

BullyBlocker: Towards an Interdisciplinary Approach to Identify Cyberbullying*

Yasin N. Silva
Arizona State University
Glendale, AZ, USA
ysilva@asu.edu

Deborah Hall
Arizona State University
Glendale, AZ, USA
d.hall@asu.edu

Christopher Rich
Arizona State University
Glendale, AZ, USA
cdrich2@asu.edu

Abstract— Cyberbullying is the deliberate use of online digital media to communicate false, embarrassing, or hostile information about another person. It is the most common online risk for adolescents, yet well over half of young people do not tell their parents when it occurs. While there have been many studies about the nature and prevalence of cyberbullying, there have been relatively few in the area of automated identification of cyberbullying that integrate findings from computer science and psychology. The goal of our work is thus to adopt an interdisciplinary approach to develop an automated model for identifying and measuring the degree of cyberbullying in social networking sites, and a Facebook app, built on this model, that notifies parents about the likelihood that their adolescent is a cyberbullying victim. This paper describes the challenges associated with building a computer model for cyberbullying identification, presents key results from psychology research that can be used to inform such a model, introduces a holistic model and mobile app design for cyberbullying identification, presents a novel evaluation framework for assessing the effectiveness of the identification model, and highlights crucial areas of future work. Importantly, the proposed model—which can be applied to other social networking sites—is the first that we know of to bridge computer science and psychology to address this timely problem.

Keywords—cyberbullying; automated identification; social networks; Facebook; psychology; cyberbullying factors; vulnerability factors

I. INTRODUCTION

Over half of adolescents have been bullied online, about the same number have engaged in cyberbullying, and more than one in three young people have experienced cyber-threats online [1]. Cyberbullying can take multiple forms such as posting hurtful or threatening messages online, spreading rumors on social networking sites, taking and posting unflattering pictures of a person, or circulating sexually suggestive pictures or messages about a person. The consequences of cyberbullying—which can include anxiety, depression, and even suicide [44]—are detrimental on both an individual and a societal level. Despite the growing prevalence of cyberbullying, well over half of young people do not tell their parents when cyberbullying occurs [1]. Moreover, while there have been many studies about the nature and prevalence of cyberbullying and even a few on cyberbullying measures for mobile and chat-based venues, there has been little work on the design and implementation of automated models and tools to identify cyberbullying that bridges findings from computer science and psychology.

The closest relevant work, which uses machine learning based models to identify cyberbullying in social networking sites, e.g., [10, 11, 12, 13, 14, 15], has several crucial limitations. For instance, these previously proposed models focus on the identification of cyberbullying in a single message or picture. While insults may appear in single messages, the models do not consider that in many cases cyberbullying occurs as a sequence of insulting or harassing events. The accuracy and reliability of a model for identifying cyberbullying will likely increase to the extent that the number and frequency of insulting messages are taken into account. Moreover, for the most part, previous work has not integrated critical findings from psychology research on the nature of cyberbullying, including risk factors and negative outcomes of cyberbullying and patterns of cyberbullying over time. Additionally, to our knowledge, previous work has yet to address key related problems like the automated generation of a list of anti-bullying resources for parents (e.g., websites, hotlines) based on specific characteristics of the cyberbullying that an adolescent is experiencing. Finally, few of the previous papers have sought to integrate a cyberbullying identification model into an actual app that can help parents and potential victims.

The goal of our work is thus to adopt an interdisciplinary framework to study, design, and implement a model to identify cyberbullying among adolescents in social networking sites. The initial model has been used to build BullyBlocker, a mobile app that identifies cyberbullying on Facebook and generates a customized list of anti-bullying resources for parents. BullyBlocker alerts parents of potential cyberbullying instances by providing them with a *Bullying Rank* (BR) that estimates the probability that their child is being bullied, allowing them to identify and address this form of online aggression. While Facebook is the most common social media platform for teens [2], the principles and design used in BullyBlocker can also be applied to other social networking platforms. Furthermore, similar models could also be used to identify a broad range of negative outcomes that may result from or be exacerbated by cyberbullying, such as depression, substance use, or self-destructive behavior. The primary contributions of this paper are:

- The design of a holistic model to identify cyberbullying that builds on previous research findings on cyberbullying in adolescents. The proposed model integrates crucial findings from psychological research on predictors of cyberbullying, as well as the relative strength of various predictors and

* This work was supported by National Science Foundation Award # 1719722, the Dion Initiative for Child Well-Being and Bullying Prevention and ASU NCUIRE awards.

temporal aspects of cyberbullying. Furthermore, rather than focusing on cyberbullying prediction or classification for a single message or picture, the proposed model considers streams or bursts of messages in conjunction with information from adolescents' social media profiles.

- The design and implementation of a mobile app (BullyBlocker) that uses the proposed model to identify cyberbullying in Facebook. We present the implemented app architecture and a discussion of the key software components.
- The introduction of a module that provides a customized list of parent/victim resources using information pertaining to the nature of specific cyberbullying instances.
- The development of an innovative framework for evaluating the accuracy of holistic cyberbullying identification models that involves a simulated social network with content from real world cyberbullying interactions and a comparison of the results of automated identification models with human assessments of cyberbullying likelihood. We present the results of evaluating the BullyBlocker identification model using this framework.
- The identification of challenges and opportunities to integrate the latest results from psychology and social network data analysis to address a problem of great social impact. We also highlight areas that warrant further study within psychology.
- The public and no-cost availability of the BullyBlocker app (1.0) in the Apple App Store [33].
- The public availability of the source code of the evaluation framework as well as the real-world datasets that it uses to generate the social network interactions [45].

This paper builds on and marks a significant extension of an abstract that appeared in [7]. In particular, the current paper expands our previous work by integrating: (1) an extended identification model that includes (a) multiple new cyberbullying factors—including insulting video comments, race and ethnicity, frequency of internet use, past bullying experiences, sexual orientation, mental health history, disciplinary problems, and substance use, (b) the use of correlation coefficients (r) identified in meta-analytic reviews of cyberbullying research to increase the accuracy of the weights assigned to the different factors in the model, and (c) the use of psychology research to estimate the temporal parameters in our model; (2) a more detailed description of our cyberbullying identification model; (3) an innovative evaluation framework for holistic models that integrate profile features and message streams; (4) an extended architecture diagram, (5) an in-depth explanation of the motivation behind our project that highlights how the present research addresses several important limitations of previous work; (6) an expanded and more detailed review of previous machine learning based cyberbullying identification models and relevant empirical findings in psychology; (7) the design guidelines for two additional identification models as future improvements; and (8) a detailed presentation of how additional results in psychology will be integrated into these identification models.

The remainder of this paper is organized as follows. Section II presents the background and related work, Section III

describes the proposed model and app design guidelines, Section IV discusses our novel evaluation framework for assessing the accuracy of holistic cyberbullying identification models and the results of our evaluation of the BullyBlocker identification model, Section V describes some key paths for future work, and Section VI concludes the paper.

II. BACKGROUND AND RELATED WORK

Prior contributions from computer science that propose models for identifying cyberbullying [e.g., 10, 11, 12, 13, 14, 15] rely primarily on machine learning classification or prediction models that analyze text features (e.g., comments and posts) [10, 13], externally annotated images or videos [11], or both [12]. Yet, the accuracy of these methods, as highlighted in [14], remains limited. Moreover, a major drawback of cyberbullying research within computer science is that it has largely ignored relevant psychology research findings. That is, while the results of studies incorporating some social network features [14] and user demographics [15] to improve accuracy are promising related efforts, a core open challenge is how to effectively integrate insights from psychology to improve automated identification models.

To develop an effective cyberbullying identification model, we are thus building on empirical work from within psychology. Although findings regarding the prevalence, determinants, and even the definition of cyberbullying vary somewhat within the psychology literature [16, 18], examination of the cumulative findings, across multiple studies, reveals some important trends and emerging areas of agreement [16, 27, 28, 35]. For example, seeming inconsistencies across studies in prevalence rates of cyberbullying among different age groups have more recently shed light on a potential curvilinear relation between age and rates of cyberbullying in children and adolescents [18, 19], with the prevalence of cyberbullying victimization at its highest during the later years of middle school [19]. Findings with respect to gender and cyberbullying are similarly complex, with some studies indicating that adolescent girls experience higher rates of cyberbullying than adolescent boys [4, 20, 21, 23] and others finding no systematic gender differences [19, 22]. There is greater consensus in the findings that cyberbullying is a stronger predictor of depression in adolescent girls than boys [16] and that girls report a stronger negative emotional impact of cyberbullying than boys [4], highlighting the value of cyberbullying identification models that draw on critical insights from psychology research.

Some of the most informative findings from within psychology come from meta-analytic reviews, which reveal the average effect of various risk factors across multiple studies and different research teams, weighted by characteristics that contribute to the accuracy and overall quality of a specific research study (i.e., the size of the research sample). For instance, one of the most robust predictors of cyberbullying victimization among teens is whether they have also been victims of traditional bullying in the past [16, 35]. In two separate meta-analyses [16, 35], previous history as a victim of traditional (i.e., offline, face-to-face) bullying emerged as the strongest of numerous risk factors for cyberbullying victimization. Furthermore, whereas symptoms of psychological distress and behavior problems are frequently

examined as outcomes associated with cyberbullying victimization, the psychology literature also indicates that they are a robust *predictor* of cyberbullying victimization. Conclusions about the direction of causality between psychological distress, behavior problems, and cyberbullying cannot be made, given the correlational and largely cross-sectional nature of the data; in fact, it seems likely that a reciprocal relation exists, such that symptoms and indicators of poorer mental health increase adolescents' risk of being targeted by cyberbullies, which then contributes to increased psychological distress and the onset of new symptoms or negative behaviors. A crucial insight, however, is that adolescents' previous history of mental health and behavioral challenges are important factors to include in an identification model.

To begin incorporating these results into our automated identification model, we divided the set of bullying factors into *warning signs* (quantifiable measures like the number of insulting wall posts) and *vulnerability factors* (risk factors and circumstances that may increase the probability of experiencing cyberbullying). The current BullyBlocker model identifies warning signs and vulnerability factors by (1) analyzing the interaction of an adolescent with his or her network through wall posts and picture or video comments, (2) obtaining information from the adolescent's Facebook profile, like age, gender, and schools attended, and (3) obtaining relevant information about additional vulnerability factors directly from parent users of the BullyBlocker app. We mention next specific research findings that provide the framework for the current version of our model.

The current BullyBlocker identification model considers as warning signs: the number of insulting wall posts, the number of embarrassing or insulting comments on photos, and the number of embarrassing or insulting comments on videos that an adolescent has received in the last 90 days. (All posts or comments written by the potential victim are excluded.) In the absence of prior work examining the relative strength of various warning signs as indicators of cyberbullying, we use, whenever available, the correlation coefficients identified in meta-analytic reviews of cyberbullying research to compute the weights of the factors. When this is not possible, we include estimated weights that will be modified in subsequent models.

Data pertaining to vulnerability factors are collected in two ways. First, some information is extracted from the adolescent's social media profile. This includes, for example, the adolescent's age and gender. Nuances in the findings across studies concerning rates of cyberbullying victimization among different age groups and between males and females underscore the importance of including these demographic factors in our identification model, albeit with relatively lower assigned weights, to better understand how they may contribute to cyberbullying risk. Because cyberbullying victims are typically adolescents on the "fringe" of various peer groups [3], two factors that can contribute to a teen's fringe status—whether the teen has recently relocated to a new neighborhood or a new school—will also be mined from adolescents' Facebook profile and included as vulnerability factors. Finally, when relevant information is provided, sexual orientation will be included as a vulnerability factor based on multiple studies indicating that non-heterosexual adolescents are more likely to experience

cyberbullying than their heterosexual peers [36]. In sum, social media data can be mined to identify the extent to which an adolescent possesses any of the vulnerability factors described above.

Additional information that parents can provide about their teens—by filling out a brief in-app "user profile" survey for the adolescent(s) they wish to monitor—will further increase the accuracy of the identification model by allowing us to include several additional vulnerability factors. Specifically, based on the two key meta-analyses described above [16, 35], parents will be asked about the frequency of their teen's internet use and whether, to their knowledge, there is any prior history of being bullied. Parents can also indicate in the user profile whether their teen has a known history of: (1) internalizing problems—symptoms of psychological distress that an individual directs inward, including anxiety, depression, and low self-esteem, and/or (2) externalizing problems—problem behaviors that are directed outward, towards an individual's environment, including disciplinary problems (e.g., suspension or expulsion from a school) and substance use. Averaging across multiple studies, these correlated yet distinct psychological factors have been identified as reliable predictors of cyberbullying victimization among adolescents [16, 35].

To summarize, the current BullyBlocker identification model considers as warning signs: the number of insulting wall posts, the number of embarrassing or insulting comments on photos, and the number of embarrassing or insulting comments on videos. As vulnerability factors, the model considers: race, age, gender, sexual orientation, having recently moved to a new neighborhood or a new school, past bullying history, frequency of internet use, internalizing problems (depression, anxiety, and low self-esteem), and externalizing problems (disciplinary problems and substance abuse). The specific way in which the identified factors are used in our model is described in Section III.

III. AUTOMATED IDENTIFICATION OF CYBERBULLYING

The main design components of the BullyBlocker app are presented in Fig. 1. The app is designed for use by the parent or guardian of an adolescent, who will be required to enter the Facebook login information of the adolescent being monitored.

The *Data Collection Module* uses the adolescent's login information to retrieve all of the required information from Facebook, i.e., the stream of recent wall posts by other Facebook users, the streams of photo and video comments made by other users, and user profile information such as age, gender, recently attended schools, and home location. An important component of this module is the *Query Completion Tracker*, which keeps track of the asynchronous and parallel information requests sent to Facebook. The stream of messages retrieved from Facebook is organized using a tree structure of pages that contains different subsets of the root level messages, their comments, and sub-comments. The tracker guarantees that all of the pages associated to a given stream are properly traversed and processed. This module also makes sure that the frequency of requests sent to Facebook is under the threshold established by Facebook. In addition to obtaining the message/comment streams and the Facebook profile information, this module also

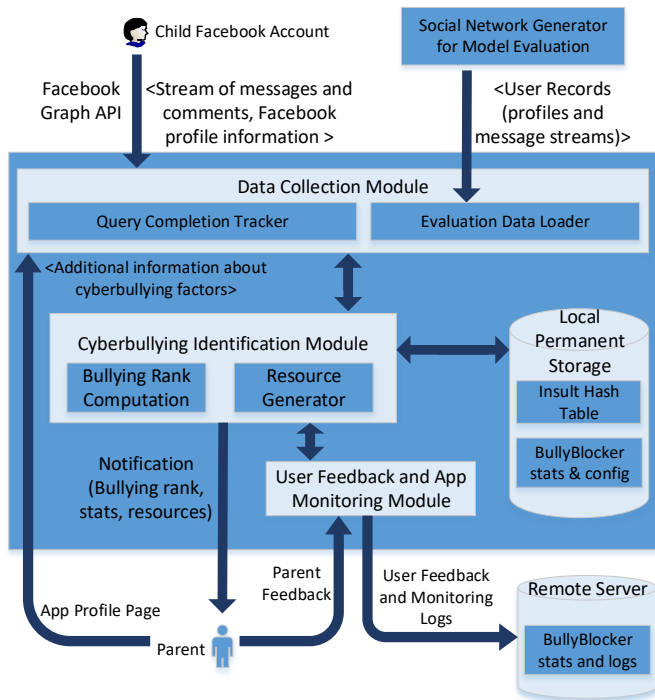


Fig. 1. BullyBlocker Architecture

collects information directly from the parent through the app using a brief survey, as shown in Fig. 3. This survey collects information related to additional vulnerability factors such as ethnicity, race, frequency of internet use, and previous bullying history.

To allow for the evaluation of the BullyBlocker identification model, we created the *Evaluation Data Loader*, a component of the Data Collection Module—used only during the evaluation process—to enable the loading of social network data from generated datasets instead of from actual users (see Section IV). The evaluation dataset is composed of user records, with each record containing a user profile and associated message stream. This dataset is generated by the *Social Network Generator* component using information from real-world cyberbullying interactions. The details of the Social Network Generator and the evaluation process are presented in Section IV.A.

The *Cyberbullying Identification Module* then uses the retrieved data to estimate the likelihood that an adolescent is a victim of cyberbullying on Facebook. To this end, the application computes a *Bullying Rank* expression, based on the identified warning signs and vulnerability factors, that aims to represent the probability that an adolescent is experiencing cyberbullying. The Bullying Rank is used to normalize the intensity of cyberbullying and to simplify the results presented to the parent. Fig. 2 shows the general approach for computing the Bullying Rank.

As shown in Fig. 2, the Bullying Rank (BR) is computed based on the values of Warning Signs (WS) and Vulnerability Factors (VF). Each part is given an appropriate weight such that the range of the BR is [0,100]. The Bullying Rank can fall into

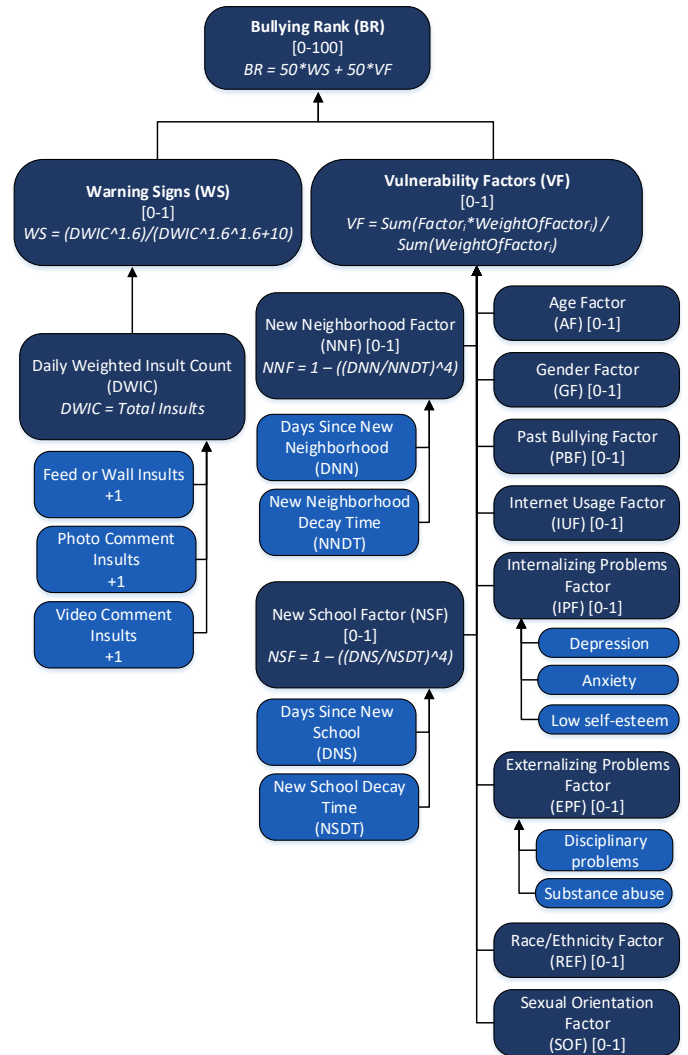


Fig. 2. Bullying Rank Factors

any of three pre-defined levels, with the respective intervals: low risk [0,33], moderate risk [34,66], and severe risk [67,100].

The Bullying Rank, together with several aggregated measures such as the number of insulting wall posts, the number of insults in photo and video comments, the number of potential bullies, and the time range of the analysis, is generated by the Cyberbullying Identification Module and loaded in the results page of the BullyBlocker app. This module also generates a customized list of resources, including websites and hotlines, that direct parents to national and local organizations that provide information about ways to address current and prevent future instances of cyberbullying. Some of the information processed and generated by this module is stored in the mobile device's permanent storage. For instance, the app records the previously computed values of the Bullying Rank and its various components and the most recent dates on which the adolescent moved to a new neighborhood or school.

The *User Feedback and App Monitoring Module* enables parents to submit a brief survey about their perceptions of the accuracy of the app. This module also includes a monitoring

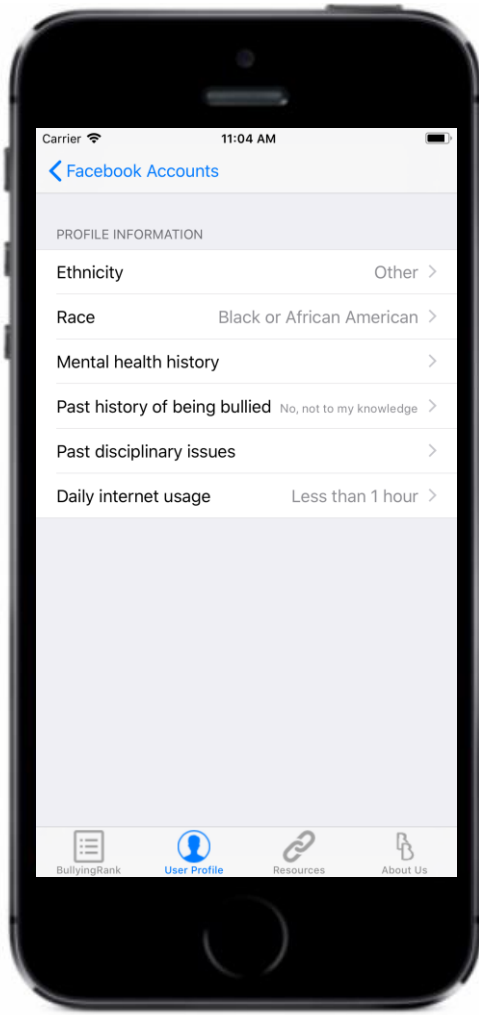


Fig. 3. User Profile Interface in the BullyBlocker App. (Note: To avoid the use of psychological terms with which parents may be unfamiliar, the internalizing problems factor is referred to as “Mental health history” and the externalizing problems factor is referred to as “Past disciplinary issues” in the in-app user profile survey.)

component that can be enabled during the app test phase with a specific set of test users to log the computed Bullying Rank values, final values of the cyberbullying factors in Fig. 2, and an encrypted version of the user IDs. The collected data is sent to a remote web and database server for app monitoring and assessment purposes.

A. Measuring Warning Signs

The Warning Signs (WS) component aims to quantify the amount of insulting content received by the monitored adolescent. This component is included to account for the *Group Effect*, as identified in [6], where the number of insults increases the severity of perceived victimization.

As shown in Fig. 2, this component is computed based on the number of feed (wall) insults, the number of insulting photo comments, and the number of insulting video comments received by the adolescent during the last N days (currently $N = 90$). To decide if a message is of insulting nature, we analyze the

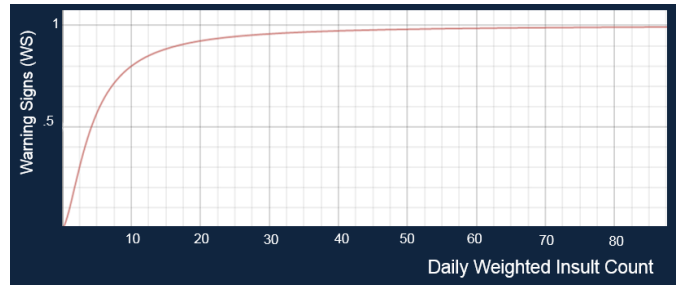


Fig. 4. Warning Signs Vs. Daily Insult Count

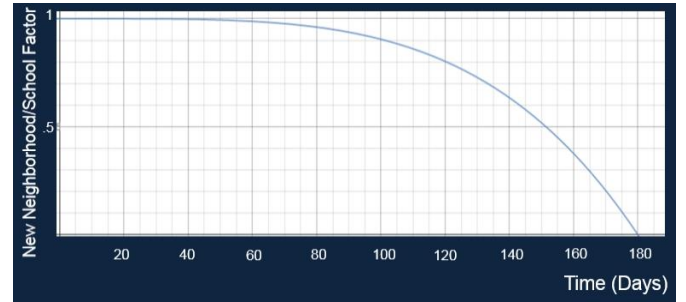


Fig. 5. New Neighborhood/School Factor Vs. Time

content in the message by running hash-based lookup operations on a dictionary of insults and their variations (variations are considered because, in many instances, adolescents use them instead of the original insulting words). The Warning Signs component could also consider the number of insulting private messages received by the potential victim when accessing this data is allowed by the social network’s query API. This is currently not the case with Facebook. The raw insult counts are combined into the *Daily Weighted Insult Count* (DWIC) by applying equal weights to all sub-components (feed, photo and video insults) and computing the average value per day. The DWIC value is then normalized to be in the $[0,1]$ range. Rather than applying uniform scaling, we use a function that assigns higher weights to initial insults, given that after a certain large value of daily insults (~ 30), additional insults tend to have a minimal effect. The function (specified in the Warning Signs box in Fig. 2) is plotted in Fig. 4. Observe that the X-axis of this graph corresponds to the values of DWIC and the Y-axis is the WS value computed using the equation previously referenced. As shown in this figure, going from 5 to 10 daily insults generates a larger increment in the function value than going from 80 to 85 insults. The current model uses 90 as the value of N (number of days). The meta-analysis of cyberbullying by Baldry et al. [34] shows that this time frame is commonly used or contains the reference period used in previous psychology studies.

While all insults receive equal weight in the current model, the model can easily be extended to assign different weights to different types of insults and to increase the weight of an insult based on properties such as the number of people who “liked” the insulting message or the message’s recency. We were not able, however, to identify studies that have directly addressed this aspect of cyberbullying.

Factor	Details	Weight
New School	# days in a new school	0.10
New Neighborhood	# days in a new neighborhood	0.10
Age	Applied if value is 11-16	0.04
Gender	Applied if value is female	0.12
Race/ Ethnicity	Applied if race is non-white or if ethnicity is Hispanic/Latino	0.02
Sexual Orientation	Applied if self-identified as LGBTQ	0.29
Past Bullying	Applied if user experienced bullying in last 1 month, 1-2 months, more than 2 months	0.42
Daily Internet Use	Considers ranges <1h, 1h-3h, 4h-6h, >6h	0.17
Internalizing Problems	Considers history of depression, anxiety, low self-esteem	0.28
Externalizing Problems	Considers history of disciplinary issues or substance use	0.21

Fig. 6. Details and Weights of Vulnerability Factors

B. Measuring Vulnerability

The Vulnerability Factors (VF) component aims to quantify the level of vulnerability of the monitored adolescent. As shown in Fig. 2, this component is computed based on the following ten factors: age, gender, sexual orientation, days since transition to a new neighborhood, days since transition to a new school, race and ethnicity, prior history of being bullied, frequency of internet use, internalizing problems (i.e., “mental health history” in the in-app survey), and externalizing problems (i.e., “past disciplinary issues” in the in-app survey). The value of each factor is in the range of [0,1]. Fig. 6 shows the details and weights of each factor. Intuitively, each factor should have a different weight in the identification model, given variability in the strength of the relation between each factor and cyberbullying risk. To capture this property, we have assigned weights primarily based on the correlation coefficients identified in previous comprehensive meta-analytic reviews. Specifically, the weights for age, prior history of being bullied, internalizing problems, and externalizing problems were based on the correlation coefficients identified in [16] and [35]. The weight for frequency of internet use was based on the correlation coefficient identified in [16] and the weights for gender and race/ethnicity were based on the correlation coefficients identified in [35]. The weight assigned to the sexual orientation factor was based on the meta-analytic effect in [36], which was reported as an odds ratio indicating that gay, lesbian, and bisexual teens were 2.24 times as likely to be a victim of

cyberbullying as heterosexual teens. We performed a transformation (see [42]) to convert the odds ratio to a correlation coefficient. In the absence of published research syntheses examining the relation of cyberbullying risk with the new neighborhood and new school factors, these factors were assigned initial estimated weights in the current identification model.

For the Age and Gender Factors (AF, GF), the model assigns a value of 1 when the age of the potential victim is between 11 and 16 years old, and the gender is female, respectively. The Race/Ethnicity Factor (REF) is set to 1 if race is non-white or if ethnicity is Hispanic/Latino. The Sexual Orientation Factor (SOF) is set to 1 when the potential victim is identified as a member of the LGBTQ group (combining the gender and “interested in” properties of the Facebook profile). The Past Bullying Factor (PBF) receives various values based on the recency of the previous history of being bullied. Similarly, the Internet Use Factor (IUF) receives different values based on the frequency of internet use. The Internalizing and Externalizing Problems Factors (IPF, EPF) are assigned different values based on the number of sub-factors (listed in the corresponding rows of Fig. 6) identified in the in-app survey. The New Neighborhood Factor (NNF) and New School Factor (NSF) weights are based on the number of days since the adolescent moved to a new neighborhood or school. The effect of these components is assumed to change over time, such that the effect should be higher if the adolescent moved recently. To represent this, the model uses a function, specified in the NNF and NSF boxes in Fig. 2 and plotted in Fig. 5, that generates a NNF or NSF value that starts at 1 (the day the adolescent moves to a new neighborhood or school) and decreases over time until it reaches a value of 0 (when the number of days is equal to a parameter value, e.g., New Neighborhood Decay Time).

The total value of Vulnerability Factors is computed by multiplying each factor by its weight and then re-scaling the result to be in the range of [0-1]. The design to compute the VF component can be extended in the future to integrate new empirical findings from both the psychology and computer science literatures. For instance, measures like the number of new friends added to an adolescent’s social network since moving to a new neighborhood or school could be incorporated into the model.

Figures 7 and 8 show the interface of the current version of the BullyBlocker app. These figures show two screenshots that correspond to the result pages generated for two monitored adolescents, Maria and Constance Winther. While the Bullying Rank of Maria is relatively low (i.e., 6), the one for Constance is significantly higher (i.e., 83) due to much higher WS and VF values.

C. Smart Generation of Parent/Victim Resources

One of the challenges that parents face upon learning that their adolescent is a victim of cyberbullying is knowing how to respond effectively; that is, how to help curb the bullying attacks, prevent future instances of cyberbullying, and provide the necessary psychological and emotional support. Critical to an effective resolution is parents’ ability to locate appropriate resources (e.g., anti-bullying organizations, hotlines, literature,

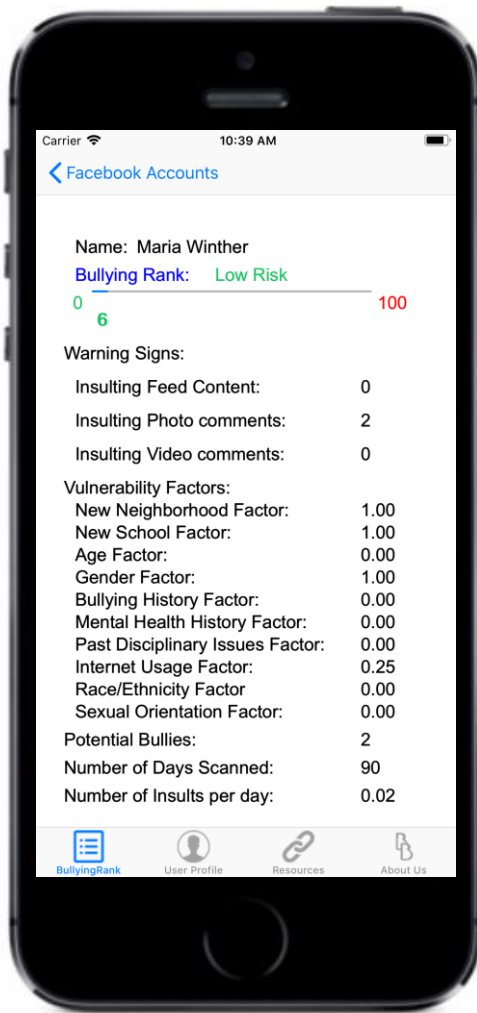


Fig. 7. BullyBlocker Results – Low Risk

etc.), yet finding the most relevant resources for a specific instance of cyberbullying can be an overwhelming task.

BullyBlocker aims to be an effective tool to address this problem by generating a customized list of resources that is tailored to the unique circumstances surrounding the bullying attack(s) and the individual needs of the adolescent. To this end, the app maintains an internal compact representation of the various factors that have been identified for the specific user being analyzed. The app also maintains a robust list of anti-bullying resources (local and national websites and hotlines) annotated with the specific groups targeted by each resource, e.g., racial and ethnic minorities, girls, members of the LGBTQ community, etc. After completing the computation of the Bullying Rank, the app uses the information pertaining to the identified factors to rank the list of resources by potential relevance. As shown in Fig. 9, the list of most pertinent resources is presented at the top of the anti-bullying resources page. In this example, the vulnerability factors that were activated are Race/Ethnicity (the monitored adolescent identifies as Hispanic) and Externalizing Problems (history of substance use was reported). Considering this information, the app recommends a customized list of resources that includes links to the websites for Drug Rehab [37], Help your Teen Now

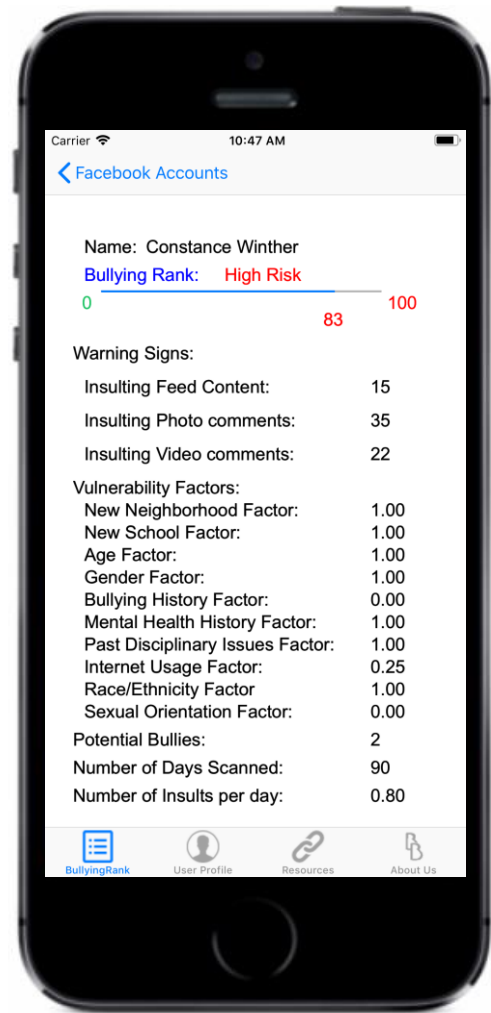


Fig. 8. BullyBlocker Results – High Risk

[38], and Drug Abuse [39]; and race/ethnicity-based resources such as Beyond Bullying [40].

IV. EVALUATION OF HOLISTIC BULLYBLOCKER IDENTIFICATION MODELS

As mentioned previously, one of the limitations of most of the prior work in this area is that the proposed models focus on the identification of cyberbullying in a single post, message, or picture. In many real-world cases, however, cyberbullying involves a repeated sequence of insults, or insult bursts [41, 43]. To our knowledge, our model is one of the first that seeks to identify cyberbullying by considering multiple streams of messages (e.g., wall posts, picture/video comments). Furthermore, our model integrates a set of vulnerability factors based on extensive empirical work in psychology and related social science fields. Finally, the evaluation of holistic models like the one presented in this paper requires a more comprehensive evaluation framework than those needed for simpler models. Yet, a key challenge stems from the difficulty of obtaining datasets that contain all of the required information. That is, whereas the generation and labeling of datasets is relatively simple with models that aim to predict cyberbullying

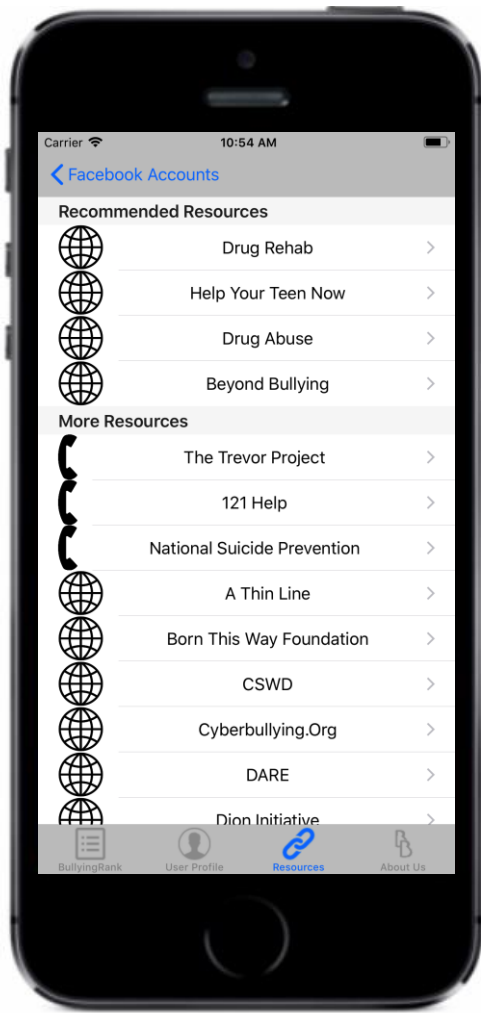


Fig. 9. BullyBlocker Personalized Anti-bullying Resources

in a single message, the complexity increases when considering message streams and multiple vulnerability factors.

To address these challenges, we propose the evaluation framework for holistic cyberbullying identification models depicted in Fig. 10. The goal of this framework is to identify whether or not cyberbullying has occurred by considering the entire user profile as well as the streams of messages received by the user. In this section, we discuss the proposed evaluation framework and its underlying hybrid social network in detail, and present the results obtained using this framework to evaluate the BullyBlocker identification model.

A. Evaluation Framework

The proposed evaluation framework uses a hybrid social network generator to create realistic datasets that are provided as input into the app and also used later in the human evaluation phase.

Hybrid Social Network Generator. This component of the framework generates a test social network composed of synthetic users and real-world interactions (i.e., messages).

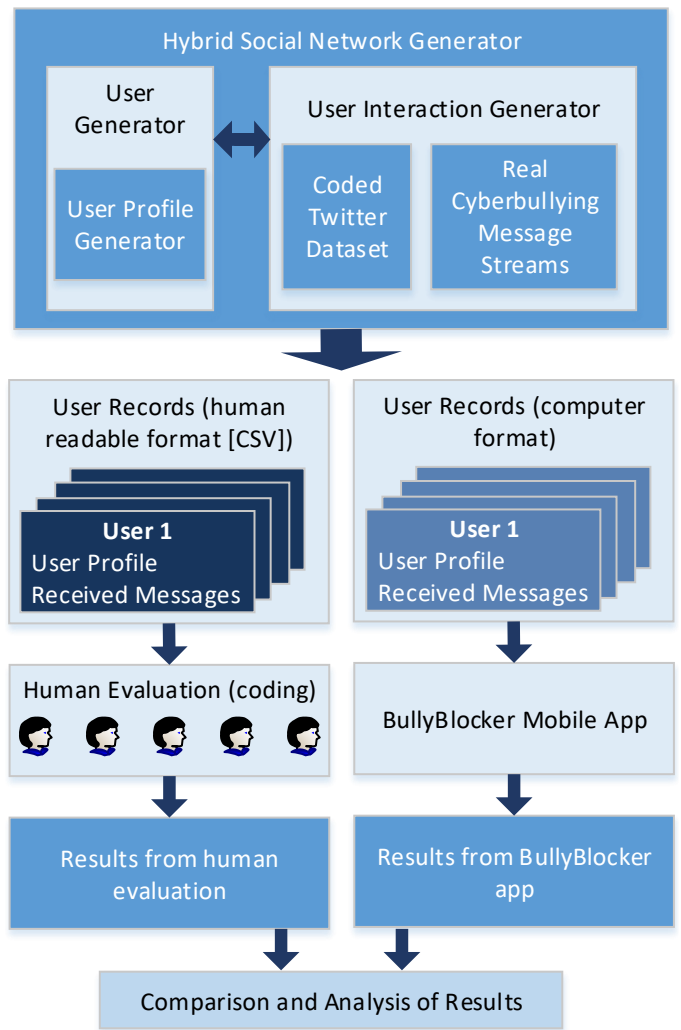


Fig. 10. Evaluation Framework using a Hybrid Social Network

The *User Generator* module outputs a set of N users and their profile information. The profile information consists of the attributes associated with the vulnerability factors presented in Fig. 6 (e.g., age, gender, race, frequency of daily internet use, bullying history, etc.). With the exception of the bullying history attribute, specific values for the different vulnerability factors were evenly distributed among all possible values or ranges. For the previous bullying history attribute, the distribution of values was: no previous bullying (50%), experienced bullying last month (16.66%), from one to two months (16.66%), and more than two months ago (16.66%).

The *User Interaction Generator* module produces a set of interactions (message streams) among the created users. The goal of this module is to create message sequences that are similar to the ones found in real-world social networks. To this end, this component uses two sources of real-world messages: (1) a coded (labeled) Twitter dataset (composed of subsets of cyberbullying and non-cyberbullying (i.e., normal) messages), and (2) real cyberbullying message streams. The coded Twitter dataset was obtained following the procedure suggested in [46].

We crawled this dataset using the Twitter streaming API [47] from September 19th to 25th, 2017 with the following keywords: *nerd, gay, loser, freak, emo, whale, pig, fat, wannabe, poser, whore, should, die, slept, caught, suck, slut, live, afraid, fight, pussy, cunt, kill, dick, bitch*. We initially obtained 4,730,766 tweets. After the initial data collection, we employed the pre-trained *Bully* classifier [48] to label each tweet in the crawled dataset and extracted a refined subset with high confidence, containing 7,500 positive samples (i.e., cyberbullying) and 7,500 negative samples (i.e., normal). Then, these 15,000 tweets were further labeled by two well-trained human annotators with backgrounds in psychology and computer science. A third trained annotator was asked to resolve any discrepancies between the ratings of the initial two annotators. After resolution of discrepancies and data cleaning, we obtained the final dataset composed of 3,647 cyberbullying tweets (referred to as TwitterCB) and 11,347 normal ones (referred to as TwitterNonCB). The second dataset (NewsCB) was composed of real cyberbullying message streams found in real social networks, e.g., Facebook, Twitter and Instagram. Many of these streams were obtained from news articles reporting well-known instances of cyberbullying. Due to required manual work to identify and extract these cyberbullying sequences, their number is relatively small (100 streams, where each stream contains between one and eleven messages). An important benefit of this dataset, however, is that it captures information about the way cyberbullying messages are distributed over time. These temporal properties are maintained in the message streams generated using this data source.

The User Interaction Generator module uses both message sources to generate the interactions among the created users. The key parameters used in this step are the total number of users ($N=400$), the number of cyberbullying streams in NewsCB ($K=100$), the number of days ($D=90$), and the maximum number of messages per user ($M=100$). The messages were generated as follows: For each of the first K users, the stream of the i th user contains all of the cyberbullying messages included in the i th stream in NewsCB. The remaining messages (totaling M messages per stream) are generated adding $BF \cdot (M - \text{NewsCB}[i].\text{length})$ cyberbullying messages and $(1 - BF) \cdot (M - \text{NewsCB}[i].\text{length})$ normal messages. BF (Bullying Fraction) is the fraction of the remaining messages that are cyberbullying interactions. This value, in the range $[0.0, 1.0]$, is randomly computed for each of the first K users. For each of the remaining $N - K$ users, the message stream of a given user is generated by interleaving $BFF \cdot D$ bullying messages from TwitterCB and $(M - BFF \cdot D)$ normal messages from TwitterNonCB. BFF (Bullying Frequency Factor) is the number of cyberbullying messages that a user receives per day and is also randomly generated for each of the remaining users in the range $[0.0, 1.0]$.

We expect that the hybrid network generator and the real-world datasets will be used by other researchers to evaluate the performance of future comprehensive and holistic identification models. To facilitate these tasks, we have made available the source code of the generator and its input datasets [45].

Generated Datasets. A final step of the Hybrid Social Network Generator is to produce two datasets capturing the information of the social network. In both datasets, each record represents

The image shows a screenshot of a user profile and a list of tweets. The profile information includes: Name: User0, Age: 17, Gender: Female, Ethnicity: Other, Race: African_American, Depression: No, Anxiety: No, Self Esteem Issues: No, History Bullied: The user was bullied last month, Disciplinary Issues: Yes, Substance Abuse: No, and Internet Usage: Weekly use between 4 and 6 hours. Below the profile is a table of tweets with columns for Timestamp, SenderID, Bullying, and Text.

Timestamp	SenderID	Bullying	Text
11/19/2017 5:19	102	B	No he sounds like a dick https://t.co/Dptnt0E3HD
11/20/2017 2:09	71	B	If you don't cry yourself to sleep every night to miserable at best, are you even emo?
11/19/2017 11:15	102	B	Calling is bad. If it wasn't for that pass interference wouldn't of scored. He suck bro
11/20/2017 22:06	139	B	Can't sling about someone's dick but you can call Marcos Alonso a murderer at the m
11/20/2017 5:16	102	B	You should eat!
11/22/2017 4:58	44	B	"If you see a good fight, get in it."
11/19/2017 17:31	16	B	TF???? EAT SOMETHING!!!!!!!!!!!!!!!
11/25/2017 3:26	107	B	Those Are Bones. Need to eat
11/22/2017 5:09	39	B	I hope your 99 kids are fine
12/12/2017 20:52	140	B	I cant beleave that fat poo ball owns the goodlooking girl contest and they say he w
11/21/2017 23:51	21	B	Word this nigga is dick eatin hard mf tags me in a post EVERYTIME the giants lose li
11/25/2017 9:05	319	B	See.... because I will fight you right now!! _Ã;Ã-Ã
11/23/2017 22:46	91	B	Too skinny
11/21/2017 0:45	39	B	Her skin is unique and beautiful, but she seems too thin.
11/28/2017 16:27	153	B	Don't hesitate to give your opinions either, you gay. I'll read them in the mornin
12/4/2017 11:49	114	B	#RETWEET IF YOU LIVE IN ENGLAND AND FOLLOW ME FOR #FOLLOWBACK !
12/8/2017 3:45	67	B	Whoever fucked me over on 2015, I hate you bitch https://t.co/ai2aOOCx4P
12/9/2017 22:50	367	B	Hour and a half, that's close-ish! I'd love to live that close to a sea or an ocean
12/12/2017 3:56	16	B	love to crack open a bundaberg with die jungen
12/14/2017 10:20	87	B	waw...not only slutty ...v...i feel that you are slut ä öï öï öï
12/2018 2:02	124	B	bitch get over that heartbreak and keep it tf movin
11/25/2017 4:07	246	B	Yes Oh.The boy dey kill us with off key here
11/22/2017 18:46	67	B	Too thin
12/13/2017 17:46	32	B	youãóre now the reason that I fight
12/11/2017 20:56	389	B	Only 30 minutes until I'm live on twitch with more Earthbound! https://t.co/uz4jmMSI
12/9/2017 15:34	202	B	But his point is lost? Troops win battles, and if they don't fight, they don't win. They I
12/2/2017 20:44	124	B	Love you Joel but I am weary fighting Thyroid #Cancer & feeling there is nobody
11/24/2017 1:42	43	B	Eat few burritos
1/12/2018 11:33	22	B	Get that son of a bitch out of the White House.
12/25/2017 15:34	151	B	because I wil fucking die if i see one i swear i'll kill everyone in the cinema
12/16/2017 20:54	107	B	That's absolute bullshit. Our troops fight for the right to protest. They defend the co
12/5/2017 21:41	96	B	KILL IT GIRL YES SLAY
12/23/2017 13:17	183	B	Why does Dysphoria make me afraid to talk about girls like I am one, why does my b

Fig. 11. Sample User Record

the social network data associated to a single user. Each user record is composed of (1) the user profile information (user ID and vulnerability related features), and (2) the set of messages received by this user. Each message contains the ID of the user sending the message, the timestamp, and the message content.

The two datasets contain the same data but use different representation formats. Each record in the first dataset uses a human-readable format. Each record in the second dataset is structured as a document intended to be processed programmatically using the BullyBlocker app. The only content difference between the two datasets is that the human-readable one includes a flag that identifies the cyberbullying-related messages. A sample of a produced user record is presented in Fig. 11.

Human Evaluation. The generated human-readable dataset containing the information of the hybrid social network was evaluated by members of our research team. Each record was assessed independently by two designated annotators and any discrepancies between the annotators' assessments were resolved by a third rater. Evaluating a user record entailed assessing the user's entire profile information and stream of messages, and assigning the record a Bullying Rank between 0 and 100 to reflect the probability that the user is experiencing or has recently experienced cyberbullying. A discrepancy in annotators' assessments was defined as two people assigning Bullying Rank values in different risk level categories (low risk [0,33], moderate risk [34,66], and severe risk [67, 100]).

Evaluation using the BullyBlocker App. The generated user records were also processed by the BullyBlocker app. To this

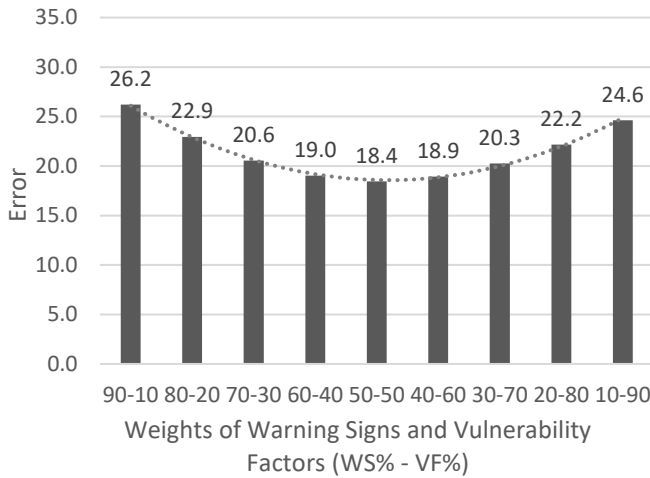


Fig. 12. Average Error for Different Weights of Warning Signs and Vulnerability Factors

end, the app was extended by a module that read from the generated dataset instead of obtaining the information from Facebook. This module also executed the Bullying Rank Computation task for each user and saved the Bullying Rank values generated by the app.

Comparison and Analysis of Results. In the last step of the process, we compared the results obtained from the human evaluation phase against the results obtained using the automated BullyBlocker cyberbullying identification model. The results of this comparison are presented in the next subsection.

B. Evaluation Results

Fig. 12 shows the average error of the Bullying Rank values computed by the proposed BullyBlocker model, with the Bullying Rank (probability that an adolescent is being cyberbullied) expressed as a percentage value (1-100). We compute the error as the absolute value of the difference between the Bullying Rank produced by the app and the average value of the human coding results. This figure presents the average error for various weight configurations of the two main components of the Bullying Rank (Warning Signs and Vulnerability Factors). As shown in Fig. 12, the BullyBlocker app produces the smallest error (18.4 percentage points) when there is an even distribution of weights between Warning Signs and Vulnerability Factors (50%-50%). The error gradually increases as either of the components is weighted more heavily than the other. Based on these results, we set the weights of both components to 50% in the latest version of the app.

Fig. 13 presents the frequency of errors for various error ranges. The results show that when the weights of Warning Signs and Vulnerability Factors in the BullyBlocker app were set to 50%, the most frequent error values fell within the lowest ranges. That is, in 33% of the cases, the error (reflecting the difference between the Bullying Rank estimated by the human annotators and the Bullying Rank calculated by the app) was smaller than 10 percentage points, while in 60% of the cases, the error was smaller than 20 percentage points.

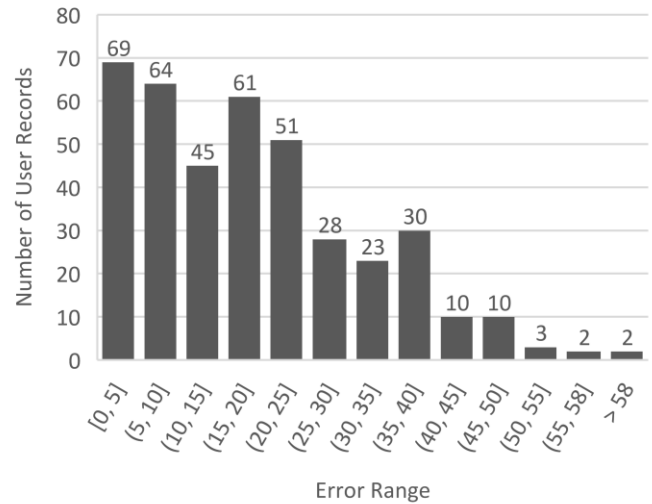


Fig. 13. Histogram of Errors

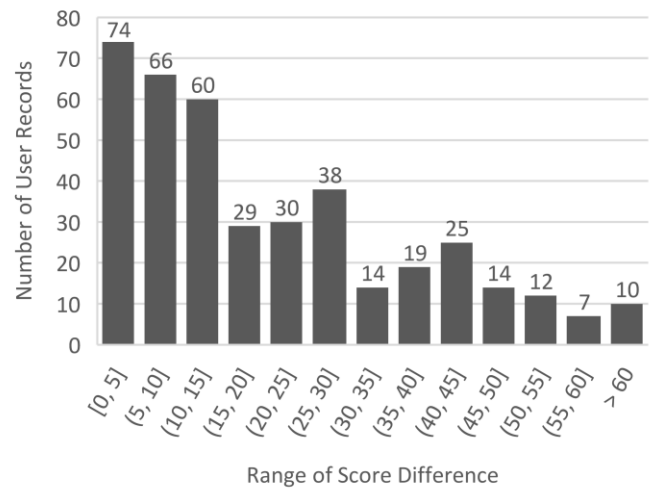


Fig. 14. Histogram of the Score Difference between Human Coders

The results in figures 12 and 13 should be interpreted considering an inherent challenge that human annotators face when estimating the probability that an adolescent is being cyberbullied, which stems from subjectivity in how annotators approach the task of quantifying cyberbullying risk. That is, we had eight human annotators each rate a subset of 100 cases from the human-readable version of the dataset; with each of the 400 cases generated by the hybrid social network rated by two annotators independently. Whereas annotators maintained a consistent strategy for estimating cyberbullying risk for the 100 cases to which they were assigned, individual differences in how each annotator interpreted the data holistically and translated their assessments into the numerical Bullying Rank index likely emerged. Fig. 14 shows the distribution of the score difference in the Bullying Rank estimates made by the two annotators who evaluated each case during the Human Evaluation phase. The distribution in this figure is similar to the distribution of the BullyBlocker model error presented in Fig. 13. Specifically, in 35% of the cases, the score difference between the two human

annotators assigned to a particular case was smaller than 10 percentage points, while in 57% of the cases, the score difference was smaller than 20 percentage points.

The results presented in this section show that the proposed model produces relatively small error in most cases. We expect that the error levels of the model will be further reduced by integrating some of the techniques described in the future work section.

V. FUTURE WORK

As an emerging sphere of research, efforts to develop and evaluate the accuracy of models for the automated identification of cyberbullying can benefit from future work in several key areas. In this section, we describe some of these areas and provide details of our team’s progress along these research paths.

A. *Alternative Identificaton Models*

This area involves the use of other computational techniques like similarity-aware data processing and machine learning to build alternative holistic cyberbullying identification models that consider both an array of profile information features and the users’ streams of messages. Two important tasks in this area are the comparison of multiple models and the study of integration mechanisms to build highly accurate hybrid models.

Similarity-aware Model for Cyberbullying Identification. To this end, we are investigating cyberbullying identification models that use the power of similarity operators [8, 9, 49, 50], i.e., data processing operators like the Similarity Join and Similarity Grouping that identify and exploit similarities in the data. An initial idea to build this model is to use a vector-based representation of a person’s behavior. Our efforts in this area are directed towards building wide feature vectors using the information on cyberbullying risk factors identified in our current BullyBlocker model. Specifically, the vector could include numeric measures for various warning signs (e.g., number of insulting messages) and vulnerability factors (e.g., a recent move to a new neighborhood or school). For example, considering an initial set of factors, the structure of the vector would be as follows: #EmbarrassingPictures, #TotalFeed Messages, #TotalPictureComments, #Pictures, #Bullies, #Friends, ..., Age, Female?, #DaysSinceNewNeighborhood, #DaysSinceNewSchool, Hispanic?, AfricanAmerican?, ...). Moreover, we are exploring the inclusion of factors that are associated with the role and position of the potential victim in his or her social network (by analyzing features like closeness centrality, betweenness centrality, degree centrality, eigenvector centrality, and clustering coefficient [32]), as well as features aimed at capturing language patterns based on the message streams. This model will also represent common cyberbullying behavior patterns as cyberbullying behavior vectors. The outcome of this sub-task will be a number of cyberbullying feature vectors that, as a group, represent the most common patterns of cyberbullying victimization. The distance of an adolescent’s feature vector to the cyberbullying behavior vectors can be used to estimate the likelihood that a person is being

cyberbullied. This approach also enables other interesting types of analyses. For instance, using clustering or similarity grouping, we can identify and study groups of adolescents who are experiencing similar types of social, emotional, or behavioral issues.

Machine Learning Model for Cyberbullying Identification.

We are also working to design and study comprehensive machine learning models [31] for cyberbullying identification. The cyberbullying identification problem can be modeled as a classification (discrete output classes, e.g., low, moderate and high cyberbullying risk) or regression (continuous output, e.g., Bullying Rank value) task and multiple strategies can be used to implement them, e.g., Logistic Regression, Sparse Neural Networks, Support Vector Machines, and Naïve Bayes. Specifically, we are in the early stages of designing a Sparse Neural Networks model. This approach will enable building a global artificial neural network by connecting smaller complete neural networks that can focus on specific classification sub-tasks, e.g., considering subsets of the warning signs and vulnerability factors. This approach will allow us to build the global model incrementally by designing, implementing, and training individual neural networks that consider well-defined subsets of the cyberbullying factors (e.g., factors pertaining to race and gender). Some of these small neural networks may be based, in fact, on previously proposed models [10-15].

B. *Integrating New Vulnerability Factors*

Another important area of future work is the integration of new vulnerability factors into cyberbullying identification models like the one presented in this paper. In this area, we plan to continue drawing on emerging research findings in psychology to guide the identification of additional factors. Among the factors we plan to integrate are physical stature/weight, disability status, and concurrent use of multiple social networking sites. Other factors that, to our knowledge, have yet to receive empirical attention pertain to an adolescent’s minority status *within* their specific school environment, neighborhood, or community. Similarly, socioeconomic status, religious identity, and immigrant status—and, importantly, the extent to which these aspects of a teen’s identity contribute to their minority or fringe status within their immediate social environment—may also provide valuable insights for the identification of cyberbullying risk. Moreover, variables like degree of parental oversight of social media use and limitations on an adolescent’s access to technology or social media can be modeled as protective factors associated with a decreased likelihood of being cyberbullied. These factors have been explored in a small handful of previous studies [16], although additional research is needed to better understand the extent to which they might buffer cyberbullying risk. Other information collected through the app, such as changes in one’s relationship status, deletions from one’s friend list, and hiding certain posts from view on one’s newsfeed or timeline in Facebook, may provide an unprecedented mechanism for tracking meaningful changes in one’s peer circle. Interestingly, most of the psychology research on cyberbullying has relied on adolescents’ self-reports [16, 18] of victimization, which may be influenced by their reluctance to report instances of cyberbullying as well as other methodological limitations [16]. Data collected through

the BullyBlocker app can thus circumvent several issues stemming from self-report measures of cyberbullying.

C. Expanding the Synergy with the Psychology Community

Applications like BullyBlocker can also help inform the work of psychologists from both a research and clinical practice perspective. For example, parent feedback regarding the perceived accuracy of the app's underlying identification model can be compared and combined with clinical experts' assessments of cyberbullying, yielding a promising avenue for future psychological research. Future studies could, for instance, examine the degree of overlap in clinicians' and parents' assessments of cyberbullying, and compare each with adolescents' self-reports. Feedback from clinicians and parents will also be beneficial for understanding the extent to which various ranges of Bullying Rank values map onto the presence and severity of clinical symptoms that are directly observed by parents and clinical experts. Furthermore, parent feedback could also provide a platform for investigating parents' more general attitudes about the use of automated tools for identifying a broad range of behavioral issues. Automated tools also have the potential to aid in the identification of symptoms of depression and anxiety, undue amounts of stress, low self-esteem, relationship violence, indicators of self-harm, and suicidal thoughts. Feedback from parents and adolescents can provide essential usability information (e.g., what level of detail about identified cyberbullying instances parents feel most comfortable receiving through the app, and what level of detail might deter adolescents from providing their parents with their social networking site login information).

VI. CONCLUSIONS

Cyberbullying is the most common online risk for adolescents. While the prevalence and determinants of cyberbullying have received considerable attention among researchers in the psychology community, there has been relatively little work on the automatic identification of cyberbullying in social networking sites, and even less work that seeks to bridge the efforts from computer science and psychology. This paper thus proposes a computational model for cyberbullying identification that builds on the research findings within the psychology literature. The paper also describes the design of BullyBlocker, an app that implements the proposed model, discusses the model's effectiveness in the context of a newly-developed evaluative framework, and presents several ways in which the model can be extended. Our hope is that BullyBlocker, which has been recently made available through the Apple App Store, will have a strong societal impact, by identifying youth most vulnerable to cyberbullying victimization and by enabling parents to help their children in time to make a difference.

ACKNOWLEDGMENT

The authors would like to thank ASU students Lisa Tsosie, Jaime Chon, Tara Tucker, Chance Brown, Liz Garcia, Bryan Sawkins, Rusty Conway, Anthony Nieuwenhuys, Tom Schenk, Lu Cheng, Ashley Trow, Ayush Sanyal, Linle Jiang, Victoria Delgadillo, and Carmen Sanchez for their contributions to the design and implementation of the BullyBlocker app.

REFERENCES

- [1] <http://www.bullyingstatistics.org/content/cyber-bullying-statistics.html>.
- [2] <http://www.pewinternet.org/2015/04/09/teens-social-media-technology-2015/>.
- [3] Piazza, J. and Bering, J. M. *Evolutionary cyber-psychology: Applying an evolutionary framework to Internet behavior*. Computers in Human Behavior, 25 (6): 1258-1269, 2009.
- [4] Ortega, R., Elipe, P., Mora-Merchin, J. A., Calmaestra, J., and Vega, E. *The emotional impact on victims of traditional bullying and cyberbullying: A study of Spanish adolescents*. Journal of Psychology, 217 (4): 197-204, 2009.
- [5] Waasdorp, T. E. and Bradshaw, C. P. *Examining student responses to frequent bullying: A latent class approach*. Journal of Educational Psychology, 103 (2): 336-352, 2011.
- [6] Dooley, J. J., Pyzalski, J., and Cross, D. *Cyberbullying versus face-to-face bullying: A theoretical and conceptual review*. Journal of Psychology, 217 (4), 182-188.
- [7] Silva, Y. N., Rich, C., and Hall, D. *BullyBlocker: Towards the Identification of Cyberbullying in Social Networking Sites*. The IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2016, 1377-1379.
- [8] Silva, Y. N., Pearson, S., and Cheney, J. A. *Database Similarity Join for Metric Spaces*. The International Conference on Similarity Search and Applications (SISAP), 2013. Springer LNCS, 8199: 266-279, 2013.
- [9] Tang, M., Tahboub, R. Y., Aref, W. G., Atallah, M. J., Malluhi, Q. M., Ouzzani, M., Silva, Y. N. *Similarity Group-by Operators for Multidimensional Relational Data*. IEEE Transactions on Knowledge and Data Engineering (TKDE), 28, 2, pp 510-523, 2016.
- [10] Dinakar, K.; Reichart, R.; Lieberman, H. *Modeling the Detection of Textual Cyberbullying*. The International AAAI Conference on Web and Social Media (ICWSM), 2011.
- [11] Rafiq, R. I., Hosseinmardi, H., Han, R., Lv, Q., Mishra, S., and Arredondo-Mattson, S. *Careful what you share in six seconds: Detecting cyberbullying instances in Vine*. The IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2015.
- [12] Hosseinmardi, H., Rafiq, R. I., Han, R., Lv, Q., and Mishra, S. *Prediction of Cyberbullying Incidents in a Media-based Social Network*. The IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2016.
- [13] Reynolds, K., Kontostathis, A., and Edwards, L. *Using Machine Learning to Detect Cyberbullying*. The 10th International Conference on Machine Learning and Applications and Workshops (ICMLA), 2011.
- [14] Huang, Q., Singh, V. K., and Atrey, P. K. *Cyber Bullying Detection Using Social and Textual Analysis*. The 3rd International Workshop on Socially-Aware Multimedia (SAM), 2014.
- [15] Squicciarini, A., Rajtmajer, S., Liu, Y., and Griffin, C. *Identification and characterization of cyberbullying dynamics in an online social network*. The IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2015.
- [16] Kowalski, R. K., Giumetti, G. W., Schroeder, A. N., and Lattanner, M. R. *Bullying in the Digital Age: A critical review and meta-analysis of cyberbullying research among youth*. Psychological Bulletin, 140(4), 1073-1137, 2014.
- [17] Cook, C. R., Williams, K. R., Guerra, N. G., Kim, T. E., and Sadek, S. *Predictors of bullying and victimization in childhood and adolescence: A meta-analytic investigation*. School Psychology Quarterly, 25(2), 65-83, 2010.
- [18] Tokunaga, R. S. *Following you home from school: A critical review and synthesis of research on cyberbullying victimization*. Computers in Human Behavior, 26(3), 277-287, 2010.
- [19] Williams, K. R., and Guerra, N. G. *Prevalence and predictors of internet bullying*. Journal of Adolescent Health, 41(6), S14-S21, 2007.
- [20] Kowalski, R. M., and Limber, S. P. *Electronic bullying among middle school students*. Journal of Adolescent Health, 41(6), S22-S30, 2007.

- [21] Ybarra, M. L., and Mitchell, K. J. How risky are social networking sites? A comparison of places online where youth sexual solicitation and harassment occurs. *Pediatrics*, 121(2), e350-e357, 2008.
- [22] Hinduja, S., and Patchin, J. W. Cyberbullying: An exploratory analysis of factors related to offending and victimization. *Deviant Behavior*, 29(2), 129-156, 2008.
- [23] Mishna, F., Khoury-Kassabri, M., Gadalla, T., and Daciuk, J. *Risk factors for involvement in cyber bullying: Victims, bullies and bully-victims*. *Children and Youth Services Review*, 34(1), 63-70, 2012.
- [24] Ang, R. P., Chong, W. H., Chye, S., and Huan, V. S. Loneliness and generalized problematic Internet use: Parents' perceived knowledge of adolescents' online activities as a moderator. *Computers in Human Behavior*, 28 (4), 1342-1347, 2012.
- [25] Law, D. M., Shapka, J. D., and Olson, B. F. *To control or not to control? Parenting behaviours and adolescent online aggression*. *Computers in Human Behavior*, 26 (6), 1651-1656, 2010.
- [26] Dooley, J. J., Pyżalski, J., and Cross, D. *Cyberbullying versus face-to-face bullying: A theoretical and conceptual review*. *Journal of Psychology*, 217 (4), 182-188, 2009.
- [27] Wolke, D., Lereya, T., and Tippet, N. *Individual and Social Determinants of Bullying and Cyberbullying*. *Cyberbullying*. Ed. T. Vollink, Ed. F. Dehue, Ed. C. Guckin. Routledge, 26-53, 2016.
- [28] Patchin, J. W. and Hinduja, S. *Cyberbullying - An Update and Synthesis of the Research*. *Cyberbullying Prevention and Response*. Ed. J. W. Patchin, Ed. S. Hindija. Routledge, 13-35, 2012.
- [29] Due, P., Merlo, J., Harel-Fisch, Y., Damsgaard, M. T., Holstein, B. E., Hetland, J., Currie, C., Gabhainn, S. N., Gaspar de Matos, M., and Lynch, J. *Socioeconomic Inequality in Exposure to Bullying During Adolescence: A Comparative, Cross-Sectional, Multilevel Study in 35 Countries*. *American Journal of Public Health*, 99(5), 907-914, 2009.
- [30] Blake, J. J., Lund, E. M., Zhou, Q., Kwok, O. M., and Benz, M. R. *National prevalence rates of bully victimization among students with disabilities in the United States*. *School Psychology Quarterly*, 27(4), 210, 2012.
- [31] Kelleher, J. D., Namee, B. M., D'Arcy, A. *Fundamentals of Machine Learning for Predictive Data Analytics: Algorithms, Worked Examples, and Case Studies*. The MIT Press, 1st edition, 2015.
- [32] Golbeck, J. *Analyzing the social web*. Morgan Kaufmann, 2013.
- [33] BullyBlocker App in the Apple Store. App <https://itunes.apple.com/us/app/bullyblocker-app/id1236410370?mt=8>
- [34] Baldry, A. C., Farrington, D. P., Sorrentino, A. *Cyberbullying in youth: A pattern of disruptive behaviour*. *Psicología Educativa*, 22(1), 19-26, 2016.
- [35] Guo, S. 2016. *A meta-analysis of the predictors of cyberbullying perpetration and victimization*. *Psychology in the Schools*, 53(4), 432-453, 2016.
- [36] Fedewa, A. L., and Ahn, S. *The effects of bullying and peer victimization on sexual-minority and heterosexual youths: A quantitative meta-analysis of the literature*. *Journal of GLBT Family Studies*, 7, 398-418, 2011.
- [37] Drug Rehab Website. <http://www.drugrehab.com/guides/bullying>
- [38] Help your Teen Now Website. <http://www.helpyourteennow.com/cyberbullying-and-addiction-in-teenagers>
- [39] Drug Abuse Website. <http://www.teens.drugabuse.gov/blog/post/four-things-know-about-cyberbullying>
- [40] Beyond Bullying Website. <http://www.beyondbullying.com/racistbullying.html>
- [41] Hinduja, S. and Patchin, J. W. *Social Influences on Cyberbullying Behaviors Among Middle and High School Students*. *Journal of Youth and Adolescence*, 42(5), 711-722, 2013.
- [42] Bonnet, D. G. Transforming odds ratios into correlations for meta-analytic research. *American Psychologist*, 62(3), 254-255, 2007.
- [43] Squicciarini, A., Rajtmajer, S., Liu, Y., and Griffin, C. *Identification and characterization of cyberbullying dynamics in an online social network*. The IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2015.
- [44] Van Geel, M., Vedder, P., Tanilon, J. Relationship between peer victimization, cyberbullying, and suicide in children and adolescents: A meta-analysis. *JAMA Pediatrics*, 168(5), 435-442, 2014.
- [45] BullyBlocker Project Website. <https://bullyblocker.project.asu.edu/data>
- [46] Nand, P., Perera, R., and Kasture, A. How Bullying is this Message?: A Psychometric Thermometer for Bullying. In COLING. 695-706, 2016.
- [47] Morstatter, F., Pfeffer, J., Liu H., and Carley, K. M. Is the Sample Good Enough? Comparing Data from Twitter's Streaming API with Twitter's Firehose. In Proceedings of the 7th International AAAI Conference on Web and Social Media (ICWSM), 2013.
- [48] Xu, J.-M., Jun, K.-S., Zhu, X., and Bellmore, A. Learning from bullying traces in social media. In Proceedings of the 2012 conference of the North American chapter of the association for computational linguistics: Human language technologies. Association for Computational Linguistics, 656-666, 2012.
- [49] Yu, C., Cui, B., Wang, S., Su, J. Efficient index-based knn join processing for high-dimensional data. *Information and Software Technology*, 49 332-344, 2007.
- [50] Silva, Y. N., Pearson, S. S., Chon, J., Roberts, R. Similarity Joins: Their implementation and interactions with other database operators. *Information Systems* 52:149-162, 2015.