

Predictive Insights into LGBTQ+ Minority Stress: A Transductive Exploration of Social Media Discourse

Santosh Chapagain
Utah State University
santosh.chapagain@usu.edu

Yuxuan Zhao
New Mexico State University
zyx8010@nmsu.edu

Taylor K. Rohleen
University of Florida
trohleen@ufl.edu

Shah Muhammad Hamdi
Utah State University
s.hamdi@usu.edu

Soukaina Filali Boubrahimi
Utah State University
soukaina.boubrahimi@usu.edu

Ryan E. Flinn
University of North Dakota
ryanelliottflinn@gmail.com

Emily M. Lund
University of Alabama
Ewha Womans University
emlund@ua.edu

Dannie Klooster
Oklahoma State University
dannie.klooster@okstate.edu

Jillian R. Scheer
Syracuse University
jrscheer@syr.edu

Cory J. Cascalheira
New Mexico State University
cjcascalheira@gmail.com

Abstract—Individuals who identify as sexual and gender minorities, including lesbian, gay, bisexual, transgender, queer, and others (LGBTQ+) are more likely to experience poorer health than their heterosexual and cisgender counterparts. One primary source that drives these health disparities is minority stress (i.e., chronic and social stressors unique to LGBTQ+ communities’ experiences adapting to the dominant culture). This stress is frequently expressed in LGBTQ+ users’ posts on social media platforms. However, these expressions are not just straightforward manifestations of minority stress. They involve linguistic complexity (e.g., idiom/lexical diversity), rendering them challenging for many traditional natural language processing methods to detect. In this work, we designed a hybrid model using Graph Neural Networks (GNN) and Bidirectional Encoder Representations from Transformers (BERT), a pre-trained deep language model to improve the classification performance of minority stress detection. We experimented with our model on a benchmark social media dataset for minority stress detection (LGBTQ+ MiSSoM+). The dataset is comprised of 5,789 human-annotated Reddit posts from LGBTQ+ subreddits. Our approach enables the extraction of hidden linguistic nuances through pretraining on a vast amount of raw data, while also engaging in transductive learning to jointly develop representations for both labeled training data and unlabeled test data. The RoBERTa-GCN model achieved an accuracy of 0.86 and an F1 score of 0.86, surpassing the performance of other baseline models in predicting LGBTQ+ minority stress. Improved prediction of minority stress expressions on social media could lead to digital health interventions to improve the wellbeing of LGBTQ+ people—a community with high rates of stress-sensitive health problems.

Index Terms—sexual and gender minority, deep learning, transductive learning, bidirectional encoder representation of transformers (BERT), graph neural networks, graph convolution networks (GCN), stress

I. INTRODUCTION

Lesbian, gay, bisexual, transgender, queer, and other sexual and gender minority (LGBTQ+) individuals experience

significantly poorer physical and mental health outcomes in comparison to heterosexual and cisgender populations: higher incidences of asthma, activity limitation, cardiovascular risk [1], human immunodeficiency virus (HIV), chronic health conditions [2], and depression, anxiety, post-traumatic stress disorder (PTSD), substance use, self-harm, and suicidality [3]–[5]. Robust evidence demonstrates that social determinants of health—particularly, minority stress—contribute to the prevailing health disparities between heterosexual/cisgender and LGBTQ+ individuals [4], [6].

Minority stress [4] theory posits that LGBTQ+ individuals face widespread environmental threats globally [7], [8]. A main rationale for minority stress theory is that, by identifying as a sexual and gender minority, LGBTQ+ individuals face criticism and discrimination by members of the dominant culture. In an effort to adapt and to assimilate to heteronormative and cissexist societal expectations, chronic, social stressors tend to emerge. These include LGBTQ+ individuals experiencing the effects of stigmatizing and discriminating social systems, including prejudiced events (i.e., acute or chronic external stressful events), identity concealment (i.e., hiding LGBTQ+ identities due to fear of harm), expected rejection (i.e., expectations accompanied with feelings of vigilance towards possible prejudiced events), and internalized stigma (i.e., making negative societal attitudes part of their nature). Moreover, transgender and non-binary subgroups tend to present with an additional indicator of minority stress, gender dysphoria (i.e., experience of stress or discomfort stemmed from misalignment between perceived gender identity and assigned sex at birth) [9]. Despite established empirical support for the substantial health disparities faced by LGBTQ+ individuals and the central role of minority stress experiences in driving these inequities, there remains a gap in computational support for validating minority stress.

Existing applications of artificial intelligence (AI) in social science and healthcare have demonstrated their value in promoting general well-being across a multitude of disciplines—most notably, the expansion of existing theoretical frameworks in the social sciences and the emersion of AI-powered interventions [10], [11]. Natural language processing (NLP), specifically, holds the capacity to foster equity for minoritized communities [12]. That is, NLP allows for an accessible and timely detection of mental stress, opening novel pathways for the delivery of unique interventions that target minoritized individuals’ nuanced experiences and, in turn, promote the well-being of minoritized communities [11], [13].

For LGBTQ+ communities, the Internet—specifically, social media platforms like Reddit, on which one can post anonymously—has become a crucial space for LGBTQ+ individuals to come out, connect with peers, and seek support with less risk of experiencing disapproval and prejudice from their peers online, due to subgroup users sharing similar identities [14], [15] (e.g., LGBTQ+ users can find Reddit forums comprised of other LGBTQ+ users). Because LGBTQ+-specific Reddit forums are available for anonymous posting by LGBTQ+ people, Reddit provides an excellent space to study minority stress experiences among LGBTQ+ people. However, the expression of minority stress on these platforms is linguistically sophisticated and dynamic, which proves to be an obstacle for developing AI algorithms to study minority stress [16]. Cascalheira et al. proposed that linguistic sophistication of expressions of minority stress is threefold: (1) minority stress entails LGBTQ+-specific semantics and pragmatics, (2) psycholinguistic permutations unique to the community, and (3) lexical density (i.e., a lot of words needed to convey minority stress) [2]. Additionally, the intersecting identities of LGBTQ+ individuals (e.g., race/ethnicity, disability status, age, geographical location, etc.) inform the development and usage of novel cultural idioms that, too, convey minority stress [11].

Since expressions of minority stress maintain such linguistic sophistication, capturing minority stress is a great challenge for traditional NLP methods. Recently, large-scale pretraining models have demonstrated effectiveness across a spectrum of NLP tasks [17], [18]. These pretraining models, developed through extensive training on vast unlabeled corpora in an unsupervised fashion, exhibit a capacity to grasp nuanced semantic structures within the text at a considerable scale. While their intrinsic advantages for transductive learning are evident, it is noteworthy that current models designed for transductive text classification [19], [20] often disregard the integration of large-scale pretraining techniques. Consequently, the potential impact of large-scale pretraining on the classification of minority stress remains uncertain within the current landscape of transductive text classification models.

Our research aims to leverage both a transductive approach and large-scale pretraining by simultaneously training Bidirectional Encoder Representations from Transformers (BERT) and graph convolutional networks (GCN) modules. This novel approach to classifying minority stress from social media posts

remains largely unexplored, and the study contributes in the following manners:

- 1) We have designed a hybrid deep learning model, specifically BERT-GCN and RoBERTa-GCN, which leverage large-scale pretraining and graph-based representation learning for understanding relationships between different samples, enhancing predictions related to health disparities.
- 2) The paper aims to enhance predictive capabilities by using a language model that is based on the transformer architecture and graph neural networks. This integration allows the model to capture relationships within the entire corpus, providing a more comprehensive understanding of the data and improving predictions related to health disparities.
- 3) We provide a comprehensive benchmarking of various machine learning and deep learning models on the LGBTQ+ MiSSoM+ dataset [21]. It highlights the performance superiority of BERT-based architectures, particularly RoBERTa-GCN ($F1 = 0.86$), in predicting minority stress labels, offering valuable insights for future research in this domain.

II. RELATED WORK

Social media serves as a vital tool for LGBTQ+ identity exploration and expression, as it allows individuals a platform to present varied facets of themselves across different networks and to ascertain support for their LGBTQ+ identity [22], [23]. While this study focuses on text data analysis to contribute to existing research [11], [24]–[26], the broader potential lies in employing a transductive approach to analyze LGBTQ+ users’ social media posts in media formats other than text (e.g., short-form and long-form video, voice recording, multimedia) to reduce LGBTQ+ health disparities. The multifaceted nature of social media content (text, image, video) suggests an opportunity for innovation in modeling social determinants beyond textual information. Our results can provide immediate value to applied social scientists for interventions in LGBTQ+ health, with the potential to extend into a multimodal prediction of minority stress and related health constructs using non-text data—a promising yet underexamined area in applied big data analytics [24]. Importantly, social media emerges as a unique data source for LGBTQ+ research, as it provides exceptional reliability for studying psychological processes, such as the establishment of belonging and community [27]. Prior studies have demonstrated the effectiveness of deep learning models, such as Bidirectional Encoder Representations from Transformers (BERT) [17], recurrent neural networks (RNNs) [28] with its variants like bidirectional long short-term memory (Bi-LSTM), bidirectional gated recurrent unit (BiGRU), convolutional neural networks (CNNs), and hybrid models, in predicting social determinants of health disparities, including cyberbullying and racism on Twitter [29]. Multichannel CNN achieved a notable F1 score of 0.97 in classifying expressions of general stress on the Twitter dataset, which consists of interview responses from 38 students and stress-related tweets

with certain hashtags. [29]. Hybrid models, such as BERT-LSTM, have shown promise in improving the prediction of misogynistic speech on Twitter (F1 = 0.81) [30].

Past attempts to classify minority stress on Reddit users' text data [11] achieved an F1 score of 0.75 with models such as logistic regression. However, their limited ability to handle sequential data led to drawbacks. Cascalheira et al. [25] introduced a Bi-LSTM but achieved only an F1 score of 0.61, indicating architectural limitations. Later, the combination of BERT-CNN exhibited excellent performance for both composite minority stress and factors of minority stress [11]. While these models have successfully predicted health disparities related to external discrimination, studies focusing on composite minority stress have shown room for improvement.

In contrast to the previous research, our purposed method aims to enhance predictive capabilities by leveraging large-scale pretraining (BERT, RoBERTa) and neural networks. It employs graph neural networks (GNNs) to model relationships between labeled and unlabeled documents, utilizing the similarity between them within the whole corpus. This approach extends beyond previous models to explore the potential of neural networks in understanding the relationships between different samples and enhancing predictions related to health disparities.

III. NOTATIONS AND PRELIMINARIES

A. Graph Definition

Formally, in a graph with n nodes, $G = (V, E)$ is characterized by a set of $V = \{v_1, v_2, v_3, \dots, v_n\}$ and a set of edges $E = \{e_{ij}\}_{i,j=1}^n$ denoting relationships between the nodes. The adjacency matrix A , with dimensions $n \times n$, serves as a fundamental representation of the graph. In unweighted graphs, A_{ij} equals 1 if an edge exists between nodes i and j , and 0 otherwise. Conversely, in weighted graphs, A_{ij} represents the weight of the relationship between nodes i and j , with 0 indicating no relationship. For undirected graphs, A is symmetric ($A_{ij} = A_{ji}$), while for directed graphs, it is asymmetric ($A_{ij} \neq A_{ji}$). In this work, we constructed an undirected heterogeneous graph for the whole corpus. There are two node types (documents and words), and the edges are formed using term frequency-inverse document frequency (TF-IDF) [31] for word-document connections and positive point-wise mutual information (PPMI) for word-word connections.

B. Graph Convolutional Networks (GCN)

Graph Convolutional Networks (GCN) [32], are neural networks that operates through neighborhood aggregation and message passing mechanism. GCNs generate embedding vectors of nodes based on the properties of their neighbor nodes. Additionally, a feature matrix $X \in \mathbb{R}^{n \times m}$ containing n nodes with m -dimensional feature vectors is defined, with each row $x_v \in \mathbb{R}^m$ representing the feature vector for node v . An adjacency matrix A and its associated degree matrix D , where $D_{ii} = \sum_j A_{ij}$, are introduced.

GCN processes information solely from immediate neighbors in one convolutional layer. By stacking multiple GCN layers, it integrates information from multi-hop neighborhoods. For a single-layer GCN, the new k -dimensional node feature matrix $L^{(1)} \in \mathbb{R}^{n \times k}$ is computed as

$$L^{(1)} = \sigma(\tilde{A}XW^{(0)}) \quad (1)$$

where $\tilde{A} = D^{-1/2}AD^{-1/2}$ is the normalized symmetric adjacency matrix, and $W^{(0)} \in \mathbb{R}^{m \times k}$ is a weight matrix. Here, σ represents an activation function such as ReLU ($\sigma(x) = \max(0, x)$). Stacking multiple GCN layers allows the incorporation of higher-order neighborhood information:

$$L^{(i)} = \sigma(\tilde{A}L^{(i-1)}W^{(i)}) \quad (2)$$

where i denotes the layer number, and $L^{(0)} = X$.

IV. METHODS

A. Neural Network Architectures

The BERT-GCN architecture, shown in Figure 1, utilizes a BERT model to initialize representations for document nodes in a heterogeneous Reddit corpus graph. These initialized representations serve as inputs to a Graph Convolutional Network (GCN), where iterative updates based on the graph structures take place. The resulting outputs are treated as final representations for document nodes, which are forwarded to a linear layer and a softmax layer for predictions. Below are detailed descriptions and the roles of each component involved in this architecture.

1) *BERT*: BERT is built upon the transformer architecture, utilizing bidirectional transformer encoder layers to learn contextual word relationships. Each transformer encoder block comprises two sublayers: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network. Residual connections are applied within each sub-layer. The architecture of the BERT transformer encoder block is shown in Figure 2. During pre-training, BERT utilizes masked language model-based and next-sentence prediction objectives [11]. Fine-tuning for specific tasks is achieved by adding task-specific layers to the pre-trained BERT model.

We represent the number of document nodes by N_{doc} , and the number of word nodes by N_{word} that includes both training and test sets. The document node embeddings, $Z_{doc} \in \mathbb{R}^{N_{doc} \times d}$, where d is the dimension of the embedding, and the initial node feature matrix is given by equation (3):

$$X = \begin{pmatrix} Z_{doc} \\ 0 \end{pmatrix}_{(N_{doc}+N_{word}) \times d} \quad (3)$$

Now, the X is fed into a dense layer with a softmax, and the output, Z_B is expressed as below:

$$Z_B = softmax(WX) \quad (4)$$

where W is the weight matrix responsible for transforming the node features.

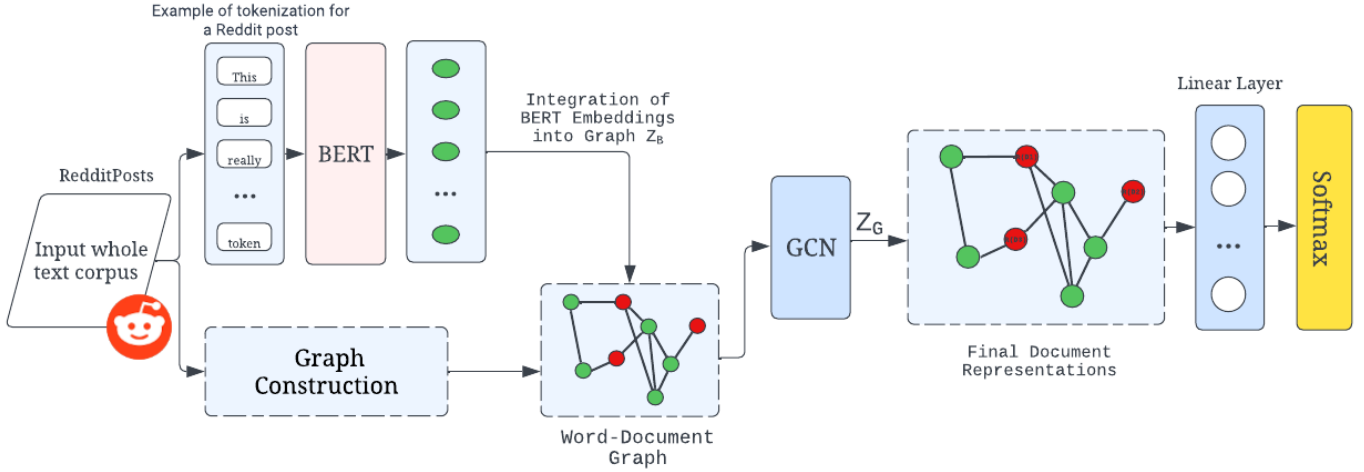


Fig. 1: BERT-GCN Network Architecture in Action.

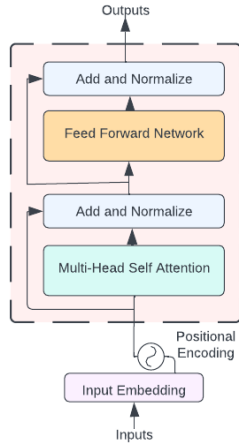


Fig. 2: The BERT layer.

2) *Graph Construction*: We constructed a heterogeneous graph following TextGCN [19], and BertGCN [33] by using TF-IDF for word-document edges and PPMI for word-word edges, we define the edge weights between nodes i and j with equation (5). Figure 3 illustrates the top five words from ten documents using the same equation to create a heterogeneous graph.

$$A_{i,j} = \begin{cases} \text{TF-IDF}(i,j), & \text{if } i \text{ is a document, } j \text{ is a word} \\ \text{PPMI}(i,j), & \text{if } i, j \text{ are words and } i \neq j \\ 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

The Positive Pointwise Mutual Information (PPMI) value of a word pair (i, j) is calculated as:

$$\text{PPMI}(i, j) = \log \frac{p(i, j)}{p(i)p(j)} \quad (6)$$

Where:

- $p(i, j) = \frac{\#W(i, j)}{\#W}$ represents the probability of occurrence of the word pair (i, j) in the corpus.
- $p(i) = \frac{\#W(i)}{\#W}$ is the probability of occurrence of word i in the corpus.
- $\#W(i)$ denotes the number of sliding windows in the corpus that contain the word i .
- $\#W(i, j)$ denotes the number of sliding windows that contain both words i and j .
- $\#W$ is the total number of sliding windows in the corpus.

A positive PPMI value indicates a high semantic correlation between words in the corpus, while a negative PPMI value suggests little to no semantic correlation. Consequently, edges are only added between word pairs with positive PPMI values.

After constructing the graph, X is fed into GCN layers, and the feature matrix of the i -th layer is given by equation (2). The output of GCN is then fed to the softmax layer as is given by equation (7):

$$Z_G = \text{softmax}(\text{gcn}(X, \tilde{A})) \quad (7)$$

where gcn is GCN model as shown in the equation (1) and (2).

3) *BERT-GCN*: BERT-GCN performs the linear interpolation of the prediction from BERT and GCN model, and the final result is shown in equation (8):

$$Z_{Final} = \lambda Z_G + (1 - \lambda) Z_B \quad (8)$$

where $\lambda \in (0, 1)$ is the hyper-parameter that controls two models. When $\lambda = 0$, we fully use the BERT module, and when $\lambda = 1$, we use the GCN component fully. This allows us to balance between two models to achieve the best result for our classification task.

V. EXPERIMENTS

A. Dataset

We used the LGBTQ+ Minority Stress on Social Media advanced version (MiSSoM+) dataset [21]. The dataset is

- Random Forest: An ensemble learning method that combines 100 decision trees trained on TF-IDF features, making predictions through a majority vote or averaging individual tree predictions [34].
- AdaBoost: An ensemble learning algorithm that iteratively trains weak classifiers on weighted versions of the data, assigning higher weights to misclassified samples in each iteration to focus subsequent classifiers on those instances [35]. The baseline used a decision tree classifier with a max depth of 1, 50 estimators, and TF-IDF features.
- Multi-Layer Perceptron (MLP): A feed-forward neural network with multiple layers of fully connected neurons, trained using TF-IDF with default hyperparameters (hidden layer: 100 neurons, adam optimizer, relu activation) [34]
- BiLSTM: The BiLSTM (Bidirectional Long Short-Term Memory) model represents a sophisticated form of recurrent neural network architecture specifically designed for the processing of sequential data [36]. It utilizes a two-layer BiLSTM with a hidden size of 256, with pre-trained Glove word embeddings (dimension 100), and a fully connected layer as in the paper [25].
- BERT-BiGRU: A hybrid model that combines BERT’s contextual embeddings with two layers of Bidirectional Gated Recurrent Units (BiGRU). Using a hidden dimension of 256 in the BiGRU layers, the model captures intricate patterns bidirectionally in input sequences. A linear layer finalizes predictions [11].
- BERT-CNN: BERT is the embedding layer with three convolutional layers consisting of 100 filter counts on each layer and a diverse kernel size of 3, 4, and 5 on each layer [11].
- BERT: A pre-trained natural language processing model that captures contextual relationships in words [17]. The baseline model configuration includes setting the batch size to 32, initializing BERT with *bert-large-uncased*, with a learning rate of 1e-05 and dropout of 0.5.
- RoBERTa: A robustly optimized BERT-based model, designed to enhance NLP tasks such as question answering, text classification, and language modeling [18]. The baseline configuration shares the same settings as BERT for most parameters.

C. Experimental Settings

We conducted all our experiments on the Linux server. The server featured dual Intel Xeon Gold 5220R processors, each comprising 24 cores clocked at 2.20 GHz, with a substantial 35.75 MB cache. Additionally, the server was equipped with four NVIDIA RTX A5000 GPUs, having 24 GB of VRAM per GPU. The training was conducted using PyTorch 1.13.0 with CUDA 11.1 to implement BERT-GCN. We opted for a two-layer GCN as it demonstrated superior performance compared to a single layer. For BERT-GCN and RoBERTa-GCN, we utilized the *bert-large-uncased* and *roberta-large*, respectively, as they outperformed *bert-base-uncased* and *roberta-base*.

The learning rate for the GCN was set to 1e-3, while the BERT module used a learning rate of 1e-5. GAT variants were also trained with the same hyperparameters, and the attention head count was set to 8. Both the BERT-GCN and BERT-GAT were trained for 50 epochs. The source code can be found on GitHub.¹

D. Performance Measures

We used accuracy and weighted F1 score as our performance metric. Accuracy is a fundamental measure in classification tasks that quantifies the overall correctness of predictions by comparing the number of correctly classified instances to the total number of instances. On the other hand, a weighted F1 score is a metric that takes into account both precision and recall across multiple classes, offering a more detailed evaluation that provides a balanced assessment that considers both false positives and false negatives.

E. Results

Table 2 presents a performance comparison of various models on the LGBTQ+ MiSSoM+ dataset, evaluating their accuracy and F1-score for the composite minority stress. Traditional machine learning models, including Naïve Bayes, Logistic Regression, SVM (Linear), Random Forest, AdaBoost, and MLP, demonstrated competitive but varied results, with accuracy ranging from 0.75 to 0.81 and F1-scores from 0.69 to 0.79.

Moving to deep learning models, BiLSTM did not perform as well compared to others, but BERT-based architectures (BERT-BiGRU, BERT-CNN, BERT, RoBERTa) showed significant improvements, with accuracy scores exceeding 0.80 and F1-scores surpassing 0.78. RoBERTa-GCN performed the best among all other classifiers, with an accuracy of 0.8624 and an F1 score of 0.8536, demonstrating its effectiveness in capturing nuanced patterns related to minority stress labels. Compared with BERT and RoBERTa, there is a slight performance boost from BertGCN and RoBERTaGCN, suggesting that the models take advantage of the graph structure.

Figure 5 displays the Receiver Operating Characteristic (ROC) diagram of various classifiers along with their respective Area Under Curve (AUC). The ROC-AUC of RoBERTa-GCN is 0.8735, slightly lower than that of RoBERTa and BERTCNN. While ROC-AUC is an important metric for illustrating classifier accuracy by distinguishing between positive and negative classes across diverse thresholds, the F1 score is more attuned to class imbalance. Our models outperform all classifiers in terms of F1 score, indicating their excellence in correctly identifying instances of the minority class and effectively addressing class imbalance.

F. Ablation Study

λ is the hyperparameter in BERT-GCN that controls the weight assigned to the BERT module and GCN module. Different values of λ lead to different accuracy and F1 Score outcomes. Figure 6 illustrates the classification results on

¹<https://github.com/chapagaisa/transductive>

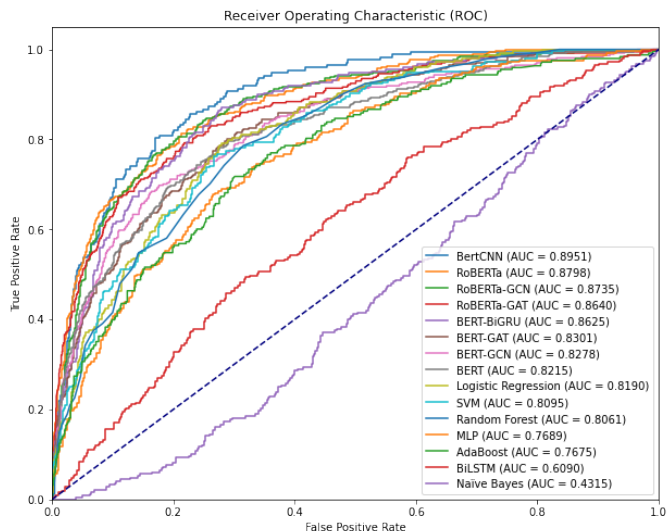


Fig. 5: ROC-AUC of different classifiers.

TABLE II: Performances comparison on minority stress label on the LGBTQ+ MiSSoM+ dataset.

| Model | Accuracy | F1-Score |
|---------------------|---------------|---------------|
| Naïve Bayes | 0.7932 | 0.6926 |
| Logistic Regression | 0.8103 | 0.7832 |
| SVM (Linear) | 0.8103 | 0.7902 |
| Random Forest | 0.7922 | 0.7018 |
| AdaBoost | 0.7902 | 0.7845 |
| MLP | 0.8026 | 0.7865 |
| BiLSTM | 0.7545 | 0.7236 |
| BERT-BiGRU | 0.8432 | 0.8311 |
| BERT-CNN | 0.8608 | 0.8422 |
| BERT | 0.8112 | 0.8106 |
| RoBERTa | 0.8457 | 0.8415 |
| BERT-GCN | 0.8198 | 0.8116 |
| RoBERTa-GCN | 0.8624 | 0.8536 |
| BERT-GAT | 0.8140 | 0.8150 |
| RoBERTa-GAT | 0.8307 | 0.8334 |

the test dataset for the LGBTQ+ MiSSoM+ dataset with different values of λ . When we set λ to 0, we only use the BERT module. When we set λ to 1, we only use the GCN module. We ran several experiments to get optimal values and observed that the best classification results, in terms of both accuracy and F1-score, were achieved when $\lambda = 0.2$. The performance boost with the addition of the GNN component is approximately 1.05% for BERT in terms of accuracy and 0.12% in terms of F1 score, while for RoBERTa, it is approximately 1.97% for accuracy and 1.44% for F1 score. It is commonly acknowledged that pretraining based on BERT holds greater significance compared to graph learning utilizing GNNs. However, it is imperative to note that employing GNN-based graph learning, when appropriately emphasized and integrated into the model architecture, has the potential to significantly enhance overall accuracy.

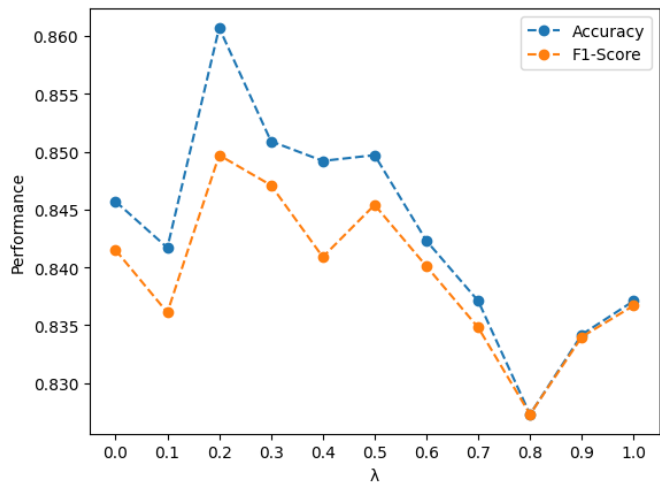


Fig. 6: Test Accuracy and F1 score on LGBTQ+ MiSSoM+ dataset with varying λ values on RoBERTa-GCN.

VI. DISCUSSION AND CONCLUSION

This is the first study to leverage transductive modeling coupled with large-scale pretraining to predict minority stress among LGBTQ+ people who use Reddit. Fifteen different models were tested, running the gamut from traditional supervised machine learning to novel neural network architectures. The transductive-pretrained hybrid model, RoBERTa-GCN, yielded superior performance. Several ethical considerations, implications for digital health interventions, and limitations are discussed to frame our experimental results.

A. Ethical Considerations

In conducting these analyses, we are committed to following Reddit’s User Agreement [37], the American Psychological Association’s Ethical Principles of Psychologists and Code of Conduct [38], and the Association for the Advancement of Artificial Intelligence’s Code of Professional Ethics and Conduct and Diversity Statement [39] throughout our study. We emphasize the need for careful consideration of end-use products when studying the social determinants of LGBTQ+ health inequities, as highlighted in past work [11], [24]–[26].

This study’s key ethical strengths include safeguarding LGBTQ+ Reddit users’ privacy and actively involving LGBTQ+ individuals in all aspects of the research. As part of our commitment to nonmaleficence, our research team is comprised of members of the LGBTQ+ community [11]. We aim to maintain Reddit user anonymity and privacy by restricting access to the datasets, and releasing only the code for analysis reproduction due to potential safety concerns. While scraping Reddit data without explicit consent is common practice, we acknowledge the need for improved communication with LGBTQ+ social media users and AI scientists. Despite not publicly releasing the data, we recognize the enduring potential if unlikely risk of malicious replication of our modeling techniques, which could lead to the identification of LGBTQ+ users and forms of online harassment (e.g. cyberbullying).

Additionally, we acknowledge the mindfulness and nuance necessitated in working with a vulnerable and historically marginalized community, as LGBTQ+ research maintains/elicits unique risks aside of common potential risks. An example of a unique risk is “outing” a closeted individual (i.e., revealing that they are not cisgender and/or heterosexual), as they may face mental decline, diminished safety, or legal consequences in certain regions of the world as a result. Our commitment to minimizing harm and maximizing benefits guides our pursuit of understanding the relationship between minority stress and LGBTQ+ health.

Further, in an effort to contribute thoughtfully to the social sciences, to data science, and to the development of digital health interventions, the MiSSoM+ dataset is well-developed in its accuracy. This advantage can be attributed to the creation of an annotating team constitution and the diverse composition of the annotating team: An annotating team constitution provokes the development of precise definitions for contextualizing LGBTQ+ Reddit user experiences—definitions which were guided by their expertise in LGBTQ+ health and their own lived experiences [40], [41]. Thus, the diverse cultural backgrounds of our team members facilitate accuracy in linguistically capturing cultural nuances. Their heterogeneity offers a variety of academic disciplines and a keen adeptness in LGBTQ+ studies, both of which yield a more accurate dataset.

The current study, due to its absence of direct interactions with human subjects and use of publicly available data, does not qualify for an ethics board revision process. Yet, in consideration of the numerous risks in working with vulnerable populations while employing a computational method, we encourage future AI and social scientists to enhance our current knowledge of confidentiality, expand on the functionality of an ethics board, and continue this commitment of systematic ethics reviews in computational social science research.

Lastly, we acknowledge the timeliness of our models with respect to constantly evolving language in social media. Common LGBTQ+-specific semantics and pragmatics remain fluid and continue to evolve with increased understanding of LGBTQ+ identities among both the LGBTQ+ communities and the general public [42]. We suggest continuous education on sophisticated linguistics of LGBTQ+ communities and timely updates of available datasets to maintain integrity in future similar computational studies. With respect to sophistication in linguistics, we also acknowledge that the annotating team for the MiSSoM+ dataset lacks expertise or lived experiences of certain cultural communities (e.g., lack of understanding of diverse queer culture and living situations on an international level due to limited experiences with diverse geographical locations of members). We encourage development of datasets with further consideration of intersectionality and contextual factors.

B. Implication for Digital Health Interventions

The precise identification of minority stress within social media language holds immense potential for impactful innovations in the healthcare sector and social policy. Our findings

have the capacity to support and instruct the development of personalized digital health interventions for LGBTQ+ individuals. These interventions might include stress-reduction smartphone apps triggering coping strategies upon detecting minority stress. Moreover, due to the uniquely challenging process of LGBTQ+ identity development, our findings intend to support the conception of identity exploration interventions.

Our models could also optimize content delivery in existing health or identity exploration apps, as well as improve safety or crisis resources through assigning AI-informed modules to a user based on detected stress levels. Additionally, the accurate detection of minority stress could prompt timely booster sessions or interventions, potentially reducing minority stress and thereby improving LGBTQ+ individuals’ mental health outcomes and, thus, safety. These applications could significantly benefit LGBTQ+ individuals’ quality of life by addressing and mitigating the impacts of minority stress, ultimately beginning to bridge the inequity of health among heterosexual/cisgender and LGBTQ+ populations.

C. Limitations and Future Research

Although our paper outperforms all the baseline models that were successful in predicting LGBTQ+ Minority Stress, several limitations point toward avenues for future work. Firstly, the use of a pre-trained BERT model without fine-tuning the embedding layer may limit the model’s task-specific performance, urging researchers to explore the potential benefits of fine-tuning BERT for classification tasks. Secondly, the transductive nature of BERT-GCN hampers its agility in processing unseen test documents, indicating a need to investigate and integrate inductive models to enhance adaptability. Thirdly, the GAT module primarily focuses on the 1-degree neighborhood, prompting future work to consider expanding the attention range and incorporating sub-graph attention learning for a more comprehensive understanding. Additionally, the graph construction process, relying solely on document statistics, may be sub-optimal compared to models that can automatically establish edges between nodes. Future research should work into methods for automated edge construction, potentially offering a more robust graph structure and enhancing overall model performance.

From the lens of social science, we applied a human-annotated dataset based on the minority stress theory. While the minority stress theory is considered a well-established theory in contextualizing social stressors LGBTQ+ individuals experience, we acknowledge that the model does not account for intersecting identities (e.g., being a gender/sexual minority while identifying with low-income social status, etc.) [43], [44]. Further research should continue expanding on minority stress with a broader consideration of social identities and their intersectionality.

D. Conclusion

In this study, we used BERT-GCN, a model that leverages the strengths of both pre-trained language models and a graph neural network for predicting LGBTQ+ minority

stress. Employing a real-world social media dataset (LGBTQ+ MiSSoM+), we created a heterogeneous graph and used BERT representation for representing documents as node embeddings. We jointly trained BERT and GCN to improve our performance and the experiment shows the RoBERTa-GCN performs the overall best to all baseline models [11], [24]–[26], improving the prediction of minority stress in comparison to BERT-CNN.

VII. ACKNOWLEDGEMENTS

Cory J. Cascalheira is supported as a RISE Fellow by the National Institutes of Health (R25GM061222). Ryan E. Flinn is supported as a Scholar/Trainee by the following training programs, each of which are funded by the National Institute on Drug Abuse (NIDA): the HIV/AIDS, Substance Abuse, and Trauma Training Program at the University of California, Los Angeles (R25DA035692); the Lifespan/Brown Criminal Justice Research Training Program on Substance Use, HIV, and Comorbidities (R25DA037190); the JEAP Initiative (R24DA051950); and the Brandeis-Harvard SPIRE Center Substance Use Disorder Systems Performance Scholars Program (P30DA035772). Shah Muhammad Hamdi is supported by the CISE and GEO directorates under NSF awards #2301397 and #2305781. Soukaina Filali Boubrahimi is supported by CISE and GEO Directorates under NSF awards #2204363, #2240022, #2301397, and #2305781. Emily M. Lund is a visiting professor at Ewha Women’s University, and resources purchased with Ewha funds were used in the preparation of this manuscript. Jillian R. Scheer is supported by a Mentored Scientist Development Award (K01AA028239-01A1) from the National Institute on Alcohol Abuse and Alcoholism.

REFERENCES

- [1] K. J. Conron, M. J. Mimiaga, and S. J. Landers, “A population-based study of sexual orientation identity and gender differences in adult health,” *American journal of public health*, vol. 100, no. 10, pp. 1953–1960, 2010.
- [2] A. Flentje, N. C. Heck, J. M. Brennan, and I. H. Meyer, “The relationship between minority stress and biological outcomes: A systematic review,” *Journal of Behavioral Medicine*, vol. 43, pp. 673–694, 2020.
- [3] C. Kelleher, “Minority stress and health: Implications for lesbian, gay, bisexual, transgender, and questioning (lgbtq) young people,” *Counseling psychology quarterly*, vol. 22, no. 4, pp. 373–379, 2009.
- [4] I. H. Meyer, “Prejudice, social stress, and mental health in lesbian, gay, and bisexual populations: conceptual issues and research evidence,” *Psychological bulletin*, vol. 129, no. 5, p. 674, 2003.
- [5] J. T. Goldbach, E. E. Tanner-Smith, M. Bagwell, and S. Dunlap, “Minority stress and substance use in sexual minority adolescents: A meta-analysis,” *Prevention Science*, vol. 15, pp. 350–363, 2014.
- [6] “Understanding the well-being of lgbtqi+ populations,” *National Academies of Sciences, Engineering, and Medicine and others*, 2020.
- [7] C. Peele, “Roundup of anti-lgbtq legislation advancing in states across the country,” Human Rights Campaign, 2023, accessed 09-Jul-2023. <https://www.hrc.org/press-releases/roundup-of-anti-lgbtq-legislation-advancing-in-states-across-the-country>.
- [8] Human Rights Watch, “#outlawed: the love that dare not speak its name,” 2023, accessed 09-Jul-2023. https://features.hrw.org/features/features/lgbt_laws.
- [9] L. Lindley and M. P. Galupo, “Gender dysphoria and minority stress: Support for inclusion of gender dysphoria as a proximal stressor,” *Psychology of Sexual Orientation and Gender Diversity*, vol. 7, no. 3, p. 265, 2020.

- [10] G. Di Franco and M. Santurro, “Machine learning, artificial neural networks and social research,” *Quality & quantity*, vol. 55, no. 3, pp. 1007–1025, 2021.
- [11] C. J. Cascalheira, S. Chapagain, R. E. Flinn, Y. Zhao, S. F. Boubrahimi, D. Klooster, A. Gonzalez, E. M. Lund, D. Laprade, J. R. Scheer *et al.*, “Predicting linguistically sophisticated social determinants of health disparities with neural networks: The case of lgbtq+ minority stress,” in *2023 IEEE International Conference on Big Data (BigData)*. IEEE, 2023, pp. 1314–1321.
- [12] D. M. Lazer, A. Pentland, D. J. Watts, S. Aral, S. Athey, N. Contractor, D. Freelon, S. Gonzalez-Bailon, G. King, H. Margetts *et al.*, “Computational social science: Obstacles and opportunities,” *Science*, vol. 369, no. 6507, pp. 1060–1062, 2020.
- [13] T. Nijhawan, G. Attigeri, and T. Ananthakrishna, “Stress detection using natural language processing and machine learning over social interactions,” *Journal of Big Data*, vol. 9, no. 1, pp. 1–24, 2022.
- [14] E. McDermott, “Asking for help online: Lesbian, gay, bisexual and trans youth, self-harm and articulating the ‘failed’ self,” *Health*, vol. 19, no. 6, pp. 561–577, 2015.
- [15] E. McDermott and K. Roen, “Youth on the virtual edge: Researching marginalized sexualities and genders online,” *Qualitative health research*, vol. 22, no. 4, pp. 560–570, 2012.
- [16] E. Brill, R. Florian, J. Henderson, and L. Mangu, “Beyond n-grams: Can linguistic sophistication improve language modeling?” in *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1*, 1998, pp. 186–190.
- [17] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [18] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” *arXiv preprint arXiv:1907.11692*, 2019.
- [19] L. Yao, C. Mao, and Y. Luo, “Graph convolutional networks for text classification,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 7370–7377.
- [20] X. Liu, X. You, X. Zhang, J. Wu, and P. Lv, “Tensor graph convolutional networks for text classification,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 8409–8416.
- [21] C. J. Cascalheira, S. Chapagain, R. E. Flinn, D. Klooster, D. Laprade, Y. Zhao, E. M. Lund, A. Gonzalez, K. Corro, R. Wheatley *et al.*, “The lgbtq+ minority stress on social media (missom) dataset: A labeled dataset for natural language processing and machine learning,” in *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 18, 2024, pp. 1888–1899.
- [22] Y. Cannon, S. Speedlin, J. Avera, D. Robertson, M. Ingram, and A. Prado, “Transition, connection, disconnection, and social media: Examining the digital lived experiences of transgender individuals,” *Journal of LGBT Issues in Counseling*, vol. 11, no. 2, pp. 68–87, 2017.
- [23] O. Haimson, “Social media as social transition machinery,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 2, no. CSCW, pp. 1–21, 2018.
- [24] C. J. Cascalheira, R. E. Flinn, Y. Zhao, D. Klooster, D. Laprade, S. M. Hamdi, J. R. Scheer, A. Gonzalez, E. M. Lund, I. N. Gomez *et al.*, “Models of gender dysphoria using social media data for use in technology-delivered interventions: Machine learning and natural language processing validation study,” *JMIR Formative Research*, vol. 7, no. 1, p. e47256, 2023.
- [25] C. J. Cascalheira, S. M. Hamdi, J. R. Scheer, K. Saha, S. F. Boubrahimi, and M. De Choudhury, “Classifying minority stress disclosure on social media with bidirectional long short-term memory,” in *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 16, 2022, pp. 1373–1377.
- [26] K. Saha, S. C. Kim, M. D. Reddy, A. J. Carter, E. Sharma, O. L. Haimson, and M. De Choudhury, “The language of lgbtq+ minority stress experiences on social media,” *Proceedings of the ACM on human-computer interaction*, vol. 3, no. CSCW, pp. 1–22, 2019.
- [27] E. Formby, *Exploring LGBT spaces and communities: Contrasting identities, belongings and wellbeing*. Taylor & Francis, 2017.
- [28] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [29] B. Shaw, S. Saha, S. K. Mishra, and A. Ghosh, “Investigations in psychological stress detection from social media text using deep architectures,”

- in *2022 26th International Conference on Pattern Recognition (ICPR)*. IEEE, 2022, pp. 1614–1620.
- [30] R. S. Angeline, D. Nurjanah, and H. Nurrahmi, “Misogyny speech detection using long short-term memory and bert embeddings,” in *2022 5th International Conference on Information and Communications Technology (ICOIACT)*. IEEE, 2022, pp. 155–159.
- [31] G. Salton and C. Buckley, “Term-weighting approaches in automatic text retrieval,” *Information processing & management*, vol. 24, no. 5, pp. 513–523, 1988.
- [32] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *arXiv preprint arXiv:1609.02907*, 2016.
- [33] Y. Lin, Y. Meng, X. Sun, Q. Han, K. Kuang, J. Li, and F. Wu, “Bertgcn: Transductive text classification by combining gcn and bert,” *arXiv preprint arXiv:2105.05727*, 2021.
- [34] S. Raschka, Y. H. Liu, V. Mirjalili, and D. Dzhulgakov, *Machine Learning with PyTorch and Scikit-Learn: Develop machine learning and deep learning models with Python*. Packt Publishing Ltd, 2022.
- [35] Y. Freund, R. E. Schapire *et al.*, “Experiments with a new boosting algorithm,” in *icml*, vol. 96. Citeseer, 1996, pp. 148–156.
- [36] M. Schuster and K. K. Paliwal, “Bidirectional recurrent neural networks,” *IEEE transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.
- [37] Reddit Inc. (2023) Reddit user agreement. Accessed on: November 21, 2024. [Online]. Available: <https://www.redditinc.com/policies/user-agreement-september-25-2023>
- [38] American Psychological Association. (Year of the latest revision) Apa code of ethics. Accessed on: November 21, 2024. [Online]. Available: <https://www.apa.org/ethics/code/>
- [39] Association for the Advancement of Artificial Intelligence. (Year of the latest revision) Aaai ethics and diversity statement. Accessed on: November 21, 2024. [Online]. Available: <https://aaai.org/about-aaai/ethics-and-diversity/#diversity-statement>
- [40] K. Saha, S. C. Kim, M. D. Reddy, A. J. Carter, E. Sharma, O. L. Haimson, and M. De Choudhury, “The language of lgbtq+ minority stress experiences on social media,” *Proceedings of the ACM on human-computer interaction*, vol. 3, no. CSCW, pp. 1–22, 2019.
- [41] J. Gray and M. Cooke, “Intersectionality, language and queer lives,” *Gender and Language*, vol. 12, no. 4, pp. 401–415, 2018.
- [42] F. G. Sicurella, “The approach that dares speak its name: queer and the problem of ‘big nouns’ in the language of academia,” *Gender and Language*, vol. 10, no. 1, pp. 73–84, 2016.
- [43] M. Rivas-Koehl, D. Rivas-Koehl, and S. McNeil Smith, “The temporal intersectional minority stress model: Reimagining minority stress theory,” *Journal of Family Theory & Review*, vol. 15, no. 4, pp. 706–726, 2023.
- [44] N. Noyola, M. Sánchez, and E. V. Cardemil, “Minority stress and coping among sexual diverse latinxs.” *Journal of Latinx Psychology*, vol. 8, no. 1, p. 58, 2020.