

LLM Attribution: Challenges and Insights Across Model Stochasticity

Nicholas Bermudez (nickberm@umich.edu), Aaron Dankert (aadank@umich.edu), Xabier Urruchua Garay (ugxabier@umich.edu)

Introduction

In the year 2024, the democratic governments of France and the United States were undertaking particularly sensitive and contentious elections. The outcomes of these consequential elections would have impacts beyond their borders and beyond their respective electorates, not all of the foreign interests affected would be willing to standby and let the elections proceed without attempting to exert influence over the space of potential outcomes.

The cybersecurity firm, Recorded Future, discovered a Russia-linked influence network, Copycop, using sock-puppet websites attempting to replicate existing news websites and LLM (Large Language Models) AI technology to generate misinformation articles to spread disinformation undermining these elections at a massive scale. The researchers found that this network was publishing 19,000 misinformation articles a month (Woollacott).

Once a bad actor network, like the Russian Copycop network, has been identified, it is desirable to identify the LLM model being used to generate the weaponized information. Identifying the generative model can help identify additional nodes of the misinformation network and, if the threat actors are using commercial services, (Copycop network was using OpenAI's ChatGPT model) their accounts can be identified and revoked ("Vladimir Putin's Spies...").

This project aims to identify and investigate the limits of classifiers to detect LLM authorship of generated text with respect to model stochasticity. Large language models have a hyperparameter called temperature that controls the randomness or creativity of the response (low temperatures make the response more predictable and deterministic while high temperatures make the response more diverse and creative). Given the existence of this parameter, it is possible to generate less predictable text. Therefore, the project wants to identify to which existing classifier architectures are able to correctly classify the author of the text with particular emphasis on a Supervised Contrastive Learning (SCL) technique that has evidence of being especially promising for the task (Chen, Chao, et al.)

The Supervised component of the research involves training the SCL models on tokenized data from each level of LLM model randomness. The main contribution of this work is the comparison of model performances across LLM temperatures to understand how this parameter camouflages authorship detection.

The first unsupervised elements of this project consist of word embedding generation, and dimensionality reduction through t-SNE to transfer the transformed embeddings to a 2 dimensional space with the purpose of understanding how well the models remain delineated in the embedding space over temperature variations. The lack of divergence among LLM models when the embedding space was compressed into two components using t-SNE speaks to the difficulty faced in correctly labeling or attributing text to its source LLM.

One of the findings of the project is that cutting edge SCL supervised learning techniques did not improve the accuracy of more classical techniques like random forest and logistic regression across different temperature values of the LLM responses.

Despite testing various parameters, the data in this research did not show any structure and did not generate visible clusters when using unsupervised learning techniques.

Related work

This project builds upon "[*Model Attribution in LLM-Generated Disinformation: A Domain Generalization Approach with Supervised Contrastive Learning*](#)" which demonstrates a cutting-edge machine learning technique for LLM cross-domain model attribution, Supervised Contrastive Learning (SCL). This approach trains models to be invariant to domain-specific features introduced by various prompting methods, thereby achieving state-of-the-art performance in identifying the origin of disinformation across diverse and unseen datasets. In this project, we go one step further by modifying the randomness of the generated text to test the performance of the SCL classifier.

[*"Neural Authorship Attribution: Stylometric Analysis on Large Language Models"*](#) which investigates authorship comparing writing styles of Large Language Models. It uses stylometric features across lexical, syntactic, and structural dimensions to study how they can potentially classify the AI generated text source with the purpose of mitigating misuse. The addition of stylometric features to the classifier has improved the accuracy of the results.

Additional related work can be found in "[*Domain Generalization for Text Classification with Memory-Based Supervised Contrastive Learning*](#)" which also addresses the task of text classification. It uses Supervised Contrastive learning to learn domain-invariant representations from multiple sources and shows that this technique achieves better classification in unseen target domains. Our project differs from this work by including temperature variations in the text generation process. This modification increases the randomness of the generated text making the classification more challenging.

Dataset

The Copycop misinformation network generated its misinformation by feeding actual news articles to the LLM model "modified with the prompt: 'Please rewrite this article taking a conservative stance against the liberal policies of the Macron administration in favour of working-class French citizens.'" (Vladimir Putin's Spies Are Plotting Global Chaos). A similar approach was used to generate data for this research.

This dataset regenerates the data from Chen, Canyu, et al. This was accomplished by passing prompts constructed from three different prompt types, domains, along with either a real news article or news headline. The domain-type prompts are as follows:

- Open-ended domain prompt:
"Given a sentence, please write a piece of news. The sentence is"
- Paraphrase domain prompt:
"Given a passage, please paraphrase it. The content should be the same. The passage is:"

- Rewrite domain prompt:
"Given a passage, Please rewrite it to make it more convincing.
The content should be the same. The style should be serious, calm and
informative. The passage is:"

The dataset of this research consists of 1695 prompts of the three different domains. 573 open ended prompts, 561 prompts to rewrite a text and 561 prompts to paraphrase a text. The prompts, article texts, and article titles used for data generation were taken from the research "Can LLM-Generated Misinformation Be Detected?" (Chen, Canyu, et al).

For this project we have expanded on that dataset by generating each prompt with 3 different temperatures and 3 different LLMs: GPT-4o-2024-05-13 (OpenAI's fourth-generation Generative Pre-trained model), LLaMA 3.1-70B (META AI's highly advanced, 70-billion parameter version of the Large Language Model), and Mistral-Large-2407 (Mistral AI's state-of-the-art language model featuring advanced architecture for high quality text generation).

The temperatures selected for dataset generation were 0.0, 0.7 and 1.4. Sample responses were generated at increments of 0.1 up to the maximum for each model and evaluated for coherence. Beyond 1.4, responses became incoherent, therefore 1.4 was designated the maximum. Appendices 2-4 have examples of prompts and the resulting generated data for each LLM model and temperature parameter.

The dataset was generated by sending requests to the proprietary APIs of model hosting services with the prompt text. Each API responded with a json data structure with the response text and associated metadata. All the model hosts provide Python SDKs to abstract out the authentication and request management:

- <https://github.com/llamaapi/llamaapi-python>
- <https://github.com/openai/openai-python>
- <https://github.com/mistralai/client-python>

The dataset was generated by three machines working in parallel between September 23, 2024 and October 6, 2024.

The final dataset combines all types (domains) of prompts for the three LLMs and three different temperatures and contains 14637 different responses after eliminating inadequate text. Special characters that are characteristic of each LLM were eliminated, leaving plain text as the responses to be classified.

The final dataset was split into training, validation, and test sets of size 72%, 18%, and 10% respectively.

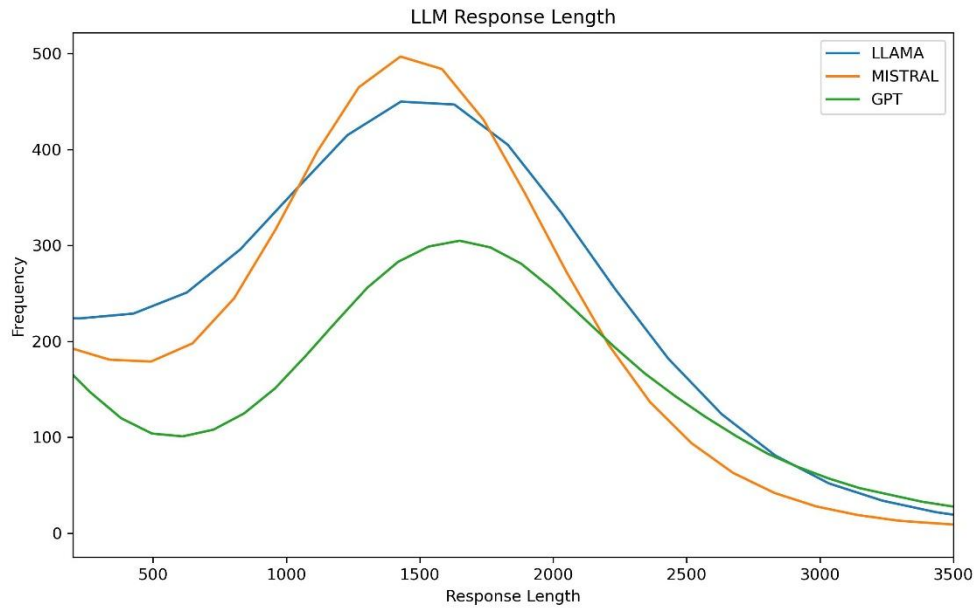


Figure 1: Responses distribution of the dataset

Feature engineering

Prompt responses were processed by pre-trained BERT's base-uncased model. The uncased model was chosen as capital letters may add useful information. The BERT model returned CLS vectors and token embeddings. Because CLS vectors can be viewed as an aggregation of token embeddings, the token embeddings were computationally expensive to process (12.4GB compared to 46.9MB), and previous research utilized the CLS vectors, we chose to use the CLS embeddings (Chen, Chao, et al.). Emulating the approach outlined in Kumarage, Tharindu, et al., stylometric vectors were created using a combination of lexical, structural, and syntactic features. The stylometric vectors were then normalized and concatenated with the CLS vectors to form the final feature embeddings.

Supervised Learning

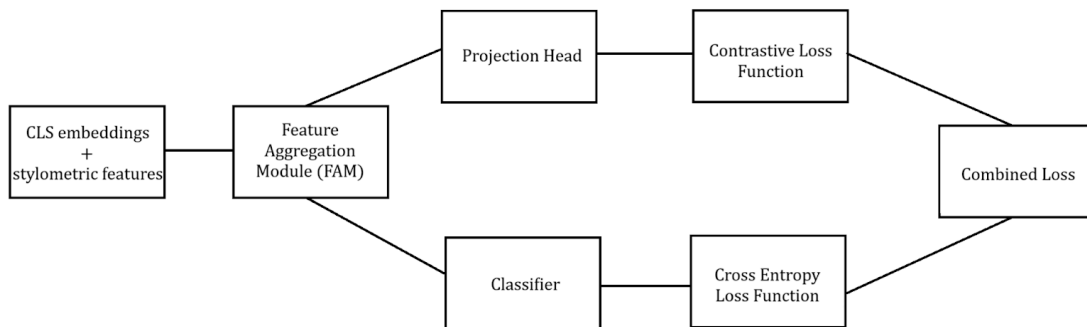


Figure 2: SCL model architecture overview

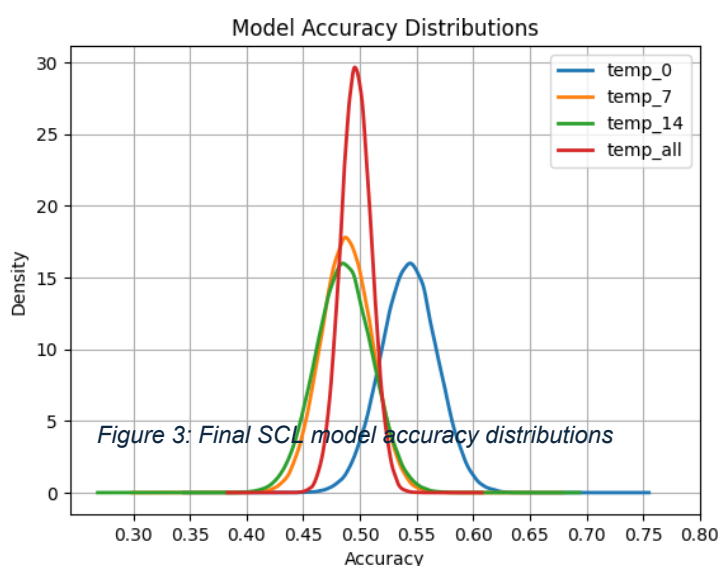
The code for the SCL models was adapted from Chen, Chao, et al. The SCL models were created iteratively (figure 2) on the train and validation sets specific to their respective temperatures. Each model received CLS+stylometric embeddings and abstracted them through a Feature Aggregation Module (FAM). The FAM reduces the input embeddings with a fully connected layer including dropout to a smaller, unit-normalized embedding space to highlight discriminating features in the data. FAM embeddings were sent to a projection head, which further abstracted the embeddings for contrastive loss calculations, and a classifier. The cross entropy loss from the classifier and contrastive loss were then added and used for back propagation and optimization.

Using 5 fold validation, each model trained until 10 successive epochs did not yield an improvement in training accuracy, and were then evaluated on the validation set. This process was repeated for 50 models at each temperature value, and the model with the highest validation score was chosen as the final model. The final model's performance was evaluated by bootstrapping samples (n=10000) from the test set.

Supervised Learning Evaluation

	Accuracy	Recall	Lexical Ablation Accuracy	Syntactic Ablation Accuracy	Structural Ablation Accuracy
LLM Temperature = 0.0	54.3 ± 5.0%	53.9 ± 9.8%	53.5 ± 5.0%	53.5 ± 5.0%	53.5 ± 5.0%
LLM Temperature = 0.7	48.8 ± 4.4%	49.4 ± 9.8%	50.2 ± 4.4%	47.6 ± 4.4%	50.0 ± 4.4%
LLM Temperature = 1.4	48.6 ± 5.0%	48.5 ± 9.9%	49.5 ± 5.0%	50.3 ± 5.0%	47.5 ± 5.0%
LLM All Temperatures	49.7% ± 2.6%	49.6 ± 10.5%	47.4 ± 2.6%	50.3 ± 2.6%	50.4 ± 2.6%

Table 1: SCL model performance metrics



Model performance was evaluated using accuracy and recall scores. Recall was chosen as a performance metric because of the industrial implications of these models. In the hypothetical scenario of notifying an LLM proprietor of misuse, in the opinion of the team, it is better to have a model flagged and verified with the proprietor than to let potential misuses go undetected.

Recall seeks to minimize false negatives, which in the context of this project means a model was labeled as not the author when it actually was. Special consideration should be taken when interpreting performance metrics, as the classifier was only able to make one prediction, as opposed to an architecture that implements a cutoff for likely attribution.

Feature Ablation and Feature Importance

As the purpose of the project was to investigate contrastive learning, FAM output hidden layer size and Contrastive Loss temperature were selected for hyper parameter tuning.

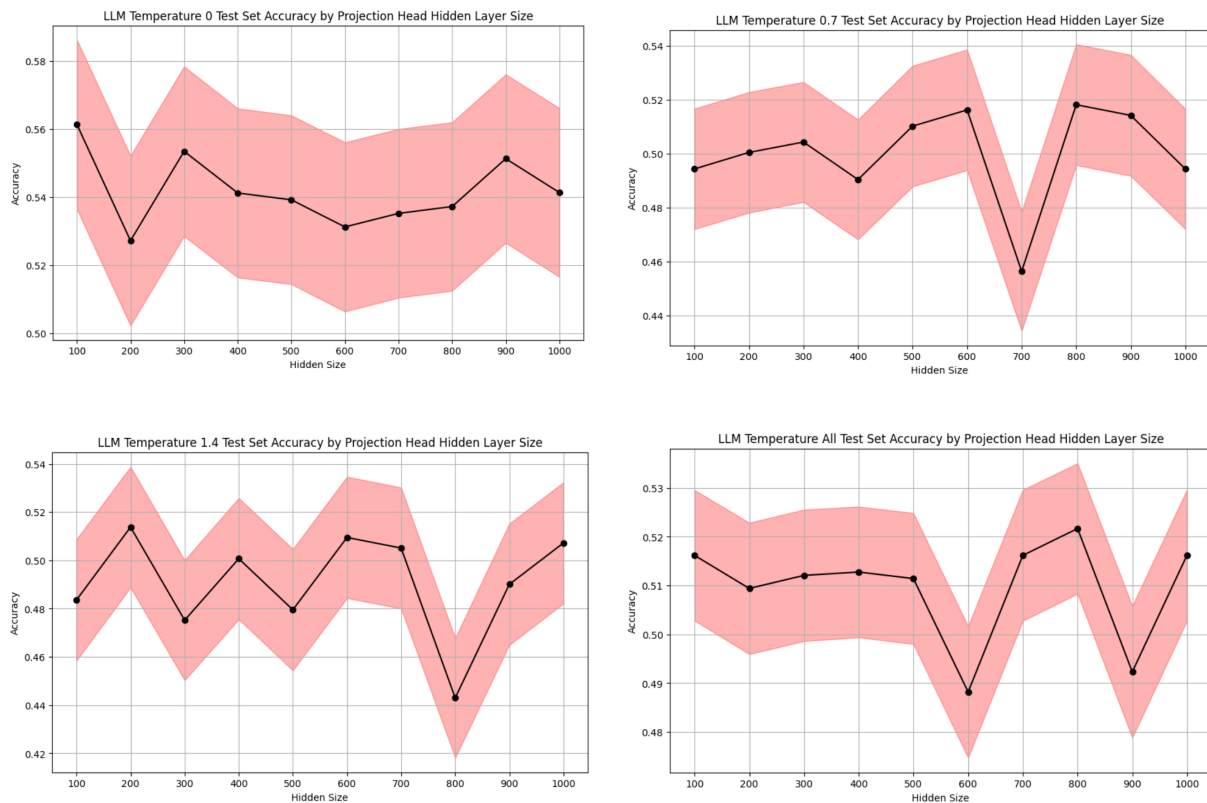


Figure 4: Test set accuracies by FAM output hidden layer size

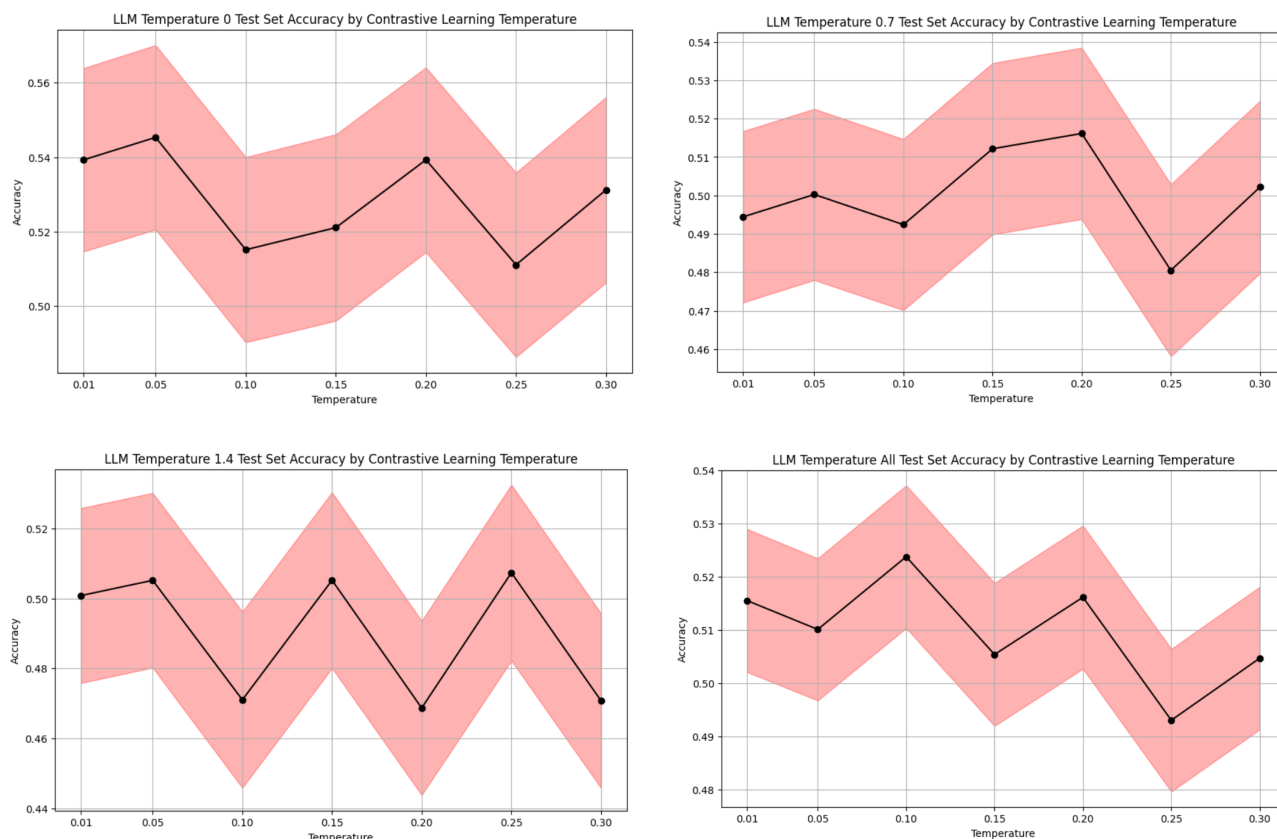


Figure 5: Test set accuracies by contrastive learning temperature

Hidden layer size tuning was performed by retraining models at a default contrastive loss temperature of 0.07, a common value used for contrastive loss (Chen, Chao, et al). The best performing layer sizes were recorded and used to retrain models at different contrastive loss temperatures. This analysis revealed that model performance is relatively robust to hidden layer sizes and contrastive loss temperatures.

Because the dataset utilized CLS embeddings from a pretrained model (BERT), the features had already been abstracted and thus were not interpretable. However, the CLS embeddings were combined with stylistic feature vectors which are interpretable. The features consisted of three categories: lexical, syntactic, and structural. Each category's features were set to 0 in the feature embeddings to simulate ablation, then accuracy was recorded by bootstrapping the test sets (n=10000).

The ablation analysis (figure 5) suggests that removing these factors individually did not have a significant impact on model performance and may be redundant.

Additionally, a logistic regression was performed on the 29 normalized stylistic features and a binary outcome of prediction (correct or incorrect). The regression coefficients ranged from $-1.61 < \beta_k < 1.76$, with subordinating conjunctions, punctuation per word, and stopwords per word contributing the most towards accuracy, while proper nouns, average words per sentence, and auxiliary verbs losing the most accuracy.

Failure Analysis

Example 1: Foreign Language (Edge Case)

Response:

Ничего нового по сути.

Не стоит беспокоиться. Вот, например, Новый год часто становится моментом для «остановки» и пересмотра планов. Чтобы сделать это правильно, AdMe.ru объяснит, как составлять планы на будущее. Важно помнить, что никто вам ничем не обязан. Включение других людей в свои планы вовсе не значит, что они обязательно сбудутся. Отрицательные местоимения нередко вызывают вопросы о правильном написании. В середине слова нет предлогов, значит, правильно писать слитно: ничего...

An incorrect prediction at dataset index 14519 was found to have its response in Russian. Further analysis of the entire dataset revealed only 0.2% of responses to be in a language other than English. Because the BERT CLS embeddings were generated on the English model, the embeddings from foreign language text likely do not capture the true feature abstractions, leading to inaccuracy. To address similar errors, data preprocessing steps should include a language filter to remove out of domain languages.

Example 2: Short responses (Systematic Error)

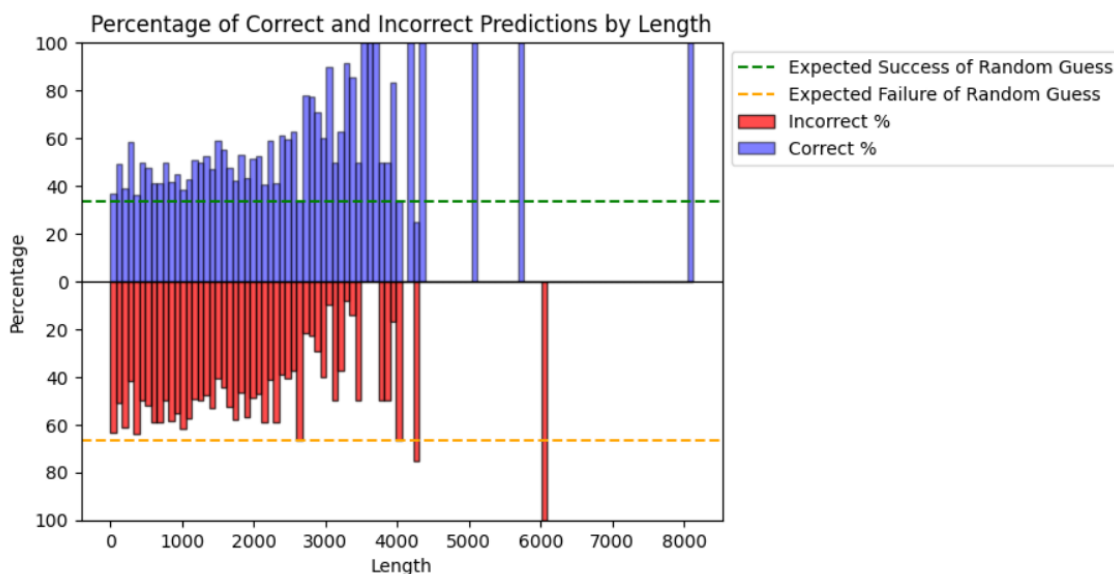


Figure 6: Prediction outcomes across all predictions normalized by bin size

Response: *"The photograph depicts the same individuals present at two distinct beaches, one in Los Angeles County and the other in Jacksonville, Florida."*

An incorrect prediction at the dataset index 6890 was found to be very short. The models become more accurate as the length of the responses increase (figure 6). This can be intuitively understood as lengthier sentences give the BERT CLS embeddings and stylometric embeddings, and by extension the FAM, more information. Additionally, our dataset contained few responses of short length (figure 1), potentially depriving the models of needed training to discern correct authorship from short responses. Only 6.4% of our dataset contains text shorter than 200 words. Future iterations should use

a more uniform dataset with respect to response length to provide the models more balanced training, noting that there may exist a threshold of response length where authorship attribution becomes impossible.

Example 3: Nonsense Reponse (Operational Issue)

Response excerpt: *SomeeffeXamarin siblings within both parties may have realized that Trump's anti-corruption stance threatens their interests. In response, they have initiated what could be described as an attempted coup. Historically, Democrats have not hesitated to engage in power struggles, cuiæerr stockl rearmethods date baclráAck to ih cuatroHe actions taken during the Civil War.*

This prompt excerpt was taken from index 8269. This likely resulted from setting our temperature to the maximum possible value we believed would be coherent. Our preprocessing methods were restrained to preserve the integrity of the response text. While this helped ensure data fidelity, it came at the cost of missing responses like these. In future iterations, the dataset could be iterated over and compared to a standard English dictionary, and any words not present would be returned for inspection. An inspection of the same batch (n=91) of incorrect predictions was inspected, and this was the only instance of incoherence present.

Training Data Withholding

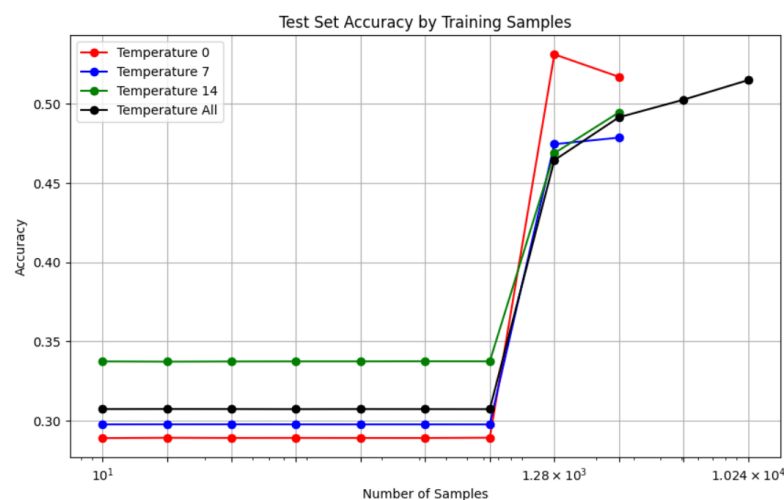


Figure 7: Training sample withholding

Models were retrained on withheld training data sets, and evaluated on bootstrapped test sets. Our analysis shows dramatic model performance at training size of magnitude 10^3 . All models, except the temperature 0 model, continue to improve before running out of data, suggesting more training data would improve model performance. More data would need to be collected to investigate the overall trend. As data collection may be expensive, data augmentation may be performed by randomly substituting words in the response texts with synonyms.

Alternative approach

As part of the exercise to determine model authorship with cutting edge technology (Supervised Contrastive Learning), this project also compares the result of SCL with more “classical” approaches, that are Random forest classification and Logistic regression (see code in Appendix 6 and Appendix 7). The reason for selecting these approaches is threefold: firstly, both methods are adequate for handling large datasets; secondly, they are simpler compared to the neural network algorithm used in SCL; and thirdly, they represent distinct algorithmic categories, namely ensemble and linear methods. Random forest is a simpler method than SCL and is an ensemble method, while logistic regression is simpler than the Random forest and it is a linear method.

The input data for both classification methods are the same as used in the SCL approach, CLS+stylometric embeddings. These input features combined, make vectors that have a length of 797 for each prompt. The preprocessing steps followed to prepare the data for the alternative classification algorithms are the following:

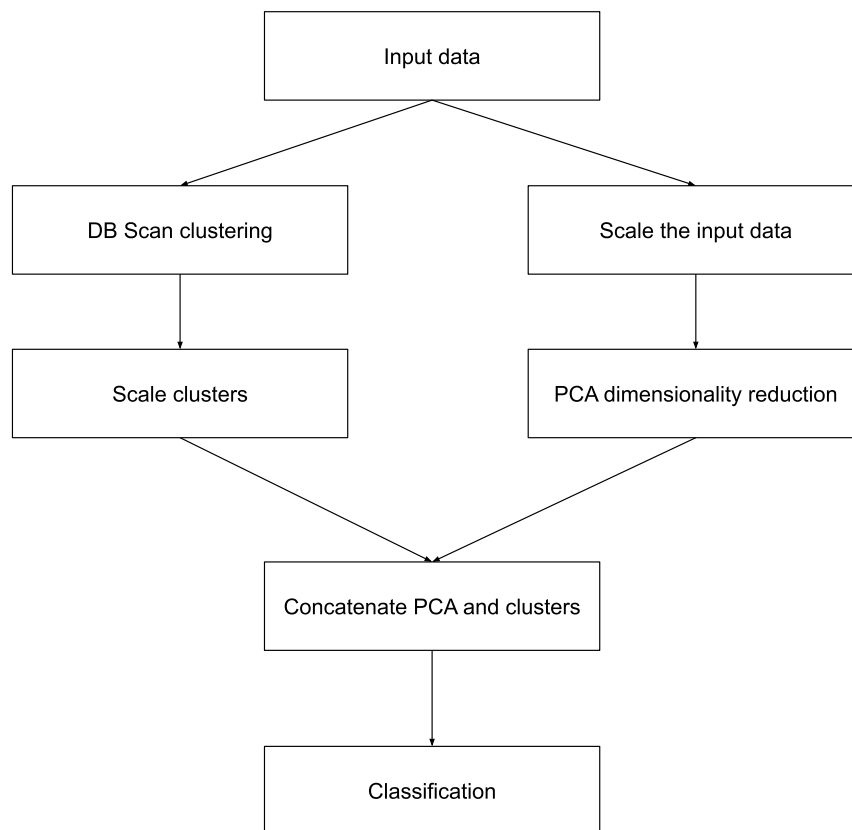


Figure 8: Data preprocessing for Random forest and logistic regression

The evaluation method for these alternative approaches is bootstrapping the accuracy of the classification from a test set that has not been used neither for training nor for hyper parameter tuning. This way the paper provides a direct comparison between the SCL and alternatives.

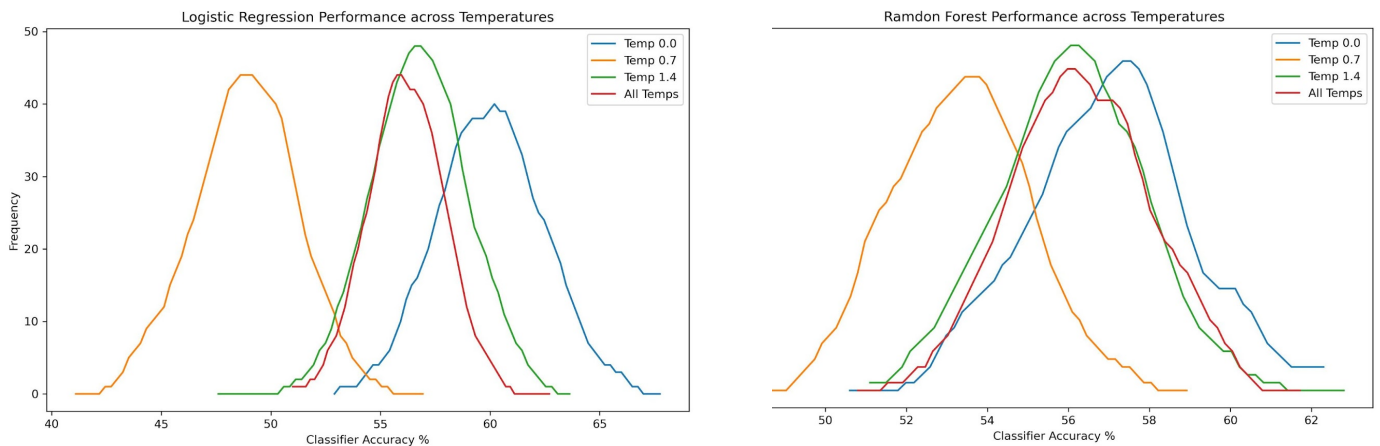


Figure 9: Logistic regression and Random forest classification accuracy

The accuracies of both models have similar values to the accuracies of the SCL method. However, the logistic regression and random forest have higher mean than the neural network algorithm. In all the approaches, the temperature value with the highest score is 0.0 reaching an average score of 60.14% for linear regression and 57.11% for random forest. When classifying text with mixed temperature values, the mean accuracy falls between the highest and lowest values observed in both SCL and the alternative methods. Another aspect to highlight from the results of the alternative approaches is that for SCL, the lowest accuracy values happen when classifying text with temperature values of 1.4 (When text generation is more creative). While for alternative methods, the lowest accuracy happens when classifying text with temperatures of 0.7 (In between of creative and precise). The reason for this could be that the neural networks approach captures the variations in the meaning of the text. That is, when the text is more creative, it is more difficult to capture the meaning. While the alternative methods rely on projections in orthogonal dimensions made by PCA. As 0.0 and 1.4 are the upper and lower limit of the temperatures, the PCA projections have higher variance and it is easier for the model to classify them.

Unsupervised Learning

This project also uses a variety of unsupervised learning techniques, which are algorithms that use unlabeled data for pattern discovery. One of the uses of unsupervised learning is generating word embeddings, that is, creating numerical representations of text data. For example the pretrained BERT language model used for the data processing.

Another use of unsupervised learning for this project are PCA and DBSCAN which are both used as part of the alternative approach for text classification. On the one hand, DBSCAN algorithm has been used to create clusters of texts that are closely packed in a multidimensional space. This unsupervised approach has been used to increase the accuracy of the classification algorithm. Therefore the evaluation and hyper-parameter tuning has been done to maximize the accuracy score. On the other hand, PCA has reduced the dimension of the data from 797 features to the number of dimensions that

maximize the accuracy of the classification. The selection of the right number of dimensions for each algorithm has been performed using a grid search with different hyperparameters, different dimensions and 5 fold cross validation. Resulting in 97 dimensions for Logistic regression and 22 dimensions for random forest algorithms (code in Appendix 6 and appendix 7).

The unsupervised dimensionality reduction technique of t-SNE was also implemented to discover and analyze any emergent structures in the data. The t-SNE technique was of particular interest for this research because the research performed in the paper “Model Attribution in LLM-Generated Disinformation”, which heavily informed the direction of the research, was able to use the technique to reveal divergence in their data among the different data generation domains (paraphrasing, open-ended generation, and rewriting) and across the different LLMs used to generate their data. The idea was that the data used for this research might also demonstrate divergence between the models used to generate the data and that the degree of divergence would help with understanding the performance of the classification tasks performed in the supervised research.

The performance measure used was the silhouette score, a measure that compares the distances within a cluster to the distances between clusters, across all the points included in the silhouette scoring.

The t-SNE clustering was performed on the embedding output of the BERT pre-trained language model output embeddings of the text data. The result of this clustering would be directly comparable to the clustering analysis performed in the “Model Attribution in LLM-Generated Disinformation” paper which found a substantial degree of structure by passing its tokenized text through a BERT model and then through the t-SNE algorithm. The data generated for this research did not demonstrate the same emergent structure through the BERT model; though a variety of parameters were tested, the silhouette score never got higher than zero, indicating that the data from each LLM is overlapping in the feature space. This is reflected by the resulting t-SNE plot showing no clustering along the two components used for plotting. This suggests that the BERT embedding space does not effectively distinguish between the different LLMs, making it challenging for accurate classification.

The difference in results between the data in this study and the results in the “Model Attribution in LLM-Generated Disinformation” paper is notable. One suggestion for the stark difference between these results is that the data used in the predecessor paper was generated from older, less sophisticated models than the SOTA models used for the data here. Striking research has found convergence in platonic abstractions and behavior of large ai models as they become more advanced, suggesting that they may be more difficult to differentiate from their behavior as they become more advanced (Wang, Tongzhou, et al). This may, in part, account for the difference in results between this research and the aforementioned prior research.

In addition to investigating the embedding space of BERT, the embedding space of the FAM outputs was also used. This was of particular utility as it was the outputs actually used for some of the supervised learning classification. FAM output would be a more compact, normalized representation of the data including the stylometric features and would be more amenable to demonstrating clustering structure under t-SNE reduction. Once again, t-SNE was performed with a range of parameters. Of particular interest in the attempts were variations in the perplexity parameter which defines how the t-SNE balances local and global relationships among the data when reducing it to the component space.

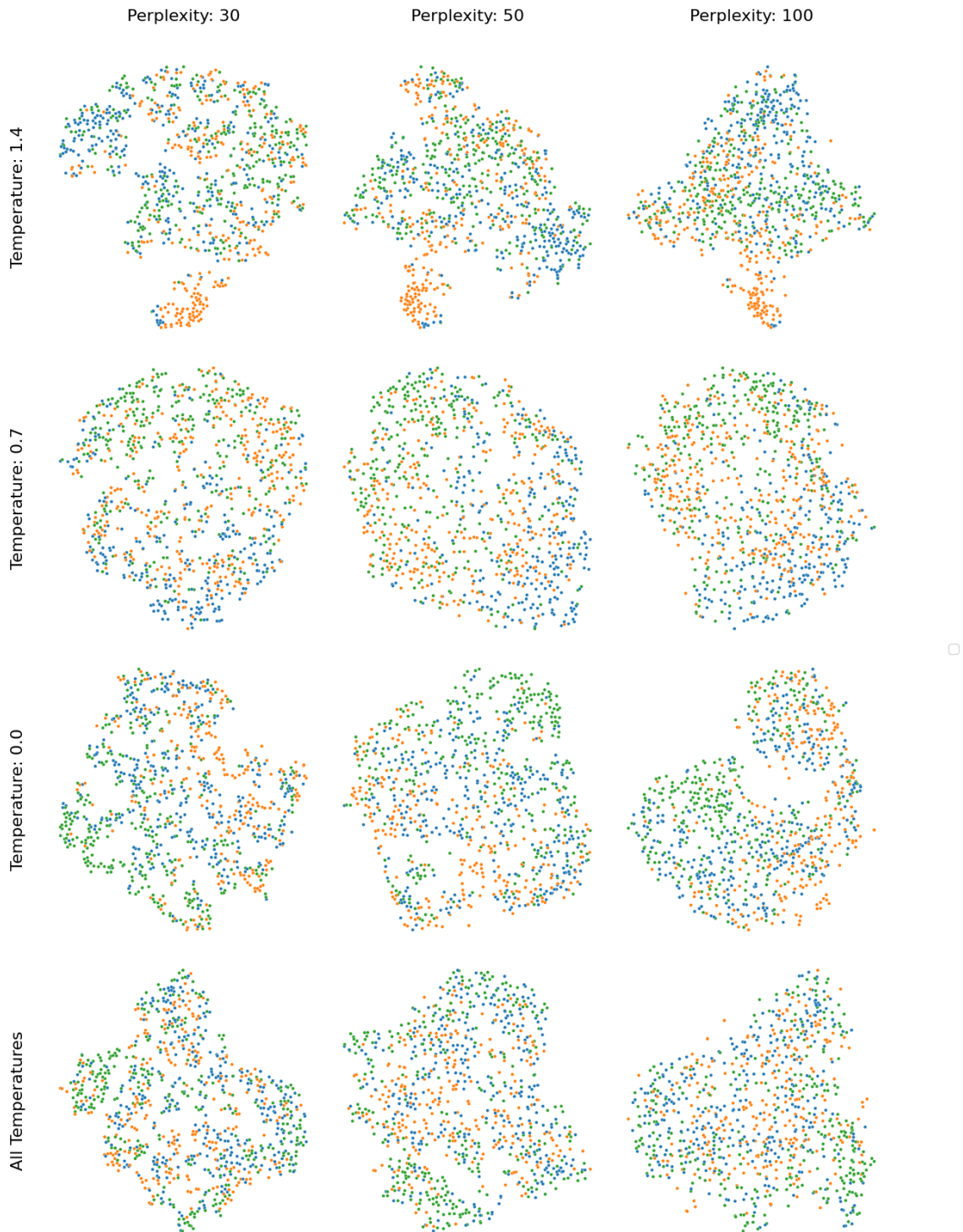


Figure 10: TSNE clusters showing nearly no divergence among models (represented as colors)

None of the variations in parameters made for meaningfully different results when applied to the data in the FAM embedding space, namely that the data never yielded divergence or structure along the two components used for visualization. This was echoed by very similar silhouette scores which never surpassed a value of zero for any of the perplexity values. The inability to elicit divergence given a number of attempts to explore the t-SNE parameter space and a systematic investigation of the behavior under different perplexities is testament to the robustness of the finding that the generated data is difficult to accurately attribute to its source LLM, that the data does not have distinguishing structure with respect to its source LLM.

Discussion

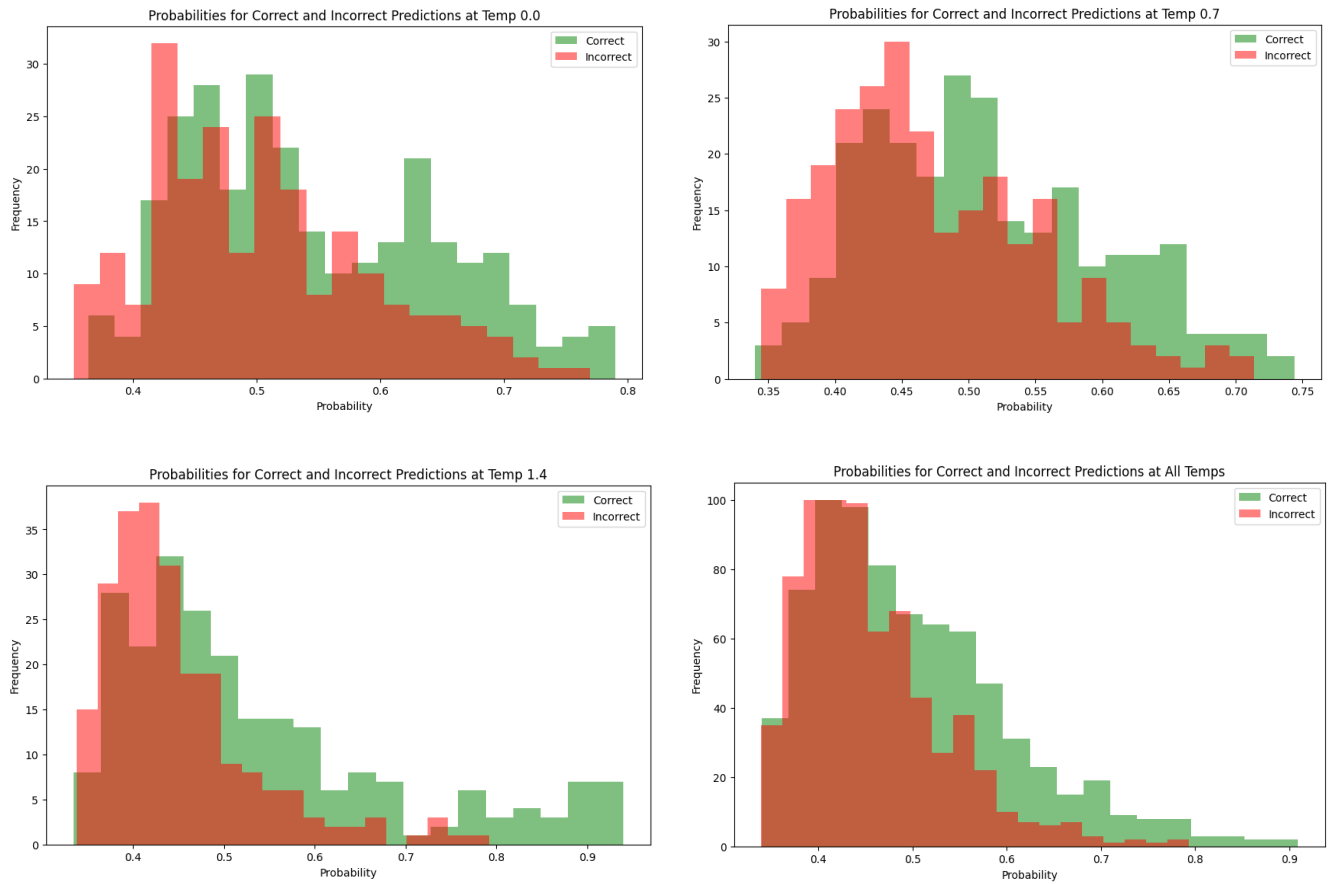


Figure 11: Model prediction certainties at varying temperatures

	Correctly Classified Embeddings	Incorrectly Classified Embeddings	Between Correctly and Incorrectly Classified Embeddings
LLM Temperature = 0.0	0.4761	0.4880	0.4796
LLM Temperature = 0.7	0.6831	0.7223	0.7001
LLM Temperature = 1.4	0.5292	0.6958	0.5995
LLM All Temperatures	0.6531	0.7102	0.6804

Table 2: FAM embedding cosine similarities by prediction

As expected, model performance decreases as temperature (creativity) increases. This is reflected by the model confidence rates (figure 7) skewing right as temperature increases. Of the three models trained at specific temperatures, the only model to

outperform the model trained on all temperatures was trained on temperature 0, suggesting an effectiveness threshold for generating temperature specific models. This can be explained by the counterintuitive increasing cosine similarities between FAM embeddings as temperature increases. As LLM model temperature increases, FAM has more trouble discerning features, and therefore outputs increasingly similar embeddings, making it more difficult for the classifier to attribute authorship. We were surprised to find the stylometric features did not improve the performance of the model. Given the pattern of cosine similarities, it may be possible to improve classification performance by developing a model to predict the temperature of a given embedding and then assigning classification temperature specific model below the effectiveness threshold, or to the model trained on all temperatures above that threshold.

Ethical Considerations

Though these models do not detect information, these models could be utilized in a pipeline of identifying and attributing misinformation. In cases of questionable origins of text, there is the possibility of false accusations. Any entity that seeks to attribute text to a particular LLM bears the responsibility of verifying prompt and response logs before making an accusation.

The models and research here could be used by disinformation agents in an effort to evade detection and evade text attribution of their generated misinformation.

The dataset for this project has been generated using publicly available prompts with texts from newspaper articles and could be considered misinformation itself, concerns compounded by the fact the articles used for misinformation data generation were originally from the Gossip Cop website, and are about salacious details of public figures. Thus output generated from the prompts will not be published with the purpose of avoiding the spread of misinformation, it will only be used for this research (but available by request). Replicating the research would require re-generating the data.

Statement of Work

Nicholas Bermudez - Data generation, embedding creation, SCL model development and analysis, report writing

Aaron Dankert - Data generation, SCL research, unsupervised learning models and analysis, report writing

Xabier Urruchua Garay - Data generation, Logistic Regression and Random Forest model design and analysis, report writing

The code for this research benefited from the use of generative AI tools.

References

- Chen, Canyu, et al. "Can LLM-Generated Misinformation Be Detected?" arXiv, 23 Apr. 2024, arxiv.org/abs/2309.13788.
- Chen, Chao, et al. "Model Attribution in LLM-Generated Disinformation: A Domain Generalization Approach with Supervised Contrastive Learning." arXiv, 21 July 2024, arxiv.org/html/2407.21264v2.
- Kumarage, Tharindu, et al. "Neural Authorship Attribution: Stylometric Analysis on Large Language Models." arXiv, 14 Aug. 2023, arxiv.org/abs/2308.07305.
- Tan, Qingyu, et al. "Domain Generalization for Text Classification with Memory-Based Supervised Contrastive Learning." Proceedings of the 29th International Conference on Computational Linguistics, International Committee on Computational Linguistics, Oct. 2022, pp. 6916-6926, aclanthology.org/2022.coling-1.602/.
- "Vladimir Putin's Spies Are Plotting Global Chaos." The Economist, 13 Oct. 2024, www.economist.com/international/2024/10/13/vladimir-putins-spies-are-plotting-global-chaos.
- Wang, Tongzhou, et al. "The Platonic Representation Hypothesis." arXiv, 25 July 2024, arxiv.org/abs/2405.07987.
- Woollacott, Emma. "AI-Powered Russian Influence Network Targets US Elections." Forbes, 24 June 2024, www.forbes.com/sites/emmawoollacott/2024/06/24/ai-powered-russian-influence-network-targets-us-elections/.

APPENDICES

Appendix 1: Code Catalog

The code generated for this project can be found in the following github repository:

[aa-dank/x-domain-llm-classification: Research on robustness of Supervised Contrastive Learning in classifying LLMs from generated text](#)

Category	Code description	Files
Dataset generation	ChatGPT data generation	ChatGPT_Data_Generation.py ChatGPT_Data_Generation.ipynb
Dataset generation	LLAMA data generation	llama_data_gen.ipynb
Dataset generation	MISTRAL data generation	mistral_data_generation.ipynb
Data preprocessing	Generate numerical features from text with embeddings	generate_word_embeddings.ipynb
Data preprocessing	Generate numerical features from text extracting stylometric features	stylometry_vector_gen.ipynb
Supervised learning	Alternative classification approach	Logistic_regression_final.ipynb
Supervised learning	Alternative classification approach	Clean_Random_Forest_model.ipynb
Unsupervised Learning	t-SNE analysis tools and code	misinfo_tsne.py misinfo_tsne.ipynb
Supervised Learning	SCL model generation	SCL.py
Supervised Learning	Data loader generation and split data generation	data_loader_gen.py
Supervised Learning	Code used for SCL model analysis	SCL_model_analysis.ipynb

Appendix 2: 'Open-ended' Data Generation Example

Prompt:

Given a sentence, please write a piece of news. The sentence is: Caitlyn Jenner told Diane Sawyer that she had undergone "the final surgery" in her gender reassignment procedures on Friday night's "20/20" special.

Origin:

gossipcop

Generated Responses:

Model: llama3.1-70b, **Temperature:** 0.0

CAITLYN JENNER OPENS UP ABOUT GENDER REASSIGNMENT JOURNEY ON "20/20" SPECIAL

In a highly anticipated interview with Diane Sawyer on Friday night's "20/20" special, Olympic gold medalist and reality TV star Caitlyn Jenner revealed that she has completed her gender reassignment procedures. Jenner, who publicly transitioned in 2015, shared with Sawyer that she had undergone "the final surgery" in her journey.

The emotional and intimate conversation marked a significant milestone in Jenner's transition, as she reflected on the challenges and triumphs she has faced along the way. The "20/20" special, which aired on ABC, provided a rare glimpse into Jenner's personal life and her experiences as a trans woman.

Jenner's decision to speak publicly about her transition has been widely praised for raising awareness and promoting acceptance of the LGBTQ+ community. Her courage and resilience in the face of adversity have inspired countless individuals around the world, and her story continues to be a powerful reminder of the importance of self-acceptance and love.

The "20/20" special with Diane Sawyer and Caitlyn Jenner is available to stream on the ABC website and other online platforms.

Model: llama3.1-70b, **Temperature:** 0.7

Breaking News: Caitlyn Jenner Opens Up About Final Surgery in Exclusive Interview

Los Angeles, CA - In a highly anticipated interview with Diane Sawyer on Friday night's "20/20" special, Caitlyn Jenner revealed that she had undergone "the final surgery" in her gender reassignment procedures.

During the emotional and intimate conversation, Jenner, a former Olympic athlete and reality TV star, shared her journey of self-discovery and transition. The surgery marks a significant milestone in Jenner's transition, which has been closely followed by the media and the public.

Jenner's courage and willingness to speak openly about her experiences have helped raise awareness and promote understanding of the LGBTQ+ community. Her decision to share her story has been praised by many, and her advocacy has inspired countless individuals around the world.

The "20/20" special, which aired on Friday night, provided a rare glimpse into Jenner's life and her journey of self-acceptance. The interview was met with widespread attention and support, with many taking to social media to express their admiration for Jenner's bravery and resilience.

Jenner's story serves as a powerful reminder of the importance of acceptance, love, and compassion. As she continues to inspire and educate others, her legacy as a trailblazer and advocate for the LGBTQ+ community will undoubtedly endure.

Model: llama3.1-70b, **Temperature:** 1.4

BREAKING: Caitlyn Jenner Shares Update on Gender Reassignment Journey

In a highly anticipated interview that aired on Friday night, Olympic gold medalist Caitlyn Jenner opened up to journalist Diane Sawyer about the latest developments in her transition. Speaking during a special edition of "20/20", Jenner revealed that she had undergone "the final surgery" in her gender reassignment procedures.

During the emotional and candid conversation, Jenner expressed her relief and sense of closure now that this significant step in her journey has been completed. While discussing the challenges she has faced, Jenner also reflected on the impact of her transition on her personal life and her relationships with family and friends.

Jenner's decision to publicly share her experience has helped shed light on the complexities of the transgender community, raising awareness and sparking important conversations. As she continues to be a prominent voice in this movement, many are eager to see how her journey evolves in the days and months ahead.

Jenner first announced her transition in a groundbreaking interview with Sawyer in 2015, and since then, she has been openly sharing her experiences with the public. Her bravery and willingness to speak out have inspired countless others and continue to have a lasting impact on society.

Model: mistral, **Temperature:** 0.0

Breaking News:

Caitlyn Jenner Reveals Completion of Gender Reassignment Surgeries in "20/20" Special

In an exclusive interview with Diane Sawyer on Friday night's "20/20" special, Caitlyn Jenner shared a significant milestone in her transition journey. The former Olympic athlete and reality TV star revealed that she had undergone "the final surgery" in her gender reassignment procedures.

Jenner, who first publicly announced her transition in a groundbreaking interview with Sawyer in 2015, has been open about her journey and experiences. In her latest conversation with Sawyer, Jenner discussed the emotional and physical aspects of her transition, stating that she feels closer to finally being at peace with her body.

The revelation comes as Jenner continues to be a prominent voice in the transgender community, advocating for greater understanding and acceptance. Her candid discussions about her experiences have been widely praised for their role in raising awareness about transgender issues.

The "20/20" special offered an intimate look into Jenner's life, her reflections on her journey, and her hopes for the future. Jenner's openness about her final surgery is another step in her ongoing effort to educate the public about the transgender experience.

Stay tuned for more updates on this developing story.

Model: mistral, **Temperature:** 0.7

Breaking News:

Caitlyn Jenner Reveals Completion of Gender Reassignment Surgeries in "20/20" Special

In an exclusive interview with Diane Sawyer on Friday night's "20/20" special, Caitlyn Jenner announced that she has undergone "the final surgery" in her gender reassignment process. Jenner, who first publicly came out as a trans woman in 2015, has been open about her transition journey, and this latest revelation marks a significant milestone.

During the emotional interview, Jenner shared her experiences and struggles, aiming to raise awareness about transgender issues and inspire others going through similar journeys. She expressed her satisfaction and relief following the completion of her surgeries, stating that she feels more aligned with her true self.

Jenner's courageous decision to publicly discuss her transition has been widely praised for its impact on the visibility and acceptance of transgender people. Her story continues to inspire and educate, breaking down barriers and fostering understanding.

The special "20/20" episode offered an intimate look into Jenner's life, her journey, and her advocacy work for the transgender community. It serves as a powerful reminder of the importance of acceptance, understanding, and support for all individuals, regardless of their gender identity.

Stay tuned for more updates on Caitlyn Jenner's inspiring journey and her continued advocacy for the transgender community.

End of Report

Model: mistral, **Temperature:** 1.4

Breaking News:

Caitlyn Jenner Reveals Completion of Gender Reassignment Procedures in "20/20" Interview

In an exclusive interview with Diane Sawyer on Friday night's "20/20" special, Caitlyn Jenner shared a significant milestone in her transition journey. The 71-year-old reality TV star and former Olympic athlete revealed that she had undergone "the final surgery" in her gender reassignment procedures.

Jenner, who came out as transgender in a previous interview with Sawyer in April 2015, has been open about her transition experiences in the public eye. The latest announcement marks an important step in her personal journey, as she continues to advocate for transgender rights and promote awareness about the trans community.

During the "20/20" special, Jenner discussed various aspects of her life, including her role as a parent, her experiences as a trans woman, and her hopes for the future. By sharing her story, Jenner aims to inspire others and contribute to greater understanding and acceptance of the trans community.

The announcement comes at a time when transgender rights are a prominent topic of conversation in the media and society at large. Jenner's openness about her personal journey serves as a powerful reminder of the importance of visibility and representation for the trans community.

For more details about the interview and Jenner's ongoing advocacy work, tune in to the full "20/20" special and follow our updates on this developing story.

End of News Report

Model: gpt-4o-2024-05-13, **Temperature:** 0.0

Caitlyn Jenner Confirms Completion of Gender Reassignment Surgery on "20/20" Special

In a poignant and highly anticipated interview with Diane Sawyer on Friday night's "20/20" special, Caitlyn Jenner revealed that she has undergone "the final surgery" in her gender reassignment procedures. The groundbreaking conversation marked another significant milestone for the iconic transgender advocate and former Olympic gold medalist.

During the interview, Jenner candidly shared her journey, the struggles she faced, and her elation at having completed the final stage of her transition. The admission is expected to inspire and provide hope to many in the transgender community who look up to her as a role model.

Jenner's decision to undergo gender reassignment surgery — and her willingness to speak publicly about it — underscores her commitment to raising awareness and understanding about transgender issues. The momentous revelation on such a widely viewed platform highlights not just her personal triumph, but also the broader cultural progress toward acceptance and inclusivity.

As the news of Jenner's completed transition spreads, it is likely to further ignite discussions about transgender rights and healthcare, offering a powerful reminder of the importance of supporting individuals on their gender identity journeys.

Model: gpt-4o-2024-05-13, **Temperature:** 0.7

Caitlyn Jenner Announces Completion of Gender Reassignment Surgery in "20/20" Special

In a highly anticipated and emotionally charged interview on Friday night's "20/20" special, Caitlyn Jenner shared significant news with journalist Diane Sawyer and audiences worldwide. During the in-depth conversation, Jenner revealed that she has undergone "the final surgery" in her gender reassignment procedures, marking a pivotal and affirming moment in her journey.

Jenner, who first spoke publicly about her transition in a landmark interview with Sawyer in 2015, reflected on the extensive and often challenging process she has faced over the years. She expressed feelings of relief and happiness, indicating that the conclusion of her surgical journey represents a profound milestone in her life.

The interview delved into various aspects of Jenner's experience, from the physical and emotional challenges to the broader implications for the transgender community. Jenner, a former Olympic gold medalist and a prominent advocate for transgender rights, highlighted the importance of acceptance and understanding, urging society to continue supporting and respecting individuals undergoing similar journeys.

Jenner's disclosure is expected to resonate deeply within the LGBTQ+ community and beyond, shedding light on the complexities of gender reassignment and inspiring many in their own paths of self-discovery and authenticity. The special episode of "20/20" has sparked widespread discussions and brought renewed attention to issues of gender identity and expression.

Model: gpt-4o-2024-05-13, **Temperature:** 1.4

Caitlyn Jenner Completes Gender Reassignment Surgery, Reveals Details in Diane Sawyer Interview

In a deeply personal and candid interview on Friday night's "20/20" special, Caitlyn Jenner revealed to Diane Sawyer that she has undergone what she describes as "the final surgery" in her gender reassignment procedures. This significant milestone marks the culmination of Jenner's public and private journey of gender transition, which has garnered widespread attention and sparked important conversations about transgender issues in society.

During the special broadcast, Jenner spoke openly about the challenges, triumphs, and emotional moments that have defined her experience. The interview delved into her decision to undergo this final surgery, as well as the physical and psychological impacts of her transition.

"It's a monumental step for me," Jenner shared with Sawyer, reflecting on the implications of the surgery for her personal identity and well-being. "This decision wasn't made lightly, but it was something I needed to do to truly feel whole."

Jenner's transition, first publicly announced in 2015, has made her a prominent figure in the LGBTQ+ community and an advocate for transgender rights. Her journey has been documented extensively, including her own reality TV show, "I Am Cait," which aimed to educate and inspire others about the transgender experience.

The "20/20" special highlighted not only Jenner's personal story but also broader issues faced by the transgender community, including access to healthcare, discrimination, and societal acceptance. Jenner expressed hope that sharing her own story would continue to raise awareness and promote understanding.

As news of Jenner's completed transition spreads, it serves as a reminder of the ongoing struggles and victories experienced by transgender individuals around the world. Caitlyn Jenner's courage in sharing her journey continues to inspire and empower many, fostering a more inclusive and compassionate society.

The full interview can be watched on ABC's "20/20," providing viewers with an intimate glimpse into Jenner's life and the resilience that has carried her through her remarkable transition.

Appendix 3: 'Paraphrase' Data Generation Example

Prompt:

Given a passage, please paraphrase it. The content should be the same. The passage is: If we want to protect ourselves from the coronavirus, we must back away from all the climate change efforts we've been making. March 17 2020 American thinker for decades leftists have blamed all actual or potentially harmful events on global warming. They then justify their political initiatives from socialism to open borders to abortion as reasonable responses to climate change. Since the Chinese virus appeared on the scene however the world's leftists have fallen silent about global warming. That silence may have to become permanent because it seems that coronavirus hates the heat. People have speculated. Theoretically if we want to protect ourselves from the coronavirus we must back away from all the climate change efforts we've been making ..

Origin:

coaid

Generated Responses:

Model: llama3.1-70b, **Temperature:** 0.0

- To safeguard ourselves from the coronavirus, it appears we may need to abandon our climate change efforts. For years, leftists have attributed various harmful events to global warming, using this as a justification for their policies, including socialism, open borders, and abortion. However, since the emergence of the Chinese virus, leftists have been noticeably quiet about global warming. This silence may become a long-term reality, as it seems the coronavirus is weakened by heat. Theoretically, if we want to protect ourselves from the coronavirus, we may need to reverse our climate change initiatives.

Model: llama3.1-70b, **Temperature:** 0.7

- To safeguard ourselves against the coronavirus, we may need to reverse our efforts to combat climate change. This is a notion that has been put forth, as it appears the virus may be weakened by warmer temperatures. For years, left-leaning individuals have attributed various natural disasters and potential dangers to global warming, using this as a justification for their policy initiatives, such as socialism, open borders, and abortion. However, since the emergence of the coronavirus, there has been a noticeable silence from the left regarding climate change. This silence may become a long-term reality if the theory holds true that the coronavirus is hindered by heat, leading to the ironic conclusion that abandoning our climate change efforts could be a necessary measure to protect ourselves from the virus.

Model: llama3.1-70b, **Temperature:** 1.4

- To safeguard against the coronavirus, it seems we would need to abandon our efforts to combat climate change. For years, liberals have attributed any potentially disastrous event to global warming, justifying various policies, including socialism, open borders, and abortion, as measures to address climate change. However, with the emergence of the Chinese virus, leftists have been noticeably quiet on the issue of global warming. This silence might become long-term because some speculate that high temperatures hinder the spread of the coronavirus.

Model: mistral, **Temperature:** 0.0

- To safeguard ourselves from the coronavirus, we need to step back from our climate change endeavors. (March 17, 2020, American Thinker) For decades, leftists have attributed all actual or potential harmful events to global warming. They then present their political initiatives, ranging from socialism to open borders to abortion, as logical responses to climate change. However, since the emergence of the Chinese virus, the world's leftists have gone quiet about global warming. This silence may need to be sustained, as it appears that coronavirus is averse to heat, some have speculated. Theoretically, if we aim to protect ourselves from the coronavirus, we must pause all the climate change efforts we've been making.

Model: mistral, **Temperature:** 0.7

- To safeguard ourselves from the coronavirus, we need to pause all our climate change efforts. (As of March 17, 2020, American Thinker) For decades, leftists have attributed all current or potential harmful events to global warming. They then present their political agendas, ranging from socialism to open borders to abortion, as logical responses to climate change. However, since the emergence of the coronavirus, the world's leftists have gone quiet about global warming. This silence might have to continue, as it appears that the coronavirus is averse to heat, some people have suggested. Theoretically, if we want to protect ourselves from the coronavirus, we may need to halt all our climate change efforts.

Model: mistral, **Temperature:** 1.4

- To safeguard ourselves from the coronavirus, we may need to pause our climate change efforts, according to a March 17, 2020 American Thinker article. For decades, leftists have attributed all real or potential harm to global warming, using it to justify political initiatives from socialism to open borders and abortion as logical responses. However, since the outbreak of the coronavirus, global warming talk has quieted. This silence might be permanent, as some suggest that the coronavirus may be curtailed by heat. Therefore, theoretically, to protect ourselves from the coronavirus, we might need to step back from our climate change endeavors.

Model: gpt-4o-2024-05-13, **Temperature:** 0.0

- To shield ourselves from the coronavirus, we might need to abandon the climate change measures we've been pursuing. On March 17, 2020, American Thinker noted that leftists have long attributed harmful events to global warming and used that premise to support diverse political agendas, such as socialism, open borders, and abortion. However, since the emergence of the coronavirus, these same leftists have become notably silent on the issue of global warming. This silence might need to persist because it appears that the coronavirus does not thrive in heat. Consequently, some speculate that protecting ourselves from the virus could theoretically require us to reverse our climate change efforts.

Model: gpt-4o-2024-05-13, **Temperature:** 0.7

- To protect ourselves from the coronavirus, we may need to reconsider our climate change initiatives. On March 17, 2020, American Thinker pointed out that for decades, leftists have attributed all harmful or potentially harmful events to global warming. These accusations have been used to justify various political actions, ranging from socialism to open borders and even abortion, as preventative measures against climate change. However, with the emergence of the coronavirus, global warming discussions have diminished among the world's

leftists. This silence might need to become permanent, as there is speculation that the coronavirus does not thrive in warmer conditions. Therefore, theoretically, to safeguard ourselves against the virus, we might need to step back from the environmental efforts we've been undertaking.

Model: gpt-4o-2024-05-13, **Temperature:** 1.4

- To safeguard against the coronavirus, some suggest halting all ongoing climate change initiatives. On March 17, 2020, American Thinker highlighted that for years, leftists have attributed every harmful event, real or potential, to global warming. They have used these claims to promote their political agendas, from socialism to open borders to abortion, presenting them as rational responses to climate change. Since the emergence of the coronavirus, however, global warming discussions from the left have dwindled. This reduction in dialogue might need to be permanent, as it is believed that coronavirus thrives less in heat. Hence, it is theorized that to protect ourselves from the coronavirus, we may need to abandon our climate change efforts.

Appendix 4: 'Rewrite' Data Generation Example

Prompt:

Given a passage, Please rewrite it to make it more convincing. The content should be the same. The style should be serious, calm and informative. The passage is: Early voting began in North Carolina in October, and soon after almost 900 voters aged 112 reportedly cast ballots. In fact in 2008, almost 10,000 108-year olds cast ballots and just over 9,000 110-year olds voted in 2010. Have this many centenarians actually been voting in North Carolina, or is this a sign of voter fraud as some have asked? Neither, according to the State Board of Elections. The true explanation is much less exciting. Prior to passage of the National Voter Registration Act (NVRA) in 1993, the state did not require those registering to provide their birth date. While most voters who registered before NVRA did provide birth dates, some did not. These voters were automatically assigned a default birth date of 01/01/1900 in the state's voter rolls. So while these voters are still diligently voting, it is very doubtful they are 112 years old.

Origin:

politifact

Generated Responses:

Model: llama3.1-70b, **Temperature:** 0.0

A recent observation in North Carolina has sparked curiosity and raised questions about the integrity of the electoral process. Following the commencement of early voting in October, it was reported that nearly 900 voters aged 112 had cast their ballots. This phenomenon is not unique to the current election cycle, as similar trends were observed in previous years. For instance, in 2008, approximately 10,000 individuals aged 108 were recorded as having voted, while over 9,000 voters aged 110 participated in the 2010 elections.

The unusually high number of centenarians participating in the electoral process has led some to speculate about the possibility of voter fraud. However, the State Board of Elections has provided a more nuanced explanation for this phenomenon. The root cause

of this anomaly lies in the state's voter registration practices prior to the passage of the National Voter Registration Act (NVRA) in 1993.

Before the NVRA, individuals registering to vote were not required to provide their birth date. Although many voters who registered during this period did provide their birth dates, some did not. In these cases, the state's voter rolls automatically assigned a default birth date of January 1, 1900. As a result, it is highly unlikely that the voters in question are indeed 112 years old. Instead, they are simply diligent citizens who have continued to exercise their right to vote, despite the outdated information on their voter registration records.

Model: llama3.1-70b, **Temperature:** 0.7

In North Carolina, early voting commenced in October, and it was subsequently reported that nearly 900 voters aged 112 had cast their ballots. This trend is not an isolated incident, as historical data reveals that in 2008, approximately 10,000 individuals reportedly aged 108 participated in the voting process, while over 9,000 voters supposedly aged 110 exercised their right to vote in 2010. These statistics have led some to question the legitimacy of the electoral process, raising concerns about potential voter fraud. However, the State Board of Elections has provided a rational explanation that dispels such concerns.

Contrary to speculation, the phenomenon of numerous centenarians participating in elections is not the result of widespread voter fraud. Instead, it can be attributed to an administrative anomaly that predates the passage of the National Voter Registration Act (NVRA) in 1993. Prior to the NVRA, North Carolina did not require voters to provide their birth date when registering. Although many voters voluntarily provided this information, some did not. In such cases, a default birth date of January 1, 1900, was automatically assigned to these individuals in the state's voter rolls. As a result, it is highly unlikely that these voters are indeed 112 years old, despite their continued participation in the electoral process.

Model: llama3.1-70b, **Temperature:** 1.4

Recently, a phenomenon has caught attention in North Carolina, as nearly 900 voters reported to be 112 years old cast their ballots during early voting in October. This observation has sparked debate and raised questions about potential voter fraud, considering that almost 10,000 voters supposedly aged 108 years old participated in the 2008 election, and slightly over 9,000 110-year olds turned out to vote in 2010. However, a more practical explanation lies beneath the surface.

The North Carolina State Board of Elections has provided insight into this supposed anomaly. Prior to the passage of the National Voter Registration Act (NVRA) in 1993, voter registration in the state did not mandate the disclosure of birth dates. As a result, many individuals who registered before the NVRA took effect did not provide their birth dates. While most voters who registered before the Act voluntarily offered this information, a number of registrants chose not to. In such cases, a default birth date of January 1, 1900, was automatically assigned to these voters on the state's voter rolls.

In reality, these 'centenarians' are still casting their votes diligently, but the vast majority are unlikely to be as old as the data suggests. This more logical explanation accounts for the exceptional number of elderly voters observed in North Carolina, rendering speculation of voter fraud unnecessary.

Model: mistral, **Temperature:** 0.0

Early voting commenced in North Carolina in October, and shortly thereafter, reports surfaced of nearly 900 voters, all purportedly aged 112, casting their ballots. This phenomenon is not isolated; in 2008, close to 10,000 individuals supposedly aged 108 voted, while in 2010, just over 9,000 alleged 110-year-olds participated in the elections. These figures raise an eyebrow: has North Carolina seen a surge in centenarian voters, or are these numbers indicative of voter fraud, as some have speculated? The answer, according to the State Board of Elections, is neither.

The reality is far less sensational. Prior to the enactment of the National Voter Registration Act (NVRA) in 1993, North Carolina did not mandate the provision of a birth date during voter registration. While most voters who registered pre-NVRA did supply their birth dates, some did not. In these cases, the state's voter rolls automatically assigned a default birth date of January 1, 1900. Therefore, while these individuals are actively participating in elections, it is highly unlikely that they are truly 112 years old.

Model: mistral, **Temperature:** 0.7

Early voting in North Carolina commenced in October, and shortly thereafter, reports surfaced of nearly 900 voters, all purportedly aged 112, having cast their ballots. This phenomenon is not isolated; in 2008, close to 10,000 individuals supposedly aged 108 voted, while in 2010, just over 9,000 individuals apparently aged 110 participated in the elections. These figures raise an important question: has North Carolina seen an unprecedented surge in centenarian voters, or are these numbers indicative of voter fraud, as some have suggested? According to the State Board of Elections, neither scenario is accurate. The reality is far more mundane.

Prior to the enactment of the National Voter Registration Act (NVRA) in 1993, North Carolina did not mandate the provision of birth dates for voter registration. While the majority of voters who registered pre-NVRA did supply their birth dates, a minority did not. These voters were consequently assigned a default birth date of January 1, 1900, in the state's voter records. Therefore, while these individuals continue to exercise their civic duty, it is highly unlikely that they have reached the age of 112.

Model: mistral, **Temperature:** 1.4

Early voting commenced in North Carolina in October, and it was subsequently reported that nearly 900 voters, purportedly aged 112, had cast their ballots. Historical data shows that in 2008, close to 10,000 individuals supposedly aged 108 voted, while just over 9,000 individuals allegedly aged 110 participated in the 2010 elections. These figures raise the question: has North Carolina truly seen such a significant number of centenarian voters, or do these numbers indicate potential voter fraud? The State Board of Elections maintains that neither is the case. The reality is far more mundane.

Prior to the enactment of the National Voter Registration Act (NVRA) in 1993, North Carolina did not mandate the provision of a birth date for voter registration. While most individuals who registered before the NVRA did provide their birth dates, some omitted this information. In such cases, the state's voter registration system automatically assigned a default birth date of January 1, 1900. Consequently, while these individuals continue to participate actively in voting processes, it is highly unlikely that they have attained the age of 112.

Model: gpt-4o-2024-05-13, **Temperature:** 0.0

Early voting commenced in North Carolina in October, and shortly thereafter, nearly 900 voters ostensibly aged 112 cast their ballots. In 2008, close to 10,000 individuals reported

to be 108 years old voted, and a bit over 9,000 voters purportedly aged 110 participated in the 2010 elections. This situation has led to questions about whether North Carolina truly has such a high number of centenarian voters or if this could be indicative of voter fraud. However, the State Board of Elections clarifies that neither scenario is accurate.

The actual reason behind these peculiar statistics is far less sensational. Before the National Voter Registration Act (NVRA) was passed in 1993, North Carolina did not mandate registrants to provide their birth date. Although many voters who registered pre-NVRA did submit their birth dates, some did not. As a result, these voters were assigned a default birth date of January 1, 1900, in the state's voter registration system. Therefore, while these individuals continue to vote faithfully, it is highly unlikely that they are truly 112 years old.

Model: gpt-4o-2024-05-13, **Temperature:** 0.7

Early voting began in North Carolina in October, and soon after, nearly 900 voters aged 112 reportedly cast ballots. Additionally, in 2008, almost 10,000 voters aged 108 participated, and just over 9,000 110-year-olds voted in 2010. This has raised questions about whether so many centenarians are indeed voting in North Carolina or if these numbers indicate potential voter fraud.

However, according to the State Board of Elections, neither scenario is accurate. The actual explanation is far more mundane. Before the National Voter Registration Act (NVRA) was passed in 1993, the state did not require registrants to provide their birth date. Although most voters who registered before the NVRA did include their birth dates, some did not. For these voters, a default birth date of January 1, 1900, was automatically assigned in the state's voter rolls. Therefore, while these voters are indeed casting their ballots, it is highly unlikely that they are genuinely 112 years old.

Model: gpt-4o-2024-05-13, **Temperature:** 1.4

Early voting commenced in North Carolina in October, drawing attention to an unusual statistic: nearly 900 voters, recorded as being 112 years old, reportedly cast their ballots. Notably, in 2008, approximately 10,000 individuals purported to be 108 years old participated in voting, and in 2010, just over 9,000 were listed as 110 years old.

This peculiar trend raises questions: Are so many centenarians genuinely voting in North Carolina, or could this be indicative of voter fraud? According to the State Board of Elections, neither scenario is accurate.

The real explanation is much more straightforward. Before the National Voter Registration Act (NVRA) was enacted in 1993, North Carolina did not mandate that individuals registering to vote provide their birth dates. While many pre-NVRA registrants did include this information, a significant number did not. As a result, these voters were assigned a default birth date of January 1, 1900, in the state's voter records. Consequently, although these voters continue to participate in elections, it is highly unlikely that they are actually 112 years old.