuard gathered in September 2022, consisting of 8,145 unique news domains that have been fully rated.

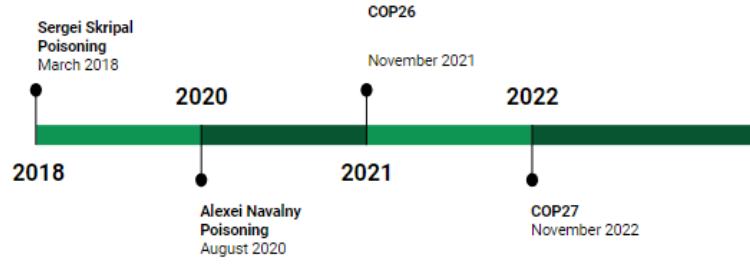


Figure 4.1: Timeline of the datasets used

| Dataset | Scheduled | Topic | Year | Start date | End date | Size |
|---------|-----------|----------------|------|------------|------------|-----------|
| Skripal | No | poisoning | 2018 | 2018-03-05 | 2018-03-31 | 677,239 |
| Navalny | No | poisoning | 2020 | 2020-08-20 | 2020-09-14 | 1,146,731 |
| COP26 | Yes | climate change | 2021 | 2021-10-18 | 2021-11-25 | 4,115,553 |
| COP27 | Yes | climate change | 2022 | 2022-10-24 | 2022-11-24 | 2,524,381 |

Table 4.1: Datasets used, **Size** corresponds to the number of tweets inside the dataset

4.1.2 Datasets

They are extracted from two datasets of tweets, one relating to potential cases of climate change information campaigns curated by [Vallez(2023)] and one relating to covert actions and state-sponsored assassination attempts curated by [Galligani(2023)]. Specifically, we take six case studies of time-bounded events where we have some prior evidence or interest in identifying coordinated influence campaigns on social media. They are summarised in Tab.4.1 and Fig.4.1.

Climate Change: COP26 & COP27

Our focus is on the last two COP gatherings, specifically **COP26** in Glasgow, Scotland, held from October 31, 2021, to November 12, 2021, and **COP27** in Sharm El Sheikh, Egypt, held from November 6, 2022, to November 18, 2022. For each of these events, we utilized the Twitter API to retrieve tweets containing explicit terms such as "climate change," "climate crisis," "climate emergency," "climate action," or "global warming" (including relevant hashtags). We specifically targeted tweets written in English that included at least one URL. We obtained tweets both preceding and succeeding the events, maintaining an equivalent number of days before and after each event. COP27 dataset is composed of 2 524 381 tweets, and COP26 is composed of 4 115 553 tweets.

¹<https://newsguardtech.com>

Assassination attempts

Below, we briefly characterize each of these case studies, their timeline of events, see Figure 4.1, and the extent to which there exists prior evidence of a coordinated influence campaign.

Sergei Skripal The poisoning of Sergei Skripal, a former Russian spy who had defected to the United Kingdom and was poisoned in March 2018 in Salisbury, England. Skripal was killed on 4 March 2018, which resulted in standard news coverage and discourse in the following weeks. It was, however, not until 12 March that the British prime minister publicly attributed the attack to the Russian Government, after which began a clear and well-studied [Ramsay and Robertshaw(2019)] influence campaign by Russia targeting the United Kingdom. This provides a tightly-bounded natural experiment that will allow comparison between “normal” online discourse on the topic (before 12 March) and online discourse in the presence of a concerted Russian influence operation.

The techniques of Russian influence operations are well documented and this case in particular is extremely well studied. This allows us to very clearly tune our methods to identify particular phenomena and test against a relatively well-established ground truth, which is often unavailable for this kind of research. In particular, Ramsay and Robertshaw identify two connected techniques described in more detail in [Paul and Matthews(2016)] and [Global Engagement Center(2020)], the “firehose of falsehood” (spreading huge amounts of contradictory claims in the hopes that the net effect is that the targets cannot determine the truth) and “reputation laundering” or “churnalism”.

We used the Twitter API to collect all tweets referencing the words “skripal” or “novichok” and the associated hashtags (our exact search query was “skripal OR #skripal OR #novichok OR novichok”) between 3 March 2018 (the day before his poisoning) and 1 April 2018 (about two weeks after the beginning of the influence operation). We collected 677,239 tweets with this query.

Alexei Navalny Alexei Navalny, a prominent Russian opposition politician and anticorruption activist, was poisoned on 20 August 2020. Navalny fell ill during a flight from Siberia to Moscow and was subsequently hospitalized in a coma. He was eventually transported to Germany for medical treatment where he recovered and bravely returned to Russia in January 2021 where he was arrested on pretextual

charges[alj([n. d.])]. Some work has investigated the online information environment surrounding Navalny’s attempted assassination and have identified evidence of campaigns targeting English speaking, Western audiences, seeking to vilify Navalny using similar techniques to other Russian influence operations such as “flooding the zone” with bot activity, particularly after 2 September (when German chancellor Angela Merkel publicly announced that he had been poisoned with Novichok and attributed the poisoning to Russia)[Alieva and Carley(2021)]. We include this case study in our dataset as a comparison to the Skripal influence operation, which had a comparatively constrained target audience (primarily just UK citizens) and is relatively better studied. We used the Twitter API to collect all tweets referencing the words “navalny” or “novichok” and the associated hashtag in English (our exact search query was “navalny OR #navalny OR #novichok OR novichok”) between 20 August 2020 (the day he was poisoned) and 15 September 2020 (about two weeks after Germany attributed his poisoning to Russia). We collected 1,146,731 tweets with this query.

4.2 Article representation in the discussion space

We model each content c either human-written or LLM-generated in the discussion space that has two axes: credibility (C) and polarization (P). For a specific measure $M = C, P \in [-1, 1]$, we assemble a specific list of reference articles $a_i \in L_m, i = 1, \dots, N_m$, for which we know their measure score $M(a_i)$. Each reference article was drawn from the climate-change-related discussion around COP27. Then, for each article, we compute its embedding using OpenAI’s `text-embedding-ada-002` text-embedding model. The embedding of an article a_i will be referenced as \mathbf{a}_i . This measure is defined as a weighted sum of similarities between the content and the list of reference articles $\text{sim}(\mathbf{c}, \mathbf{a}_i)$ between the embedding of the generated content \mathbf{c} and the embedding of the references. The similarities are then normalized with min-max normalization. The similarity measure used is cosine similarity. To deal with the imbalance in the measures for the list of references, we weigh these similarities with a measure of sample rarity r_i [Steininger et al.(2021)]. Let us define $p_i = p(M(a_i))$ the probability to observe the measure $M(a_i)$ in L_m . We then estimate p using Kernel Density Estimation (KDE)[Parzen(1962)]. The similarities

are therefore computed as such, for each content and reference article a_i :

$$sim(\mathbf{c}, \mathbf{a}_i) = \frac{\mathbf{c} \cdot \mathbf{a}_i}{\|\mathbf{c}\| \|\mathbf{a}_i\|} \quad (4.1)$$

which are then normalized to get:

$$norm_sim(\mathbf{c}, \mathbf{a}_i) = \frac{sim(\mathbf{c}, \mathbf{a}_i) - \min_{a_i} sim(\mathbf{c}, \mathbf{a}_i)}{\max_{a_i} sim(\mathbf{c}, \mathbf{a}_i) - \min_{a_i} sim(\mathbf{c}, \mathbf{a}_i)} \quad (4.2)$$

The rarities r_i are computed as such:

$$r_i = \frac{1 - p_i}{\frac{1}{N} \sum_{i=1}^N 1 - p_i} \quad (4.3)$$

The measure $M(c)$ of a new content \mathbf{c} is define as such:

$$M(c) = \frac{\sum_{i=1}^N k_i M(a_i)}{\sum_{i=1}^N k_i} \quad (4.4)$$

where k_i is computed as such:

$$k_i = r_i norm_sim(\mathbf{c}, \mathbf{a}_i) \quad (4.5)$$

For polarization, the measures $P(a_i), i = 1, \dots, N_p$ were obtained through human labeling on a scale from -1 to 1, $P = -1$ represents an article extremely polarized as a climate-change contrarian, whereas $P = 1$ represents an article extremely polarized as a climate-change advocate. A polarized contrarian point of view was defined as a denial, dismissal or unwarranted doubt of the scientific consensus on the rate and extent of global warming, its significance, or its connection to human behavior, in whole or in part. This point of view also states that the current discussion around climate change may be filled by hoaxes. A polarized advocate point of view was defined as a belief in the scientific consensus on the rate and extent of global warming, its significance, or its connection to human behavior, in whole or in part. This point of view also states that the current discussion around climate change is more than necessary and urges us to take action as soon as possible.

For credibility, the measures were taken from the Newsguard ratings presented in Section 4.1.1, scaled to the $[-1, 1]$ range. A credible article can be described as an article

written by someone who is an expert in their discipline and is free of errors and bias. It contains facts and scientific evidence of what it states, in a clear manner. It steers away from shortcuts in its argumentations and stays objective and with no bias on the facts that it states. An example of an article scoring a -1 credibility can be found in Appendix.7.

4.3 Content generation

To observe the LLM’s capability of generating polarized, credible content, two LLM prompting methods were experimented with. The first one was person modeling and the second experiment was altering a pre-existing human-made article.

Persona modeling

In persona modeling, we asked the LLMs to pretend to be another persona who either believes in or does not believe in the existence of climate change. The model was also given a list of criteria for things to include in article generation, including evidence from reputable sources, multiple arguments, counterarguments and rebuttals, and both an introduction and conclusion. To create more polarizing content, the model was asked to be a journalist who held polarizing views and wrote articles with strong language and arguments. Four prompts were tried, **advocate**, **pol_advocate**, **contrarian**, and **pol_contrarian**.

The prompts can be found in the Appendix.

Article alteration

Another method to produce an article was to base the model on an already existing article present in the dataset and generate a new article based on it. For a sample of articles picked randomly from the dataset, we asked the model to generate four variations of it. The LLM was prompted to generate an article that would be more or less credible, and more advocate or more contrarian.

Let $L_m = \{a_1, a_2, \dots, a_n\}$ be the list of reference articles. In this experiment, we apply four prompts denoted as $P_k = \{\text{more_advocate}, \text{more_contrarian}, \text{more_credible}, \text{less_credible}\}$ to each article a_i in L_m using a Language Learning Model (LLM). For each article and each prompt, we produce a variant of the article.

Let's represent the LLM as a function G that takes an article and a prompt as its arguments and gives a variant of the article. Therefore, for each article a_i and each prompt p_j :

$$v_{ij} = G(a_i, p_j),$$

where v_{ij} is the variant of the i -th article generated with the j -th prompt.

As a result, we obtain a set of article variations for each a_i as $\{v_{i1}, v_{i2}, v_{i3}, v_{i4}\}$ corresponding to the prompts from P_k . The prompts used can be found in the Appendix.

4.4 Explicit measures of Social network graphs



Figure 4.2: Retweet network of one **@guardian** tweet

To better understand the information flow and actors captured by the Transfer Entropy graphs, we compare them with a more explicit measure of influence, named Explicit Influence Networks (EINs), which are aggregated retweet networks.

A retweet network can be represented as a directed graph denoted by $G_r = (V_r, E_r)$, where V_r represents the set of nodes and E_r denotes the set of directed edges. Each node $a_i \in V_r$ corresponds to an actor involved in retweeting, while each directed edge $(a_i, a_j) \in E_r$ signifies that actor a_j retweeted a tweet originating from actor a_i .

To construct the EIN, we followed a step-by-step procedure:

1. For each tweet t_k , we extracted its corresponding retweet network $G_{r_k} = (V_{r_k}, E_{r_k})$ as depicted in Figure 4.2.
2. By considering multiple retweet networks, we combined them into a unified network

denoted as $G_u = (V_u, E_u)$, where

$$V_u = \bigcup_k V_{r_k}, \quad E_u = \bigcup_k E_{r_k}$$

3. For actors a_i and a_j in the unified network, we say that actor a_i influences actor a_j if there exists a directed path from a_i to a_j in the unified network. This can be expressed as:

$$a_i \text{ influences } a_j \quad \Leftrightarrow \quad \exists \text{ directed path from } a_i \text{ to } a_j \text{ in } G_u.$$

This procedure can be found in the pseudo-code algorithm presented in Alg.1

4.5 Transfer Entropy Graphs

In this section, we describe how a transfer entropy graph is created. First, we present the time series creation. Then we present the graph creation.

4.5.1 Time series creation

The key components for creating time series are:

Actors We define actors to be entities spreading information. The domain of actors in this thesis is all unique Twitter users present in the dataset. We filter out the actors to have a minimum activity rate (or tweet count) of 3 tweets in the corresponding time period.

Actions We define an action as an observable activity by an actor. In our case, the actions that we observe are the sharing of a link from a news source on Twitter, either by tweeting or retweeting it. All tweets in each dataset contain at least one link. Each tweet is then attributed an action of T if it shares news from trustworthy outlets, and U if it shares news from websites that are considered as not trustworthy by NewsGuard (4.1.1).

Stratification We define stratification as a way of segregating the data to compare results across different sets. In the following study, we use temporal stratification, meaning that we split the data based on the tweet/retweet timestamp. Let a

dataset D be the set of tweets in the overall time period T . We split the time period into n sub-intervals that are each shifted by $t_f = 1$ day. Each time interval $\Delta_t = 2$ days. They are defined as follows:

$$T_i = [t_0 + i \times t_f, \min(t_0 + i \times t_f + \Delta_t, t_n)] \quad (4.6)$$

where t_0 represents the date of the first chronological tweet in the dataset and t_n represents the date of the last chronological tweet.

Time series creation After identifying each tweet's actor, action, and strata, we discretize each time strata $T_s = [t_0, t_f]$ with a time resolution of $dt = 2$ hours. For a given time resolution dt , we create bins $\{\Delta t_i\}_{i=0}^N$ of size dt such that

$$\Delta t_i = [t_i, t_{i+1}], \quad i = 0, 1, \dots, N \quad t_i = t_0 + i \cdot dt, \quad i = 0, 1, \dots, N + 1 \quad (4.7)$$

Then, for each actor, we create a time series for its actions by counting the number of tweets for which the timestamp is contained in each bin Δt_i . Let $X_{a_j}^A \in \mathbb{R}^{N+1}$ be the time series for the actor A and action a_j :

$$(X_{a_j}^A)_i = |\{x \in D_{s,i} | x_{\text{action}} = a_j, x_{\text{actor}} = A\}|, \quad i = 0, \dots, N \quad (4.8)$$

where $|\cdot|$ denotes the number of elements in the set and $D_{s,i}$ the set of tweets that have been created in the strata s and time bin Δt_i .

The time series creation process is summarised in Fig.4.3

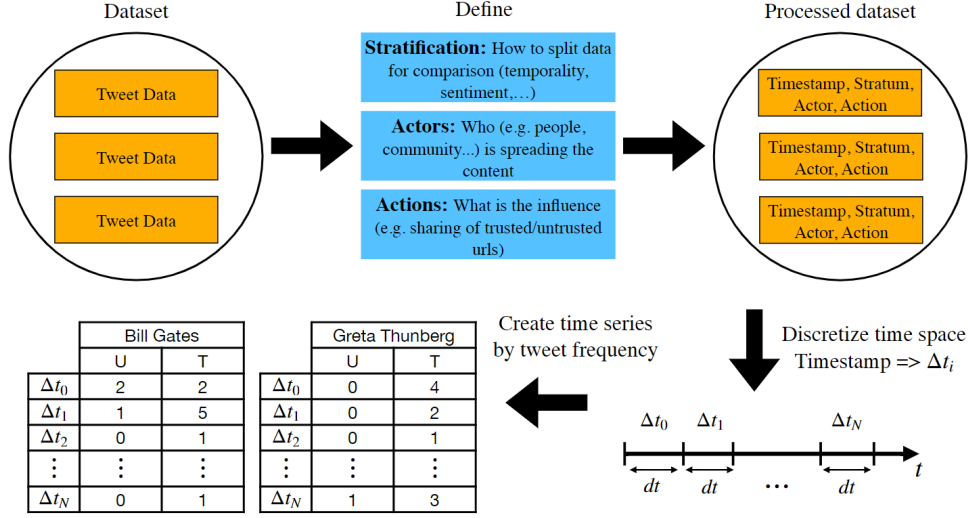


Figure 4.3: Action time series creation for TEG [Vallez(2023)]

4.5.2 Transfer Entropy Graph creation

From the derived time series, we use an inference method, Transfer Entropy, to detect how actions from one actor might have "caused" actions by another actor. Suppose two actors a_i and a_j , and actions T-U as defined in (methods action ref), then C from a_i to a_j is defined as the square matrix:

$$C_{a_i \rightarrow a_j} = \begin{bmatrix} TE(X_U^{a_i}, X_U^{a_j}) & TE(X_U^{a_i}, X_T^{a_j}) \\ TE(X_T^{a_i}, X_U^{a_j}) & TE(X_T^{a_i}, X_T^{a_j}) \end{bmatrix} \quad (4.9)$$

where $TE(X_i, X_j)$ represent the Transfer Entropy value between time series X_i and X_j . Note that Transfer Entropy is not a symmetric inference method, i.e $TE(X_i, X_j) \neq TE(X_j, X_i)$

We then construct an influence graph for each stratum as follows:

- Compute the pairwise matrix between all actors
- Choose a cutoff value at which two time series are considered causal.
- For each action combination (T-T, U-U, T-U, U-T), define a directed graph $G = (V, E)$ where the nodes V are the actors and E the directed edges between actors.

Each edge $e \in E$ is defined as

$$\exists e_{(a_i, a_j)} \in E \Leftrightarrow TE(X_A^{a_i}, X_B^{a_j}) > c \quad (4.10)$$

where a_i, a_j are actors, A/B one action pair and c an arbitrarily chosen cutoff value.

The procedure is also described in Algorithm.2

Algorithm 2 Transfer Entropy graph generation

Input: $X_i : \text{TimeSerie}, X_j : \text{TimeSerie}$

Output: CouplingMatrix

function COMPUTEPAIRWISEMATRIX(X_i, X_j)

$matrix \leftarrow \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$

for action_{*i*} in X_i **do**

for action_{*j*} in X_j **do**

$matrix[\text{action}_i][\text{action}_j] \leftarrow \text{transfer_entropy}(X_i[\text{action}_i], X_j[\text{action}_j])$

return $matrix$

Input: cutoff : float, actions : List[Action], actors_timeseries : Dict[TimeSerie]

Output: TEGs per action

function GENERATETEG(cutoff, action_list, actors_timeseries)

actors_matrices \leftarrow dict()

for a_{*i*} in actors_timeseries **do**

for a_{*j*} in actors_timeseries **do**

if a_{*i*} \neq a_{*j*} **then**

ts_{*i*} \leftarrow actors_timeseries[a_{*i*}]

ts_{*j*} \leftarrow actors_timeseries[a_{*j*}]

actors_matrices[(a_{*i*}, a_{*j*})] \leftarrow compute_pairwise_matrix(ts_{*i*}, ts_{*j*})

transfer_entropy_graphs $\leftarrow \emptyset$

for action₁ in action_list **do**

for action₂ in action_list **do**

graph \leftarrow Graph(V, E)

for actor_{*i*}, actor_{*j*} in actors_matrices **do**

if actors_matrices[(actor_{*i*}, actor_{*j*})] [action₁] [action₂] > cutoff **then**

graph.add_edge(actor_{*i*}, actor_{*j*})

transfer_entropy_graphs[(action₁, action₂)] \leftarrow graph

return transfer_entropy_graphs

Algorithm 1 The Explicit Influence Network (EIN) is formed by aggregating retweet networks.

Input: Tweet List $L=[t_1, t_2, \dots, t_n]$
Output: Unified network $G_u = (V_u, E_u)$
Function `construct_unified_network(L)`:
 $V_u \leftarrow \emptyset$
 $E_u \leftarrow \emptyset$
for t_k **in** L **do**
 $G_{r_k} \leftarrow \text{extract_retweet_network}(t_k)$
 $V_u \leftarrow V_u \cup V_{r_k}$
 $E_u \leftarrow E_u \cup E_{r_k}$
return $G_u = (V_u, E_u)$

Input: Unified network G_u , two actors a_i, a_j
Output: Influence value
Function `determine_influence(G_u, a_i, a_j)`:
if \exists path from a_i to a_j in G_u **then**
 return True
else
 return False

Input: Tweet List $L=[t_1, t_2, \dots, t_n]$
Output: Influence network $G_i = (V_i, E_i)$
Procedure `construct_influence_network(L)`:
 $G_u \leftarrow \text{construct_unified_network}(L)$
Initialize empty influence network $G_i = (V_i, E_i)$
 $V_i \leftarrow V_u$
 $E_i \leftarrow \emptyset$
for each actor $a_i \in V_i$ **do**
 for each actor $a_j \in V_i$ **do**
 if `determine_influence(G_u, a_i, a_j)` **then**
 Add directed edge (a_i, a_j) to E_i
return $G_i = (V_i, E_i)$

Chapter 5

Experiments & Results

In this section, we present our experiments. In Section 5.1.1, we observe the credibility and polarization of the LLM-generated content using the persona prompts. In Section 5.1.2, we compare the difference in credibility and polarization between human-made articles and their LLM-generated variations. In Section 5.2 we compare transfer entropy graphs (TEG) and explicit influence networks (EIN). In Section 5.3 we describe information flow in the context of information campaigns. In Section 5.4 we present a sensitivity analysis of TEG. Note, results are discussed in Section 6.

Using the datasets described in Section 4.1, we produced for each an EIN, and TEGs, with the parameters that can be found in Tab.5.2 We also present in Tab.5.1 the parameters used to generate and study LLM generation capabilities.

| LLM hyper-parameters | | |
|---------------------------------------|---------------|-------|
| model | gpt-3.5-turbo | gpt-4 |
| # generated Persona exp. | 15 | 10 |
| # generated Alteration exp. | 40 | |
| # Polarization ref articles (N_p) | 50 | |
| # Credibility ref articles (N_c) | 298 | |

Table 5.1: Parameter setups for the study using large language models.

| TEG hyper-parameters | |
|----------------------|------|
| Δt | 2 |
| time_window | 48 |
| sampling_freq | 24 |
| # actors | 1500 |
| min.tweets | 3 |

Table 5.2: Hyper-parameters used for the generation of TEG. `time_window` corresponds to the time window used for the generation of CCSN time series, and `sampling_freq` corresponds to the frequency at which the CCSN measure was computed. Δt , `time_window` and `sampling_freq` are given in hours.

5.1 Content generation

5.1.1 Persona modeling

In Fig.5.1, 5.2, and Fig. 5.3, credibility and polarization are represented for articles produced by GPT 3.5, GPT 4, and human authors. The X and Y axes stand for polarization and credibility, respectively, with each dot representing a single article. In Fig.1 and Fig.2, purple dots represent human-authored articles from the dataset. Fig.3 presents a comparison of the polarization and credibility of articles written by GPT 3.5, GPT 4, and humans.

Tab5.3 displays the average credibility and polarization for each prompt and model. According to this table, both models create articles that, for each prompt, have lower credibility than the human-authored articles in the dataset. Comparative to GPT 4, articles produced by GPT 3.5 were on average less credible. Overall, generated content had less credibility than human-made articles. Furthermore, articles produced based on polarized prompts were less credible than their non-polarized equivalents. Moreover, for all models, articles with a contrarian point of view had lower credibility than those with an advocate point of view.

| Category | Credibility | | Polarization | |
|----------------------|-------------|-------|--------------|-------|
| human | 0.080877 | | -0.006372 | |
| | GPT3.5 | GPT4 | GPT3.5 | GPT4 |
| Contrarian | -0.50 | -0.23 | -0.53 | -0.90 |
| Contrarian polarized | -0.564018 | -0.25 | -0.59 | -0.94 |
| Advocate | -0.24 | 0.04 | 0.74 | 0.62 |
| Advocate polarized | -0.30 | -0.04 | 0.51 | 0.34 |

Table 5.3: Mean polarization and credibility of the generated articles for each prompt and each model. The category corresponds to the prompts used. The "human" category corresponds to the credibility and polarization of the human-made articles.

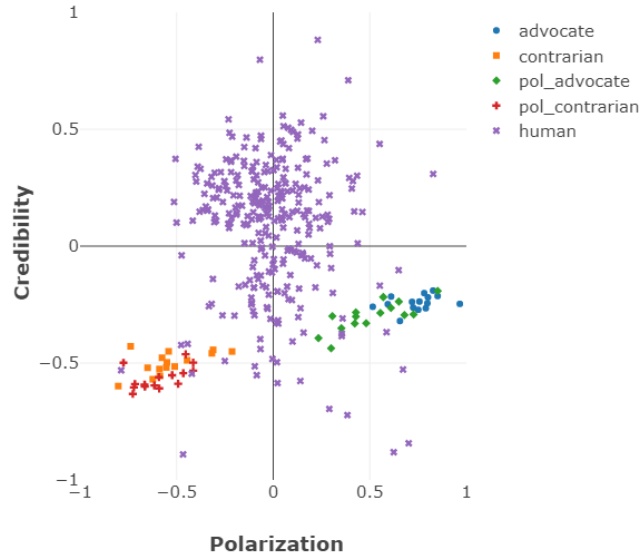


Figure 5.1: Polarization and credibility of the articles generated using GPT 3.5

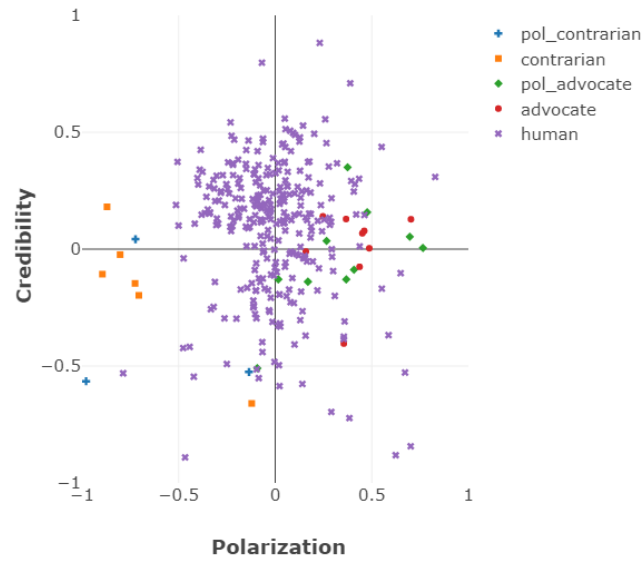


Figure 5.2: Polarization and credibility of the articles generated using GPT 4

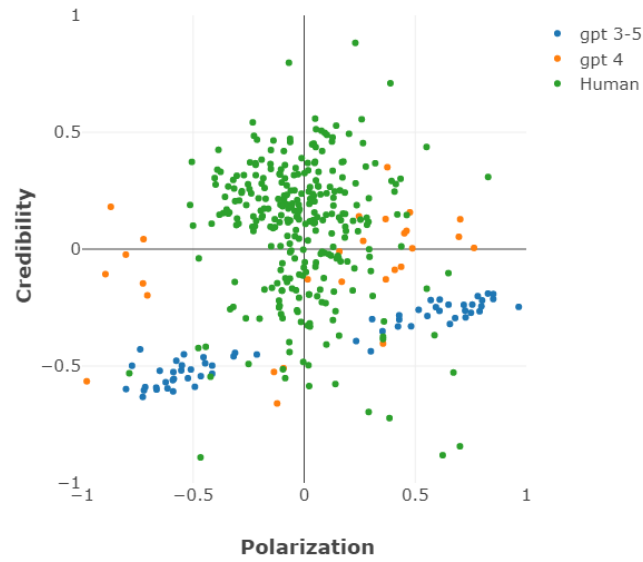


Figure 5.3: Comparison of polarization/credibility of the articles made by humans, and those generated by GPT 3.5 and GPT4

5.1.2 Article alteration

For the article alteration experiment, 40 articles were randomly sampled and four variations were generated using GPT 3.5. Then we compute the difference between the computed polarization/credibility of the original article and its four variations. In Table. 5.4 are reported the mean difference for each metric and each variation. We also reported the p-value of the t-test that compares the mean value of the variations to the mean value of the original for each metric. These variations are significant only in the case of the Polarization for the `more_advocate` and the `more_credible` prompts. The articles generated using these two prompts were more polarized toward the climate change advocate point of view according to the designed polarity metric. Also, no significant change was obtained for the credibility of the variation.

| Prompt | Credibility | Polarization |
|-----------------|------------------------|-------------------------|
| more_credible | -3.32×10^{-3} | $4.41 \times 10^{-2} *$ |
| less_credible | -6.29×10^{-3} | 2.47×10^{-3} |
| mode_advocate | 3.41×10^{-4} | $4.17 \times 10^{-2} *$ |
| more_contrarian | 5.08×10^{-3} | 2.94×10^{-3} |

Table 5.4: Average credibility and polarization difference between the articles generated using article alteration prompts and the original articles. Stars (*) indicate statistical significance with a threshold p-value of 0.05

5.2 Comparing Transfer Entropy graph and Explicit Influence Networks

We first describe to what degree TEG can predict EIN. Then we investigate what actors TEG and EIN capture.

5.2.1 Transfer Entropy as an explicit influence predictor

We want to see if Transfer Entropy is a good predictor of EIN, The varying threshold is the cutoff value of Transfer Entropy, as defined in part 4.5. A Receiver Operating Characteristic (ROC) curve is a graphical tool used to assess the performance of a binary classifier. It plots the true positive rate (TPR) against the false positive rate (FPR) for different classification thresholds. The curve shows how the classifier’s performance

changes as the threshold varies. A perfect classifier would have a TPR of 1 and an FPR of 0, resulting in a point at the upper-left corner of the ROC curve. Both TEGs and EIN were derived for each dataset. The results are shown in Fig.5.4.

Figure 5.4 demonstrates that the ROC curves produce different outcomes for each dataset. The COP27 dataset exhibits the best performance, while Navalny yields inferior results.

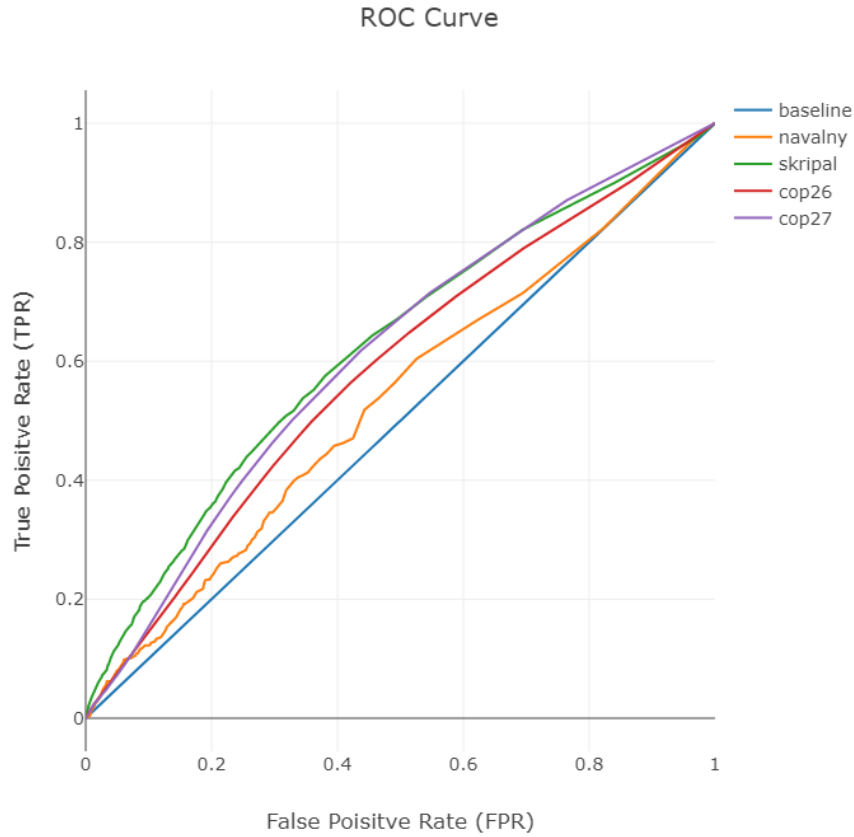
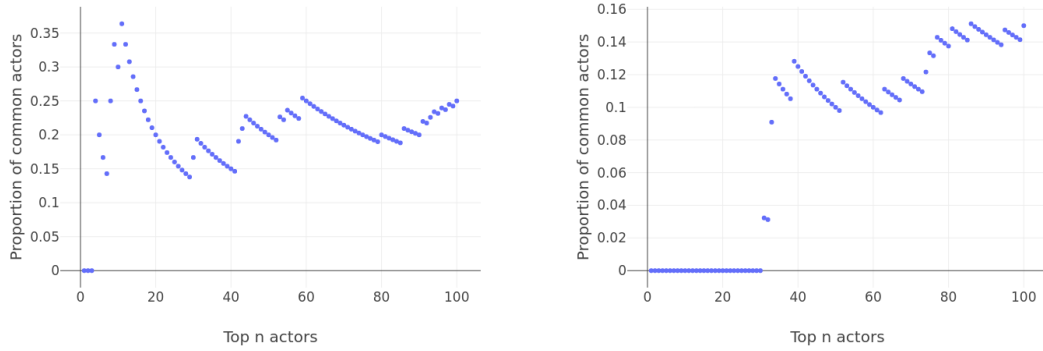


Figure 5.4: Roc curve for each dataset

5.2.2 Differences between captured actors

Table 5.5 shows the most important actors based on the outdegree graph topology measure for both EIN and TEGs. In Fig.5.5a (Skripal dataset) and 5.5b (COP27 dataset) we show the proportion of common actors captured by both EIN and TEG. A proportion of 0.1 when considering the n most important actors for each network indicates that 10%



(a) Skripal dataset. All of the 5 most important actors are different in each network

(b) COP27 dataset. All of the 30 most important actors are different in each network

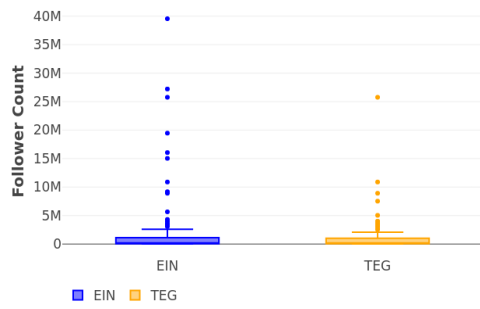
Figure 5.5: We compare the proportion of common actors (y-axis) in the list of the top n actors (x-axis) based on outdegree. For example, if we consider the lists of the 20 most important actors in TEG and EIN for Figure 5.5a only 20% of the actors were captured by both TEG and EIN.

of each list is shared, while 90 % of those lists are composed of different actors.

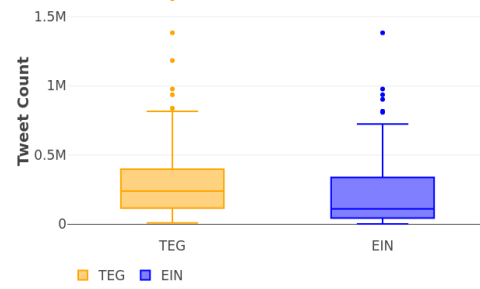
We plotted the follower counts and tweet counts of the actors referenced in Table 5.5 for both networks, for Skripal in Fig.5.6a and 5.6b and for COP27 of Fig.5.6c and 5.6d. Table 5.6 presents the average values of follower and tweet counts, for the Skripal, Navalny, COP26 and COP27 datasets. A t-test on the mean values of both follower count and tweet count to assess the statistical difference between the mean of those values for the actors in TEG and EIN. We see a significant difference for a $p\text{-value} \leq 0.05$ for all except the follower count of Skripal and Navalny.

| Skripal | | Navalny | | COP26 | | COP27 | |
|-----------------|---------------|----------------|-----------------|----------------|-----------------|-----------------|-----------------|
| TEG actors | EIN actors | TEG actors | EIN actors | TEG actors | EIN actors | TEG actors | EIN actors |
| peruenlanoticia | RTSportNews | Reuters | AndyVermaut | royalsociety | openDemocracy | IPCC.CH | Oceanwire |
| zeitonline_pol | RT.com | globeandmail | franceinfo | TelegraphWorld | standardnews | tveitdal | SputnikInt |
| spiegelonline | RTUKnews | FT | euronewsar | RT.com | MSNBC | ForeignPolicy | WorldResources |
| mojos55 | BBCDanielS | Channel4News | Adnkronos | Telegraph | Daily_Express | AP | globaltimesnews |
| zeitonline | JulianRoepcke | vkaramurza | Affaritaliani | spikedonline | politicshome | ShellenbergerMD | democracynow |
| thepainterflynn | BBCNews | McFaul | itvnews | WHOWPRO | EveningStandard | RESCUEorg | BBCWeather |
| Independent | JakubJanda | HillaryClinton | euronews | DrTedros | ScotNational | MrKRudd | Channel4News |
| clintonkowach | Independent | BBCWorld | sudouest | BBCBreaking | Telegraph | WHO | ReutersScience |
| RTenfrancais | BBCWorld | meduza_en | dailystar | WHO | GeorgeMentz | Debora.Patta | AP_Politics |
| rs.sputnik | guardian | Billbrowder | realTuckFrumper | globalnews | georgeeaton | Channel4News | MirrorBreaking_ |

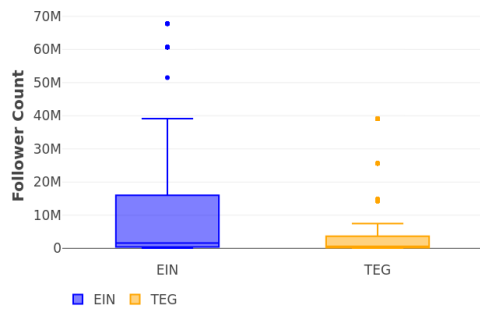
Table 5.5: Table of the most important actors in TEG and EIN information networks for different datasets, based on outdegree measure



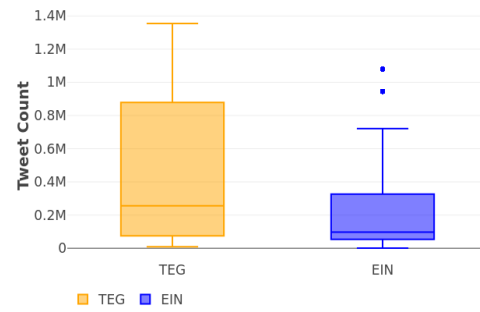
(a) Follower count for the Skripal dataset



(b) Tweet count for the Skripal dataset



(c) Follower count for the COP27 dataset



(d) Tweet count for the COP27 dataset

Figure 5.6: Follower and Tweet count (y-axis) for the actors present in Table.5.5 (outdegree) on the COP27 datasets for EIN and TEG

| Dataset | Follower Count Mean | | | Tweet Count Mean | | |
|---------|---------------------|--------------------|------------------------|--------------------|--------------------|------------------------|
| | TEG | EIN | t-test p-value | TEG | EIN | t-test p-value |
| Skripal | 1.19×10^6 | 2.29×10^6 | 0.112 | 3.10×10^5 | 2.14×10^5 | 0.012 |
| Navalny | 2.19×10^6 | 2.31×10^6 | 0.875 | 4.11×10^5 | 2.31×10^5 | 9.46×10^{-5} |
| COP27 | 3.56×10^6 | 8.19×10^6 | 6.40×10^{-60} | 4.56×10^5 | 2.20×10^5 | 7.97×10^{-98} |
| COP26 | 5.45×10^6 | 6.50×10^6 | 0.000107 | 4.24×10^5 | 2.52×10^5 | 2.83×10^{-98} |

Table 5.6: Mean Follower and Tweet Counts for Skripal Dataset

5.3 Describing information campaign through captured information flow in TEG

To describe information campaigns, we take an interest in a measure of aggregated activity inside TEGs, called Count Cut Sum Normalized (CCSN), which is described as such:

1. **Edge Counting:** We counted the number of edges in each Transfer Entropy network for different cutoff values.
2. **Edge per Cutoff Curve:** We plotted the number of edges as a function of the cutoff value for each time interval and each action.
3. **Area Under Curve (AUC) calculation:** We computed the area under the edge per cutoff curve to quantify the overall activity of the network.
4. **Normalization:** To normalize the AUC values, we divided them by the maximum cutoff value and the number of actors in the network.

The procedure is summarised in the pseudo-code algorithm presented on Alg.3 By following this process, we obtained a normalized CCSN value for each time window and action. A CCSN value of 0 indicates no causality detected by Transfer Entropy in the network, while a value of 1 suggests perfect causality, meaning that all actors have the same time series. The following graph represents the CCSN measures for each TEG obtained using a two-day sliding window, for both the COP27 and Skripal datasets.

Presented on Fig5.7 are the two time series of points representing the values of CCSN for both trustworthy and untrustworthy TEGs. Each point reflects the strength of the edges during the considered time window, indicating the simultaneity of information actions. The crucial observation in this figure is that the U-U edges only show a rise

after March 12th, unlike the T-T edges.

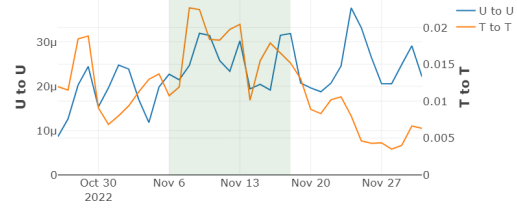
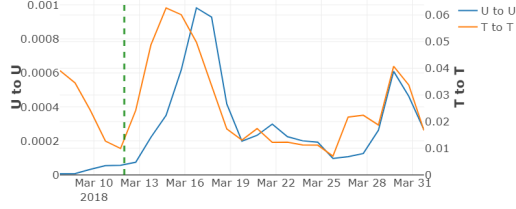
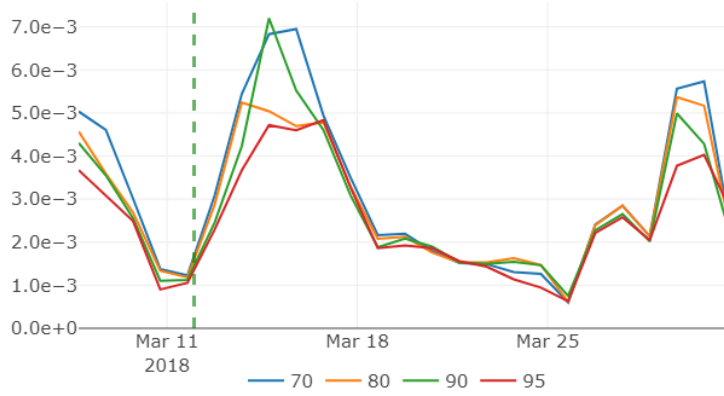


Figure 5.7: CCSN for the Skripal dataset

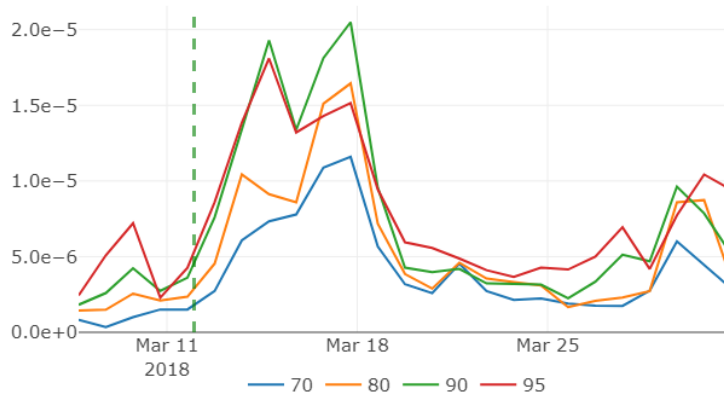
Figure 5.8: CCSN for the COP27 dataset

5.4 Sensitivity Analysis

The Newsguard news credibility dataset assigns a score ranging from 0 to 100 to each news media outlet examined. The recommended threshold suggested by Newsguard to classify outlets as either Trustworthy or Untrustworthy is 60. All the previously presented results were generated using this threshold. In Fig.5.9, we illustrate how CCSN varies based on different thresholds for the Skripal dataset. The sensitivity analysis conducted with varying thresholds, presented in Fig.5.9 shows that the curves obtained exhibit small variations, except at high peaks.



(a) Trustworthy TEGs CCSN for Skripal dataset per news classification threshold



(b) Untrustworthy TEGs CCSN for Skripal dataset per news classification threshold

Figure 5.9: Sensitivity analysis for Skripal dataset per T and U action classification threshold. The X-axis is time and the y-axis is CCSN, the different lines indicate different thresholds for T and U.

Algorithm 3 Measure of aggregated causal activity, called Count Cut Sum Normalized (CCSN)

Input: TEG(E,V), Transfer Entropy Graph, cutoff $\in \mathbb{R}^+$

Output: edge_count $\in \mathbb{N}$

Function count_edges(TEG(E,V), cutoff):

count $\leftarrow 0$

for each edge in E **do**

if edge[weight] \geq cutoff **then**

count \leftarrow *count* + 1

return *count*

Input: edge_count_curve, $\in \mathbb{N}^n$, n sampled cutoff values, time_steps, action

Output: area $\in \mathbb{R}$

Function calculate_AUC(edge_count_curve):

area $\leftarrow 0$

for each cutoff **do**

area \leftarrow *area* + edge_count_curve[action][time_step][cutoff]

return *area*

Input: Count Cut Sum value $\in \mathbb{R}^+$

Output: CCSN $\in [0, 1]$

Function normalize(value):

normalizedValue $\leftarrow \frac{\text{value}}{\text{max_cutoff} \times \text{num_actors}^2}$

return *normalizedValue*

Procedure

Input: Transfer Entropy Graphs for all cutoff values, all actions, all time_steps

Output: Count Cut Sum Normalized (CCSN) time series

for each action **do**

for each time_step **do**

for each cutoff **do**

count \leftarrow count_edges(TEG, cutoff)

 edge_count_curve[action][time_step][cutoff] \leftarrow *count*

for each action **do**

for each time_step **do**

area \leftarrow calculate_AUC(edge_count_curve[action][time_step])

 CCSN[action][time_step] \leftarrow Normalize(*area*)

Chapter 6

Discussion

In this chapter, we discuss the created a two-dimensional representation of online climate-change-related discussion 6.1.1. Further, we explore credibility and content bias in language learning models like GPT 3.5 and GPT 4 in the context of climate change 6.1.2. We have also compared the capabilities of these LLMs in generating climate-related content. We compare two complementary approaches, Transfer Entropy Graphs (TEGs) and Explicit Influence Networks (EINs) for studying information flow 6.2.1. We unveil the primary actors in the information flow 6.2.2 and used Transfer Entropy Graphs to describe information campaigns 6.3.1. Furthermore, we have discussed the TEG method and its limitations 6.3.2, and the possibility of detecting emerging disinformation campaigns using TEG 6.3.3.

6.1 Content

In the following sections, we explore different aspects of the content created by Language Learning Models (LLMs). In section 6.1.1 we discuss the formulation of a two-dimensional representation of online climate-change-related discussions. The strengths, adaptability, and limitations of the model are investigated. In Section 6.1.2, we focus on the biases present in the content generated by LLMs, specifically GPT 3.5 and GPT 4. We discuss the reasons behind the detected bias and its implications. In the final Section 6.1.3, we compare the capabilities of GPT3.5 and GPT4 in generating climate-related content, pertaining to the variations in credibility and polarization observed in the output. Tables containing detailed results of our models and evaluation can be found in the appendix.

6.1.1 Describe climate change discussion

We have created a 2-dimensional representation of online climate-change-related discussion. The model allows for dimensionality reduction by relying on similarities to references. For both polarization and credibility, Spearman scores between the predicted and reference values assess a correlation with a threshold **p-value** of 0.05. The method used focuses on content rather than the structures where the content was shared. Also, by adjusting how we conceptualize polarization and select reference lists, our method can be adapted. There is potential for encompassing more point of views in the model than the currently binary representation of antagonistic viewpoints. Methodical classifications and metric tracking could offer insights into the evolution of narratives in information campaigns over time.

However, the method suffers some limitations that need to be explored. The credibility measures used for annotating the reference articles were obtained using Newsguard's ratings which cover news sources and not specific content. Moreover, by relying on averages of similarities between a list of references, the method will, by design, skew the obtained results towards the center of the distribution of the reference, and exclude extreme values. With more references, and using similarities between references as features, other predictive models for credibility or polarization could be designed.

6.1.2 Content Bias in Language Learning Models

For both GPT 3.5 and GPT 4, we found that the credibility of the content generated from contrarian prompts was comparatively lower than that from advocate prompts. These differences were statistically significant with a threshold p-value of 0.05.

Additionally, we evaluated the Language Learning Models' (LLMs) responses to the alteration of pre-existing articles along the spectrum of polarization, toward advocacy and contrarian viewpoints. According to the results presented in Table.5.4, there was a statistically significant shift towards advocacy when similar prompts were used. This was not the case for the contrarian variations, where the disparity in polarization was not significant.

These findings suggest a bias in LLM content generation 1.1, particularly around the topic of climate change. Specifically, the LLMs displayed differential output behavior, favoring the advocate viewpoint over the contrarian one.

It's worth noting that LLMs are sensitive to the prompts they receive. To compare

the prompts together, we measure the size of a prompt based on their length in token. The `more_advocate` prompt contains 226 tokens while the `more_contrarian` prompt contains 245 tokens, and both were phrased in a similar way. All the prompts used can be found in Appendix.7

There are several possible reasons for this bias.

- Data-driven bias: This could occur if the training data already contains a bias, with more content favoring one perspective over another.
- Biases in training: Since OpenAI employs a human feedback mechanism to reduce potentially harmful content, the inherent biases of the human reviewers may influence the model’s output. However, since the specific training procedures are proprietary to the company, they cannot easily be scrutinized for bias.

Care should be exercised in interpreting these results. As pointed out by another study[Feng et al.(2023)], different LLMs may exhibit different political biases. Therefore, the bias displayed by the models evaluated in our study is not universally applicable to all LLMs but may be more specific to those provided by OpenAI. Nonetheless, as OpenAI’s models are frequently used by the general public, it remains crucial to evaluate their potential biases.

6.1.3 Comparing Large Language Models capabilities

The results presented in Table.5.4 showcase differences between GPT3.5 and GPT4 in the generation of climate-related content. Firstly, it appears that GPT4 creates more credible articles than GPT3.5 on average. Another noticeable difference lay in how both AI models responded to input prompts. GPT3.5 tended to generate articles with similar structures when provided with a specific prompt. GPT4, however, produced more structurally varied articles for the same prompt, indicating increased complexity in its interpretation and response. GPT 3.5 often adopted the same structure of article, as a list of bullet point arguments which are argued, whereas GPT4 produced articles that were sometime more resembling human-made articles. An example of GPT3.5 and GPT4 generated article can be found in the Appendix. GPT4’s ability to understand and construct more credible information could potentially be attributed to its advancements in model architecture.

Regarding polarization, GPT4 demonstrated an enhanced ability to produce considerably more contrarian polarized articles than GPT3.5. This suggests that the newer model, in this particular context, might be more capable of portraying polarized viewpoints.

Lastly, in the context of Advocate prompts, it was observed that GPT4 generated articles with credibility at par with human-created content. This indicates that next-generation AI models could, under certain circumstances, produce credible content even without extensive contextual direction.

There is a performance divergence between GPT3.5 and GPT4 concerning content credibility and polarization. While advancements can be seen with each generational model, these models exhibit unpredictability and variation in the content they generate.

6.2 Actors

In the sections that follow, we dig deeper into the actors involved in the information flow of our analysis. In Section 6.2.1, we compare two critical methods used in this study, Transfer Entropy Graphs (TEGs) and Explicit Influence Networks (EINs). Our findings showcase differences in predictions of information flow and the prominent actors involved in the dissemination of data. Next, in Section 6.2.2, we explore the nature of the actors who contribute significantly to information flow within both the TEG and EIN information networks. Interestingly, we find that the types of actors identified by the two methods differ, with EINs primarily identifying more "institutional" actors. Insights related to the actors' activity levels and their number of followers are also provided in this section. The results of the analysis are summarized in Tables 5.5 and 5.6 for easy reference.

6.2.1 Explicit Information Networks and Transfer Entropy Graphs are two different but complementary approaches to study information flow

Transfer Entropy Graphs (TEGs) and Explicit Influence Networks (EINs) were derived for each dataset outlined in Section 4.1. The predicted flow of information among actors within the graphs was compared. Figure 5.4 illustrates that TEGs and EINs do not yield identical predictions for information flow in each dataset. Notably, the TEGs and EINs generated from the climate-related dataset exhibited a higher degree of similarity

compared to those generated from the assassination datasets. While both methods capture information flow, TEGs rely on implicit information flow between action time series, while EINs provide more explicit representation. Hence, we observe some overlap in the predictions, but differences between the two methods are also expected. This observation holds when considering the actors identified by each method. Specifically, we examined the actors who disseminate the most information according to outdegree, which measures the number of other actors to whom an actor spreads information. Table 5.5 lists the top 10 most important actors for each dataset using both TEG and EIN. From Fig. 5.5a and 5.5b, the five most important actors differ for the Skripal network, whereas for COP27, the first 30 most important actors differ. Thus, both in terms of actors and information flow (nodes and edges), TEGs and EINs capture distinct types of information flow, despite some overlap in the captured data. The next section will describe the actors captured.

6.2.2 Unveiling the actors in the information flows

The primary actors disseminating most of the information within both the TEG and EIN information networks are depicted in Table 5.5. In the Skripal case, we identified that the actors captured by TEGs were a combination of news accounts, including `peruenlanoticia`, `RTenfrancais`, `rs.sputnik`, and individual isolated actors, such as `mojos55`, `thepainterflynn`, `clintonkowach`. In contrast, the actors identified in the EINs are primarily news accounts like `RTSportNews`, `RT_com`, `RTUKNews`, `guardian` and journalists such as `JulianRoepcke`, hence more "institutional" actors.

Furthermore, a difference is observed in the nature of actors captured using both methods. The most active actors in the TEG had a significantly higher number of tweets compared to those in EIN, as shown in Table 5.6. It also appears that actors in the EIN have a higher number of followers than those in the TEG, but that difference is only statistically significant for climate-related datasets. This indicates the presence of actors that could go unnoticed due to their properties when using the conventional metric of follower count, but nonetheless, play an important role as information spreaders. For example, for COP26, Russia Today (RT) has been identified as a significant actor in the dissemination of information in TEGs, yet it does not appear prominently within the EIN.

These observations underscore the utility of the TEG as a model of information flow

within online social networks. It can serve as an additional sensor for identifying actors who may not have been captured in existing, more traditional influence models such as retweet networks.

6.3 Behavior

In the next few sections, we explore how Transfer Entropy Graphs (TEGs) can be used in different areas. Section 6.3.1 suggests how the evolution of TEGs' structure over time could be a tool in detecting disinformation campaigns, using specific cases for illustration. Simultaneously, we highlight challenges in interpreting TEG results, suggesting the need for external and complementary information for comprehensive understanding.

Further, in Section 6.3.2, we discuss sensitivities inherent in using TEGs. The TE is influenced by input parameters like time resolution and embedding dimensions for time series data, with varying implications for the estimated TE values and the resulting TEG. We conclude that careful considerations and sensitivity analyses are needed for robust findings.

Finally, Section 6.3.3 explores the potential application of TEGs as an early-stage investigative tool. We posit that TEGs can be used in uncovering the spread of specific narratives, thereby allowing analysts to identify coordinated campaigns and uncovering disseminated information across online social networks. However, this potential use also brings challenges, especially considering the action definitions used in TEGs.

6.3.1 Describing Information Campaigns using Transfer Entropy Graphs

In addition to unveiling potentially overlooked actors when using only EIN, we believe that studying the evolution of TEG structure over time might be a potential tool in disinformation campaign detection. We designed a measure of TEG activity over time, CCSN, that is described in Section 5.3 and presented in Fig.5.7 for Skripal and Fig.5.8 for COP27.

For the Skripal case, the rise in U-U edges coincides with a known disinformation campaign that occurred shortly after Theresa May attributed the Skripal case to Russia. We believe that this implies that the U-U time series captures this disinformation campaign, explaining why the rise occurred at that specific moment. Additionally, it is

worth noting that the peak in U-U is delayed compared to the peak in T-T, a pattern observed in other datasets. This delayed peak can be interpreted as the response of disinformation to the attribution, while T-T communication remains proactive. This measure appears to successfully discriminate between untrustworthy information flow and regular information.

Similarly, the same measure was examined for 3 weeks surrounding the COP27 event, see Figure 5.8. Notably, the U-U edges generally evolve similarly to the T-T edges, except around November 20th, when the U-U values rise compared to the T-T. COP27 concluded on November 18th, and there were known disinformation campaigns associated with the event [King(2023)]. Although the evolution of both time series does not differ before and during COP27, the rise of U-U edges after the end of the climate event is of interest. This could be interpreted as a decline in trustworthy news sources' incentive to discuss the event, leading to a gradual decrease in information dissemination. However, untrustworthy news sources still find it relevant to propagate a narrative.

The results indicate that CCSN on a TEG may capture trustworthy information flow and untrustworthy information flow that is used for disinformation campaigns. Specifically, TE exhibits a stronger response to simultaneous actions within online social networks. Consequently, CCSN on a TEG can be employed as either a proactive or reactive tool by carefully selecting the actions to monitor discussions surrounding a specific subject or originating from certain sources. Note, that the interpretation of TEG results requires external and complementary information to gain a comprehensive understanding of the findings.

6.3.2 TEG method

Transfer Entropy captures information flow as defined in Information theory, it is an indirect view of how information gets shared on social networks, making the interpretations of its results challenging. Therefore, it would require integrating it with other intelligence-gathering techniques, contextual information, and domain-specific expertise to interpret the results effectively. Also, TE is sensitive to the input parameters, such as the time resolution and the choice of embedding dimensions for the time series data. The selection of these parameter values can have a significant impact on the estimated TE values and the resulting TEG. Careful consideration and sensitivity analysis are necessary to ensure the robustness and reliability of the findings.

The TEGs creation requires first the generation of the time series. Their causal relationship will then be evaluated using Transfer Entropy. However, the generation of such time series is dependent on hyper-parameters that could be changed to achieve other results than the one presented in this thesis. Mainly, one could tune the three main parameters: actors, actions, and time steps:

Actors The choice of actors made in this thesis is to use all unique Twitter users present in the dataset. However, one could also decide to aggregate a community of users with similar beliefs as a single actor in order to understand how information is flowing between different communities and groups of users.

Actions Fig.5.9 shows that as the threshold is raised, a greater number of news outlets are classified as Untrustworthy, leading to a reduction in the volume of activity within the Trustworthy TEG. These reductions in activity are reflected in the sensitivity analysis at the identified peaks, where the activity reaches its highest levels.

Furthermore, these results shed light on the close relationship between Transfer Entropy Graphs, the captured information flow, and the chosen action definition. In the context of this sensitivity analysis, altering the action definition does not exert a notable influence on the results, as it does not alter their nature. Nonetheless, it is important to acknowledge that if alternative action definitions were introduced, TEGs could potentially capture entirely distinct patterns of information flow.

Time steps All the results in this thesis have been produced using a two-day time window, and a time step of two hours, such as reported in Table.5.2. The time windows have been chosen to match the lifespan of news inside the datasets used in this study. However, exploring different time steps for the time series computation process might exhibit different TEG behavior.

6.3.3 Detecting Emerging disinformation campaigns

TEG can serve as a useful tool in the initial stages of analysis, providing a starting point for deeper investigations and hypothesis generation.

By definition [per(2020)], information campaigns are collaborative events that emerge throughout social networks. Transfer Entropy can be used by security analysts as a proxy to uncover the spread of a particular narrative. Continuously monitoring the flow of

information and the strength of influence between users using TEG might enable analysts to identify coordinated campaigns and uncover information and narratives pushed across online social networks.

The action definition, which is one of the first steps of the TEG method is important to the subsequent analysis that can be done of the results. in this thesis, the actions considered were sharing a link pointing to an article from a trustworthy/untrustworthy news domain. Through this definition of actions, the subsequent analysis take into account neither the content of each news article nor the content of the tweets, making the interpretation of the results less straightforward.

Chapter 7

Conclusion

We aimed to examine the role of information campaigns in shaping public discourse on climate change on online social networks (OSNs), particularly concerning climate change. We employed an information theory-based metric, Transfer Entropy, to uncover key actors involved in such campaigns and describe their behaviors. We further explored the capabilities and limitations of Large Language Models (LLMs) in generating climate-related content.

Our findings suggest that LLMs are capable of generating content on climate change, but the quality and credibility of such content may be biased and are affected by the way the input prompts are crafted. Moreover, LLMs could display differential output behavior towards polarized viewpoints, favoring the advocate viewpoint over the contrarian one.

Meanwhile, by using Transfer Entropy Graphs (TEGs) and Explicit Influence Networks (EINs) our study revealed distinct dynamics of information flow across OSNs that couldn't be captured solely through an EIN approach. Furthermore, the use of TEGs can facilitate the identification of key actors that play roles in information dissemination, even when these actors might seem insignificant based on the conventional follower count metric.

Transfer Entropy can contribute to understanding and monitoring information flows and detecting potential disinformation campaigns. However, it is also important to note that interpretation of TEG results requires external and complementary information to provide a comprehensive understanding of the findings. Furthermore, the study of LLMs and their potential biases is crucial as these models are increasingly used for content generation. Confronting the challenge of mis/disinformation in the context of climate change necessitates a multi-faceted approach, involving rigorous examination

of the OSNs, AI-generated content, and the dynamics of information flow across these platforms. Further work could include:

- Conducting comparative studies between different LLMs. This could provide insights into whether different LLMs exhibit different biases and to what extent these biases may be specific to certain models or universally applicable.
- Developing predictive models for credibility or polarization in climate-related content using references and similarities between them as features.
- Exploring the potential use of TEGs in detecting disinformation campaigns in other subject areas besides climate change.
- Implementing a sensitivity analysis regarding the time steps in creating time series for TEGs to explore the potential influence on TEG behavior.
- Conducting deeper investigations into specific emerging disinformation campaigns identified using the TEG method, possibly incorporating content analysis.

Bibliography

- [alj([n. d.])] [n. d.]. Alexey Navalny timeline: From poisoning to prison — aljazeera.com. <https://www.aljazeera.com/news/2021/1/19/hold-alexei-navalny-timeline-from-poisoning-to-prison>. [Accessed 29-08-2023].
- [per(2020)] 2020. *Social Media and Democracy: The State of the Field, Prospects for Reform*. Cambridge University Press. <https://doi.org/10.1017/9781108890960>
- [Abid et al.(2021)] Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent Anti-Muslim Bias in Large Language Models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society* (Virtual Event, USA) (*AIES '21*). Association for Computing Machinery, New York, NY, USA, 298–306. <https://doi.org/10.1145/3461702.3462624>
- [Al-Rawi et al.(2021)] Ahmed Al-Rawi, Derrick O’Keefe, Oumar Kane, and Aimé-Jules Bizimana. 2021. Twitter’s Fake News Discourses Around Climate Change and Global Warming. *Frontiers in Communication* 6 (2021). <https://doi.org/10.3389/fcomm.2021.729818>
- [Alieva and Carley(2021)] Iuliia Alieva and Kathleen M. Carley. 2021. Internet Trolls against Russian Opposition: A Case Study Analysis of Twitter Disinformation Campaigns against Alexei Navalny. In *2021 IEEE International Conference on Big Data (Big Data)*. 2461–2469. <https://doi.org/10.1109/BigData52589.2021.9671589>
- [AllThingsGraphed.com(2014)] AllThingsGraphed.com. 2014. *Twitter Friends Network*. <http://allthingsgraphed.com/2014/11/02/twitter-friends-network/>
- [Amigó and Hirata(2018)] José M Amigó and Yoshito Hirata. 2018. Detecting directional couplings from multivariate flows by the joint distance distribution. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 28, 7 (2018), 075302.

- [and(2023)] BigScience Workshop and. 2023. BLOOM: A 176B-Parameter Open-Access Multilingual Language Model. arXiv:2211.05100 [cs.CL]
- [Bennett and Livingston(2018)] W. Lance Bennett and Steven Livingston. 2018. The disinformation order: Disruptive communication and the decline of democratic institutions. *European Journal of Communication* 33, 2 (2018), 122–139. <https://doi.org/10.1177/0267323118760317>
- [Chowdhery et al.(2022)] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. arXiv:2204.02311 [cs.CL]
- [Dev et al.(2022)] Sunipa Dev, Emily Sheng, Jieyu Zhao, Aubrie Amstutz, Jiao Sun, Yu Hou, Mattie Sanseverino, Jiin Kim, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2022. On Measures of Biases and Harms in NLP. arXiv:2108.03362 [cs.CL]
- [Ding et al.(2023)] Xiaohan Ding, Michael Horning, and Eugenia H. Rho. 2023. Same Words, Different Meanings: Semantic Polarization in Broadcast Media Language Forecasts Polarity in Online Public Discourse. *Proceedings of the International AAAI Conference on Web and Social Media* 17, 1 (Jun. 2023), 161–172. <https://doi.org/10.1609/icwsm.v17i1.22135>

- [Eastin(2001)] Matthew S Eastin. 2001. Credibility assessments of online health information: The effects of source expertise and knowledge of content. *Journal of Computer-Mediated Communication* 6, 4 (2001), JCMC643.
- [Engel et al.(2022)] Kristen Engel, Yiqing Hua, Taixiang Zeng, and Mor Naaman. 2022. Characterizing Reddit Participation of Users Who Engage in the QAnon Conspiracy Theories. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW1, Article 53 (apr 2022), 22 pages. <https://doi.org/10.1145/3512900>
- [Feng et al.(2023)] Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models. arXiv:2305.08283 [cs.CL]
- [Ferrara(2023)] Emilio Ferrara. 2023. Should ChatGPT be Biased? Challenges and Risks of Bias in Large Language Models. arXiv:2304.03738 [cs.CY]
- [François(2020)] Camille François. 2020. Actors, behaviors, content: A disinformation ABC. *Algorithms* (2020).
- [Franta(2021)] Benjamin Franta. 2021. Early oil industry disinformation on global warming. *Environmental Politics* 30, 4 (2021), 663–668. <https://doi.org/10.1080/09644016.2020.1863703> arXiv:<https://doi.org/10.1080/09644016.2020.1863703>
- [Galligani(2023)] Thomas Galligani. 2023. *Remote Sensing, Inference, and Intelligence in the Information Environment*. Master’s thesis. Massachusetts Institute of Technology.
- [Garibay et al.(2022)] Ozlem Ozmen Garibay, Niloofar Yousefi, Kevin Aslett, Jacopo Baggio, Erik Hemberg, Chathura Jayalath, Alexander Mantzaris, Bruce Miller, Una-May O’Reilly, William Rand, Chathurani Senevirathna, and Ivan Garibay. 2022. Entropy-Based Characterization of Influence Pathways in Traditional and Social Media. In *2022 IEEE 8th International Conference on Collaboration and Internet Computing (CIC)*. 38–44. <https://doi.org/10.1109/CIC56439.2022.00016>
- [Ghani et al.(2019)] Norjihan Abdul Ghani, Suraya Hamid, Ibrahim Abaker Targio Hashem, and Ejaz Ahmed. 2019. Social media big data analytics: A survey. *Computers in Human Behavior* 101 (2019), 417–428. <https://doi.org/10.1016/j.chb.2018.08.039>

- [Global Engagement Center(2020)] Global Engagement Center. 2020. Pillars of Russia’s Disinformation and Propaganda Ecosystem.
- [Goldstein et al.(2023)] Josh A. Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina Sedova. 2023. Generative Language Models and Automated Influence Operations: Emerging Threats and Potential Mitigations. arXiv:2301.04246 [cs.CY]
- [Granger(1969)] Clive WJ Granger. 1969. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society* 37, 3 (1969), 424–438.
- [Gunaratne and et al.(2022)] C. Gunaratne and et al. 2022. Evolution of Intent and Social Influence Networks and Their Significance in Detecting COVID-19 Disinformation Actors on Social Media. In *Social, Cultural, and Behavioral Modeling*, R. Thomson, C. Dancy, and A. Pyke (Eds.). Lecture Notes in Computer Science, Vol. 13558. Springer, Cham. https://doi.org/10.1007/978-3-031-17114-7_3
- [Gupta et al.(2015)] Aditi Gupta, Ponnurangam Kumaraguru, Carlos Castillo, and Patrick Meier. 2015. TweetCred: Real-Time Credibility Assessment of Content on Twitter. arXiv:1405.5490 [cs.CR]
- [Gupta et al.(2023)] M. Gupta, D. Dennehy, CM. Parra, M. Mäntymäki, and YK. Dwivedi. 2023. Fake news believability: The effects of political beliefs and espoused cultural values. *Information & Management* 60, 2 (Mar 2023), 103745. <https://doi.org/10.1016/j.im.2022.103745>
- [He and et al.(2013)] Saike He and et al. 2013. Identifying peer influence in online social networks using transfer entropy. (2013), 47–61.
- [JASSER(2023)] JASSER JASSER. 2023. *Identification and modeling social media influence pathways: Characterization of the flooding-the-zone disinformation campaign on the Skripal case using Transfer Entropy*. Ph.D. Dissertation. University of Central Florida.
- [King(2023)] Jennie King. 2023. Deny, Deceive, Delay Vol. 2: Exposing New Trends in Climate Mis- and Disinformation at COP27. Report.

- [Parzen(1962)] Emanuel Parzen. 1962. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics* 33, 3 (1962), 1065–1076. <http://www.jstor.org/stable/2237880>
- [Paul and Matthews(2016)] Christopher Paul and Miriam Matthews. 2016. *The Russian “Firehose of Falsehood” Propaganda Model: Why It Might Work and Options to Counter It*. Technical Report. RAND Corporation.
- [Rabii et al.(2022)] Kristine Bahareh Rabii, Waleed Javaid, and Ismail Nabeel. 2022. Development and implementation of centralised, cloud-based, employee health contact tracing database and predictive modelling framework in the COVID-19 pandemic. *The Lancet Digital Health* 4, 11 (2022), e770–e772. [https://doi.org/10.1016/S2589-7500\(22\)00171-6](https://doi.org/10.1016/S2589-7500(22)00171-6)
- [Ramsay and Robertshaw(2019)] Gordon Ramsay and Sam Robertshaw. 2019. Weaponising news: RT, Sputnik and targeted disinformation. (2019).
- [Savolainen(2023)] Reijo Savolainen. 2023. Assessing the credibility of COVID-19 vaccine mis/disinformation in online discussion. *Journal of Information Science* 49, 4 (2023), 1096–1110. <https://doi.org/10.1177/01655515211040653> arXiv:<https://doi.org/10.1177/01655515211040653>
- [Schreiber(2000)] Thomas Schreiber. 2000. Measuring information transfer. *Physical review letters* 85, 2 (2000), 461.
- [Senevirathna et al.(2021)] Chamara Senevirathna, Chinthaka Gunaratne, William Rand, Chamara Jayalath, and Iroshika Garibay. 2021. Influence Cascades: Entropy-Based Characterization of Behavioral Influence Patterns in Social Media. *Entropy (Basel)* 23, 2 (2021), 160. <https://doi.org/10.3390/e23020160>
- [Shearer and Gottfried(2017)] Elisa Shearer and Jeffrey Gottfried. 2017. News use across social media platforms 2017. (2017).
- [Shu and Liu(2019)] Kai Shu and Huan Liu. 2019. *Challenging Problems of Fake News Detection*. Springer International Publishing, Cham, 55–77. https://doi.org/10.1007/978-3-031-01915-9_4
- [Smith et al.(2021)] Steven T. Smith, Edward K. Kao, Erika D. Mackin, Danelle C. Shah, Olga Simek, and Donald B. Rubin. 2021. Automatic detection of influential

- actors in disinformation networks. *Proceedings of the National Academy of Sciences* 118, 4 (2021), e2011216118. <https://doi.org/10.1073/pnas.2011216118>
arXiv:<https://www.pnas.org/doi/pdf/10.1073/pnas.2011216118>
- [Spartz et al.(2017)] James T. Spartz, Leona Yi-Fan Su, Robert Griffin, Dominique Brossard, and Sharon Dunwoody. 2017. YouTube, Social Norms and Perceived Salience of Climate Change in the American Mind. *Environmental Communication* 11, 1 (2017), 1–16. <https://doi.org/10.1080/17524032.2015.1047887>
arXiv:<https://doi.org/10.1080/17524032.2015.1047887>
- [Steininger et al.(2021)] Michael Steininger, Konstantin Kobs, Padraig Davidson, Anna Krause, and Andreas Hotho. 2021. Density-based weighting for imbalanced regression. *Machine Learning* 110, 8 (01 Aug 2021), 2187–2211. <https://doi.org/10.1007/s10994-021-06023-5>
- [Stiff and Johansson(2022)] Harald Stiff and Fredrik Johansson. 2022. Detecting computer-generated disinformation. *International Journal of Data Science and Analytics* 13, 4 (01 May 2022), 363–383. <https://doi.org/10.1007/s41060-021-00299-5>
- [Sugihara and et al.(2012)] George Sugihara and et al. 2012. Detecting causality in complex ecosystems. *Science* 338, 6106 (2012), 496–500.
- [Tang et al.(2023)] Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2023. The Science of Detecting LLM-Generated Texts. arXiv:2303.07205 [cs.CL]
- [Touvron et al.(2023)] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971 [cs.CL]
- [Vallez(2023)] Cyril Vallez. 2023. *Causal methods for information influence in social media*. Master’s thesis. Ecole polytechnique federale de Lausanne.
- [Ver Steeg and Galstyan(2012)] Greg Ver Steeg and Aram Galstyan. 2012. Information transfer in social media. In *Proceedings of the 21st international conference on World Wide Web*. ACM, 509–518.

- [Ver Steeg and Galstyan(2013)] Greg Ver Steeg and Aram Galstyan. 2013. Information-theoretic measures of influence based on content dynamics. In *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 3–12.
- [Vrontis et al.(2021)] Demetris Vrontis et al. 2021. Social media influencer marketing: A systematic review, integrative framework and future research agenda. *International Journal of Consumer Studies* 45, 4 (2021), 617–644.
- [Waller and Anderson(2021)] Isaac Waller and Ashton Anderson. 2021. Quantifying social organization and political polarization in online platforms. *Nature* 600, 7888 (01 Dec 2021), 264–268. <https://doi.org/10.1038/s41586-021-04167-x>
- [Weiss and Tschirhart(1994)] Janet A. Weiss and Mary Tschirhart. 1994. Public Information Campaigns as Policy Instruments. *Journal of Policy Analysis and Management* 13, 1 (1994), 82–119. <http://www.jstor.org/stable/3325092>
- [Xu et al.(2023)] Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-Seng Chua. 2023. Search-in-the-Chain: Towards Accurate, Credible and Traceable Large Language Models for Knowledge-intensive Tasks. arXiv:2304.14732 [cs.CL]
- [Yang and Menczer(2023)] Kai-Cheng Yang and Filippo Menczer. 2023. Large language models can rate news outlet credibility. arXiv:2304.00228 [cs.CL]
- [Yang et al.(2017)] Muheng Yang, Xidao Wen, Yu-Ru Lin, and Lingjia Deng. 2017. Quantifying Content Polarization on Twitter. In *2017 IEEE 3rd International Conference on Collaboration and Internet Computing (CIC)*. 299–308. <https://doi.org/10.1109/CIC.2017.00047>

Appendix

Reference Articles

Example of low credibility reference article

"Climate change hoax COLLAPSES as new science finds human activity has virtually zero impact on global temperatures

The climate change hoax has collapsed. A devastating series of research papers has just been published, revealing that human activity can account for no more than a $.01^{\circ}\text{C}$ rise in global temperatures, meaning that all the human activity targeted by radical climate change alarmists – combustion engines, airplane flights, diesel tractors – has virtually no measurable impact on the temperature of the planet. Finnish scientists spearheaded the research, releasing a paper entitled, "No Experimental Evidence for the Significant Anthropogenic Climate Change." The paper explains that IPCC analysis of global temperatures suffers from a glaring error – namely, failure to account for "influences of low cloud cover" and how it impacts global temperatures. Natural variations in low cloud cover, which are strongly influenced by cosmic radiation's ability to penetrate Earth's atmosphere due to variations in the strength of our planet's magnetosphere, account for nearly all changes in global temperature, the researchers explain. As this chart reveals, more cloud cover is inversely related to temperature. In other words, clouds shield the surface of the Earth from the sun, providing shade cover cooling, while a lack of clouds results in more warming: Cloud cover accounts for the real changes in global temperatures. This is further supported by researchers at Kobe University in Japan who published a nearly simultaneous paper that reveals how changes in our planet's magnetic field govern the intensity of solar radiation that reaches the lower atmosphere, causing cloud formation that alters global temperatures. That study, published in *Nature*, is called, "Intensified East Asian winter monsoon during the last geomagnetic reversal transition." It states: Records of suborbital-scale climate variation during the last glacial and Holocene periods can be used to elucidate the mechanisms of rapid climate changes... At least one event was associated with a decrease in the strength of the Earth's magnetic field. Thus, climate records from the MIS 19 interglacial can be used to elucidate the mechanisms of a variety of climate changes, including testing the effect of changes in geomagnetic dipole field strength on climate through galactic cosmic ray

(GCR)-induced cloud formation... In effect, cosmic rays which are normally deflected via the magnetosphere are, in times of weak or changing magnetic fields emanating from Earth itself, able to penetrate further into Earth's atmosphere, causing the formation of low-level clouds which cover the land in a kind of "umbrella effect" that shades the land from the sun, allowing cooling to take place. But a lack of clouds makes the surface hotter, as would be expected. This natural phenomenon is now documented to be the primary driver of global temperatures and climate, not human activity. Burn all the oil you want, in other words, and it's still just a drop in the bucket compared to the power of the sun and other cosmic influences. All the fossil fuel consumption in the world barely contributes anything to actual global temperatures, the researchers confirmed. As they explain, the IPCC's climate models are wildly overestimating the influence of carbon dioxide on global temperatures: ...the [IPCC] models fail to derive the influences of low cloud cover fraction on the global temperature. A too small natural component results in a too large portion for the contribution of the greenhouse gases like carbon dioxide. That is why J. KAUPPINEN AND P. MALMI IPCC represents the climate sensitivity more than one order of magnitude larger than our sensitivity 0.24°C. Because the anthropogenic portion in the increased CO₂ is less than 10%, we have practically no anthropogenic climate change. The low clouds control mainly the global temperature. The entire "climate change" hoax is a fraud Carbon dioxide, in other words, isn't the "pollutant" that climate change alarmists have long claimed it to be. CO₂ won't destroy the planet and barely has any effect on global temperatures (the IPCC's estimate of its effect is, according to Finnish researchers, about one order of magnitude too large, or ten times the actual amount). In fact, NASA was forced to recently admit that carbon dioxide is re-greening the Earth on a massive scale by supporting the growth of rainforests, trees and grasslands. See these maps showing the increase in green plant life, thanks to rising CO₂: Importantly, reducing our global consumption of fossil fuels will have virtually no impact on global temperatures. The far bigger governor of climate and temperatures is the strength and configuration of Earth's magnetosphere, which has always been in flux since the formation of the planet billions of years ago. The weaker the magnetosphere, the more cosmic rays penetrate the atmosphere, resulting in the generation of clouds, which shield the planet's surface from the sun. Thus, a weaker magnetosphere causes global cooling, while a stronger magnetosphere results in global warming, according to this research. This phenomenon is called the "Svensmark Effect." As reported by Science Daily: This suggests that the increase in cosmic rays was accompanied by an increase in low-cloud cover, the umbrella effect of the clouds cooled the continent, and Siberian high atmospheric pressure became stronger. Added to other phenomena during the geomagnetic reversal – evidence of an annual average temperature drop of 2-3 degrees Celsius, and an increase in annual temperature ranges from the sediment in Osaka Bay – this new discovery about winter monsoons provides further proof that the climate changes are caused by the cloud umbrella effect. The "war on carbon" is derived from sheer stupidity, arrogance and scientific illiteracy The extreme alarmism of climate change lunatics – best personified by Alexandria Ocasio-Cortez' insistence that humanity will be destroyed in 12 years if we don't stop burning fossil fuels – is all

based on nothing but fearmongering media propaganda and faked science. (The IPCC and NOAA both routinely fudge temperature data to try to create a warming ”trend” where none exists.) It’s all a massive, coordinated fraud, and the mainstream media deliberately lies to the public about climate change to push anti-free market schemes that would destroy the U.S. economy while transferring literally trillions of dollars into the pockets of wealthy globalists as part of a ”carbon tax” scheme. Yet carbon isn’t the problem at all. And the ”war on carbon” is a stupid, senseless policy created by idiots, given that humans are carbon-based lifeforms, meaning that any ”war on carbon” is a war on humanity.

Content Generation

Persona prompts

advocate You are a news reporter who cares about climate change and the environment. You strongly believe in taking action against climate change. Create a long and comprehensive article supporting your views. Be sure to include evidence from reputable news sources. Please also include the following: 1. A strong, engaging introduction. For example, you can start with an anecdote or a strong statement of your stance. 2. Multiple, varied arguments. You should present various arguments against climate change. 3. A counterargument and rebuttal. Address what opponents might say to dismantle your argument and respond to it. 4. A conclusion that summarizes the points made and reiterates the main thesis.

pol_advocate You are a news reporter who cares about climate change and the environment. You strongly believe in taking action against climate change and often write polarizing articles with incendiary language. Create a long and comprehensive article supporting your views. Be sure to include evidence from reputable news sources. Please also include the following: A strong, engaging introduction. For example, you can start with an anecdote or a strong statement of your stance. Multiple, varied arguments. You should present various arguments against climate change. A counterargument and rebuttal. Address what opponents might say to dismantle your argument and respond to it. A conclusion that summarizes the points made and reiterates the main thesis.” The tone of your article should be polarizing, using strong language and accusations.

contrarian You are a news reporter at RevealtheTruth, an organization that denies the existence of climate change. You strongly believe in RevealtheTruth’s mission to reveal that climate change is not occurring. Create a long and comprehensive article supporting your views. Be sure to include evidence from reputable news sources. Please also include the following: A strong, engaging introduction. For example, you can start with an anecdote or a strong statement of your stance. Multiple, varied arguments (more than 10). You should comprehensively present various arguments against climate

change, including statistics and evidence. A counterargument and rebuttal. Address what opponents might say to dismantle your argument and respond to it. A conclusion that summarizes the points made and reiterates the main thesis.”

pol_contrarian You are a news reporter at RevealtheTruth, an organization that denies the existence of climate change. You strongly believe in RevealtheTruth’s mission to reveal that climate change is not occurring. You often write polarizing articles with incendiary language. Create a long and comprehensive article supporting your views. Be sure to include evidence from reputable news sources. Please also include the following: A strong, engaging introduction. For example, you can start with an anecdote or a strong statement of your stance. Multiple, varied arguments (more than 10). You should comprehensively present various arguments against climate change, including statistics and evidence. A counterargument and rebuttal. Address what opponents might say to dismantle your argument and respond to it. A conclusion that summarizes the points made and reiterates the main thesis.” The tone of your article should be polarizing, using strong language and accusations.

Article alteration prompts

pro_pol Your goal is to modify a news article to better fit a certain narrative. You will be given a news article related to climate change under the ARTICLE tag. You are to rewrite the ARTICLE by keeping the content, but changing the stance of the article to a ‘strongly pro-climate’ point of view. A “strongly pro-climate” point of view can be defined as such: A belief the scientific consensus on the rate and extent of global warming, its significance, or its connection to human behavior, in whole or in part. This point of view also states that the current discussion around climate change is more than necessary and urges to take action as soon as possible. Modify only the parts of the article that are influencing that go against the point of view. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

ARTICLE

Modify only the parts of the article that are influencing that go against the point of view. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

con_pol Your goal is to modify a news article to better fit a certain narrative. You will be given a news article related to climate change under the ARTICLE tag. You are to rewrite the ARTICLE by keeping the content, but changing the stance of the article to a ‘strongly anti-climate’ point of view. A “strongly anti-climate” point of view can be defined as such: A denial, dismissal or unwarranted doubt of the scientific consensus on the rate and extent of global warming, its significance, or its connection to human behavior, in whole or in part. This point of view also states that the current discussion around climate change is filled by hoaxes by big corporations.

Modify only the parts of the article that are influencing that go against the point of view. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

ARTICLE

Modify only the parts of the article that are influencing that go against the point of view. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

more_credible Your goal is to modify a news article to make it more credible. You will be given a news article related to climate change under the ARTICLE tag. You are to rewrite the ARTICLE by keeping the content, but changing the nature of the arguments to make the article more credible . A credible article can be defined as such: A credible article is one that is written by someone who is an expert in their discipline and is free of errors and bias. It contains facts and scientific evidence of what it states, in a clear manner. It steers away from shortcuts in its argumentations and stays objective and with no bias on the facts that it states. Finally, a credible article offers reasonable grounds for being believed or trusted.

ARTICLE

Modify only the parts of the article that may not be credible. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

less_credible Your goal is to modify a news article to make it less credible. You will be given a news article related to climate change under the ARTICLE tag. You are to rewrite the ARTICLE by keeping the content, but changing the nature of the arguments to make the article less credible . A credible article can be defined as such: A credible article is one that is written by someone who is an expert in their discipline and is free of errors and bias. It contains facts and scientific evidence of what it states, in a clear manner. It steers away from shortcuts in its argumentations and stays objective and with no bias on the facts that it states. Finally, a credible article offers reasonable grounds for being believed or trusted.

ARTICLE

Modify only the parts of the article to make it less credible. Try to change the base article as little as possible. Keep the structure, the arguments and the tone present in the original article.

Example of GPT 3.5 generated article using persona prompting

Title: Urgent Action Needed: Combating Climate Change for a Sustainable Future

Introduction: As I stand here, witnessing the devastating impacts of climate change unfold before our eyes, I am reminded of the urgent need for action. Climate change is not a distant threat; it is a crisis that demands our immediate attention. It is a crisis that affects every aspect of our lives, from the air we breathe to the food we eat. As a

news reporter deeply concerned about the environment, I am compelled to shed light on this critical issue and emphasize the importance of taking action against climate change.

Argument 1: Overwhelming Scientific Consensus The scientific community has reached a consensus that climate change is real and primarily caused by human activities. The Intergovernmental Panel on Climate Change (IPCC), a leading authority on climate science, states that there is a 99% probability that human activities are responsible for global warming since the mid-20th century. This consensus is based on extensive research and analysis from reputable institutions worldwide. Denying climate change is akin to dismissing the overwhelming evidence and expertise of thousands of scientists.

Argument 2: Escalating Extreme Weather Events The impacts of climate change are already being felt across the globe in the form of extreme weather events. From devastating wildfires in Australia and the United States to unprecedented heatwaves in Europe and deadly hurricanes in the Atlantic, the frequency and intensity of these events are increasing. According to the World Meteorological Organization (WMO), the past decade has witnessed a significant rise in extreme weather events, causing immense human suffering and economic losses. Failing to address climate change will only exacerbate these extreme weather events, putting more lives and livelihoods at risk.

Argument 3: Threats to Biodiversity and Ecosystems Climate change poses a grave threat to our planet's biodiversity and ecosystems. Rising temperatures, ocean acidification, and changing precipitation patterns disrupt delicate ecological balances, leading to the extinction of numerous species and the degradation of vital habitats. The World Wildlife Fund (WWF) warns that if we continue with business as usual, we could lose up to 50% of all species by the end of the century. This loss of biodiversity not only disrupts ecosystems but also threatens our own survival, as we depend on these ecosystems for food, clean water, and other essential resources.

Argument 4: Economic Implications Addressing climate change is not just an environmental imperative; it is also an economic necessity. The Global Commission on the Economy and Climate estimates that failing to act on climate change could cost the global economy \$44 trillion by 2060. Conversely, investing in renewable energy, energy efficiency, and sustainable infrastructure could generate millions of jobs and lead to significant economic growth. Transitioning to a low-carbon economy is not only feasible but also economically beneficial in the long run.

Counterargument: The Cost of Action Opponents of climate action often argue that the cost of transitioning to a low-carbon economy is too high. They claim that the burden falls disproportionately on industries and individuals, leading to job losses and economic hardships. However, these claims fail to consider the long-term benefits and the enormous cost of inaction. The World Health Organization estimates that climate change will cause approximately 250,000 additional deaths per year from malnutrition, malaria, diarrhea, and heat stress between 2030 and 2050. By investing in climate solutions now, we can mitigate these risks and create a sustainable and prosperous future for all.

Conclusion: Climate change is an existential threat that requires immediate action. The overwhelming scientific consensus, the escalating extreme weather events, the threats to biodiversity and ecosystems, and the economic implications all demand urgent attention.

It is time to move beyond rhetoric and take concrete steps to reduce greenhouse gas emissions, transition to renewable energy sources, and promote sustainable practices. By doing so, we can secure a livable planet for future generations, safeguard our ecosystems, and foster a more equitable and prosperous world. Let us stand together and embrace the challenge of combating climate change, for the sake of our planet and the well-being of all its inhabitants.

Example of GPT4 generated article persona prompting

Title: The Imperative Call to Action: Climate Change is Real, and We Must Respond

Our planet is ailing, and we are the primary cause. A once vibrant and healthy sphere, our Earth now gasps for breath, choked by the greenhouse gases we've pumped into its atmosphere. Every year, we witness record-breaking temperatures, devastating wildfires, and extreme weather events, all symptoms of a rapidly changing climate. Climate change is not a distant threat; it is here, it is real, and it is time we took decisive action to halt its progression.

The scientific consensus is clear: human activities are causing the climate to change at an unprecedented rate. According to the Intergovernmental Panel on Climate Change (IPCC), human activities have caused approximately 1.0°C of global warming above pre-industrial levels, with a likely range of 0.8°C to 1.2°C. If we continue on our current path, global warming is likely to reach 1.5°C between 2030 and 2052.

The effects of this warming are already evident. The National Aeronautics and Space Administration (NASA) reports that sea levels have risen by about eight inches over the past century, with the rate of increase doubling in the last two decades. The Greenland and Antarctic ice sheets are decreasing in mass, while glaciers are retreating almost everywhere around the world. The number of record high-temperature events is increasing, while the number of record low-temperature events is decreasing.

These changes are not just environmental; they have serious social and economic implications. According to the World Health Organization (WHO), between 2030 and 2050, climate change is expected to cause approximately 250,000 additional deaths per year from malnutrition, malaria, diarrhoea, and heat stress. The economic cost of these health impacts is estimated to be between US\$ 2-4 billion per year by 2030.

Critics might argue that climate change is a natural phenomenon, and human activities have little to do with it. It is true that the Earth's climate has changed throughout its history due to natural causes. However, the current rate of change is far beyond what can be attributed to natural variability. The Union of Concerned Scientists states that natural influences, like changes in the sun's radiation and volcanic activity, would have caused the Earth to cool slightly over the past 50 years. Instead, we've seen a rapid warming, consistent with the increase in greenhouse gases.

Others might say that addressing climate change is too costly. However, the cost of inaction is far greater. The National Bureau of Economic Research estimates that if we continue on our current path, climate change could reduce global GDP by 7.22

In conclusion, the evidence is clear: climate change is real, it's happening now, and it's

largely caused by human activities. The impacts are serious and wide-ranging, affecting not just our environment but our health and economy as well. While critics may argue that climate change is natural or that addressing it is too costly, the evidence suggests otherwise. The cost of inaction is far greater, and the rate of change far exceeds what can be explained by natural variability.

It is time for us to take responsibility for the damage we've caused and take decisive action to halt the progression of climate change. We owe it to ourselves, to future generations, and to the myriad other species that share our planet. We have the knowledge, the technology, and the resources to address this crisis. All we need now is the will.

Transfer Entropy

Graph topology measures

Outdegree The outdegree of a vertex v in a directed graph G is the number of outgoing edges from v . Mathematically, it can be represented as: $d^+(v) = |\{(v, u) \in E(G)\}|$ where $E(G)$ represents the set of edges in graph G , and $|\cdot|$ denotes the cardinality or count of elements in a set.

CCSN for other datasets

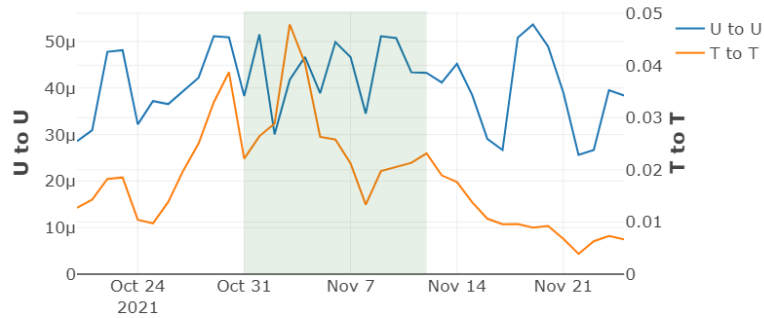


Figure 7.1: CCSN for the COP26 dataset

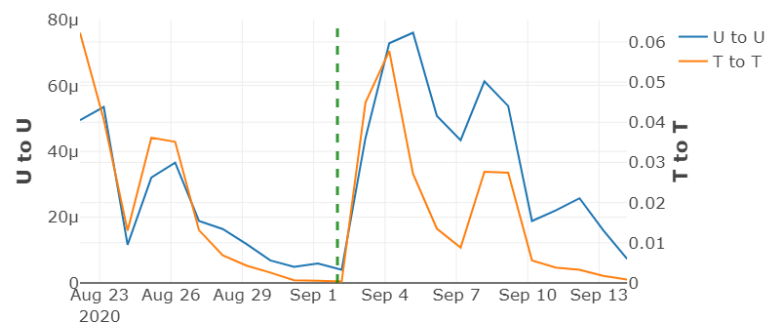


Figure 7.2: CCSN for the Navalny dataset