# Using Machine Learning, Image Processing & Neural Networks to Sense Bullying in K-12 Schools

Lalit Kumar[1], Palash Goyal[2], Karan Malik[3].Rishav Kumar[4]

[1]Gazelle Information Technologies, Delhi, India [2]Mount Carmel School, Delhi, India [3]Departement of Information Technology, Indraprastha University, Delhi, India [4]Mount Carmel School, Delhi, India

[1]lalit.k1997@gmail.com [2]palashgoyal1608@mountcarmeldelhi.com [3]karanmalik2000@gmail.com

*Abstract—We all have heard about bullying and we know that it is an immense challenge that schools have to tackle. Many lives have been ruined due to bullying and the fear it implants into students' mind has caused many of them to go into depression which can lead to suicide. Traditional methods [1] need to be accompanied with modern technology to make the method more effective and efficient. If real time alerts are to school staff, they can identify the perpetuator and extricate the victim swiftly. It this proposed method an AI based solution is implemented to monitor students using standard school surveillance technologies and CCTV to maintain a decorum and safe environment in the school premise. Also the proposed method utilizes other unstructured sources such as attendance records, social media activity and general nature of the students to deliver quick response. Artificial Intelligence (AI) techniques like Convolutional Neural Networks (CNN), which includes image processing capabilities, logistic regression methods, LSTM (Long short-term memory), and pre-trained model Darknet-19 is used for classification. Further, the model also included sentiment analysis to identify commonly used abuse terms and noisy labels to improve overall model accuracy. The model has been trained and validated with the realistic data from all the sources mentioned and has achieved the classification accuracy of 87% for detecting any sign of bullying.*

*Index Terms-- AI, CNN, Class Entropy Loss, Data-pre-processing, Data pipeline, Facial Recognition, NLP, NN, Sigmoid, Sentiment analysis, LSTM, Darknet-19.*

## 1. Introduction

What is bullying? For starters, bullying is aggressive behaviour among teenagers. This comes
\

from many factors which involve ego, power imbalance, upbringing etc. This behaviour has no bounds and can impact a victim's physical and mental health conditions. Bullying can take many forms from verbal to physical. With the introduction of social media, social bullying, a new type of bullying which involves harming one's reputation or relationships, has sprung up. A study was conducted by Symantec Reports with the help of parents of many victims. They noted that almost 24% of the students were involved in some shape or form of bullying [2]. Despite such a high percentage of victims, most methods to keep bullying in check have not borne any fruit and mitigation against bullying remains an enigma.

The traditional methods with their limited use of technology are proved not be very effective. It is time to integrate modern technology to develop more intelligent solutions. The proposed method in this paper will use specific elements of AI tools which exploit the already available school infrastructure to make it a new means of keeping taps on bullying.

*Methodology of research and solution development:*

First of all, we have identified all the physical, mental, emotional, social and cyber bullying types and parameters that are prevalent. For this extensive study of existing literature has been done [1] [2] [3]. Subsequently these parameters were studied to identify how they can be objectively analysed. This involved identifying the data sources that can give input signals for that bullying parameter. Once the data sources were identified actual data was recovered from these sources. Again, existing literature was studied [4], [5], [6] extensively to identify AI/ML algorithms that can be used to classify them with binary outputs linked to whether this is a bullying situation or not. The authors then

put together all these elements to develop a working model of their solution with a proof of concept done in a school. The research methodology followed in this is as given below in figure 1.
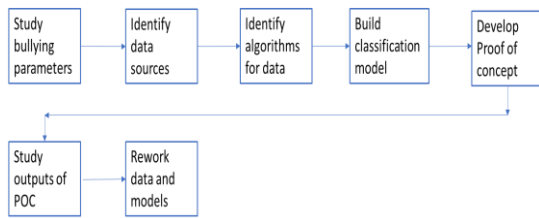


Fig. 1: Research methodology

### Bullying Parameters/Features

To identify the signs of bullying, we have taken multiple parameters under consideration through which we can filter-out real bullies from other behaviours which could be mistaken for bullying. Features like rude behaviour, troubling authorities regularly, low attendance, use of foul language, low scores, drug use, explicit content etc. are used to filter the data. In the proposed model it is planned to integrate AI tools with school infrastructure and data like cameras with microphones, student portals containing community forums where students and teachers interacts, attendance of students, their score-cards etc. Data used for analysis and a complete infrastructural setup is discussed  below in tables 1 and 2.

*Table 1: How to identify a victim of bullying*

| Victim Parameter | Data Source |
|---|---|
| Show deteriorating grades | Academic report |
| Regular absenteeism | Attendance records |
| Being picked on, pushed, punched etc. (physically harassed) | CCTV images |
| Downfall of social skills | Registration for school events |
| | and extra-curricular groups |
| Suffering from learning or mental disabilities | Counsellor reports |
| Unwilling to go to school regularly | Attendance records |
| Random bruises, missing belongings or torn clothing | CCTV images |
| Prone to attacks of anxiety | School Medical Reports |
| Alone at lunch-breaks. | CCTV images |
| Experiences regular nightmares | Medical room report |
| Starts bullying younger or weaker kids to vent out the frustration | CCTV images |

*Table 2: Bullying perpetrator parameters and data sources*

| Perpetrator parameters | Data source |
|---|---|
| Rude behaviour with students as well as teachers | CCTV images, audio, school report |
| Low Grades | School report |
| Lacks empathy or guilt | School report |
| Feeling of entitlement because of being good in school, sports or belonging to a prominent family | School report |
| Short tempered and having emotional outbursts | CCTV images, audio, school report |
| Usually popular or among a big group | CCTV images, audio, school report |
| Regularly get into trouble with authority | School report |

### Integration with School Infrastructure:

Recordings of playground, common-area (like corridor, locker-rooms etc.) and classes will be taken along with their audios via cameras that are installed. More detailed information associated to students will be extracted from student portals. This will further help in understanding and analysing student behaviour and personality. This raw information is then reworked into structured data, which will supplement the learning algorithm in predictions and analysis. Based on this, appropriate action can be taken by school authorities. This method can also be inverted and be used to discern the victims of harassment.

Data Analysis:
Photos uploaded by the students on community forums along with acquired video footage from CCTV, split into images, and ran through the algorithm to identify: drugs, number of faces, anxiety attacks, crying, isolation, fighting, torn clothes, bruises, sleeping, smoking, hard drinks, gore, explicit and adult content. The image classifier uses CNN and pretrained network call Darknet-19 (trained on ImageNet dataset), along with LSTM which provide the capability to process sequence of image data with feedback connections.

Audios which are mapped to text using Google Cloud Speech API alongside comments at student community are further used to determine the following features: tone, amplitude and pitch of voice, language used (explicit or not), uppercase text, text length and sentiments, threatening statements, trolling, unpleasant comments and distasteful words. Other attributes like low-attendance, incompetent grades, enrolment in extracurricular activities, teacher-student interaction, frequency of councillor appointments and behaviour report by staff will also be considered. Data for these inputs is converted from physical form into digital form first.

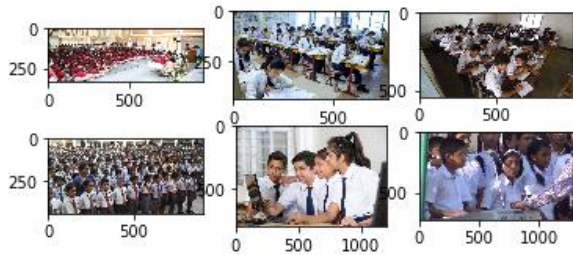Fig. 2.  Sample images (in case of bullying) from dataset



Fig. 3.  Sample Images (non- bullying) from dataset

### Techniques Involved

Convolution Neural Network is a type of neural network (NN) that provides the capability to convert pixels into well structured data [5]. CNNs replicate the function of the frontal lobe of the human brain (cerebral cortex), which is responsible for processing audio visual stimulus in humans. To process images, CCTV footage is spliced into still images and then each frame is analysed to extract the crucial and important features which can be further analysed for more refined results.

Proposed network architecture for video classification is shown in Fig 5, It has been already shown that adding LSTM (which extract global temporal features) after CNN, the local temporal features obtained from optical flow are also of great importance [7].

Optical flow is due to displacement of boundaries or movement of object across two consecutive frames. Optical Flow is best used for action-recognition, so it's functionality is mimicked by taking two consecutive frames as input to training model. These consecutive frames are fetched to pre-trained CNN model (Darknet-19). Then both of the frames-values from the bottom layer of pre-trained model are fed as an input to additional CNN. Additional CNN network is now to supposed to learn from the local motion features including the invariant features by comparing the both frame. The top layer of the pre-trained network is also fetched to another additional CNN to learn from the comparison of high-level features of both the frames. Furthermore,

the output of both the additional CNN is fully-connected with LSTM cell, which further classify the images as bully & non-bully.
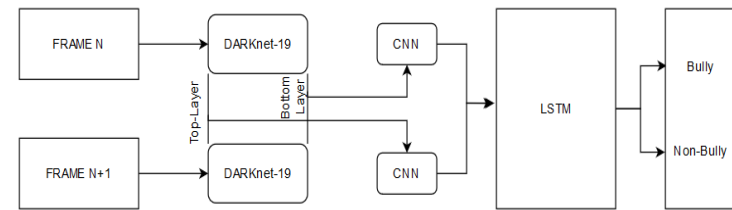


Fig 5

For audio analysis, we have mapped CCTV voice-output with Google Cloud Speech API which uses CNN and provides real-time streaming of speech recognition and conversion from audio to text. This lies within our CCTV and Google Cloud Storage for sentiment analysis to be done on it.

For text analysis, mentioned words are mapped with sentiment analysis dictionary ('Liu and Hu opinion lexicon' containing 6800 positive and negative words [4]). To extract feature from the text, multinomial Naïve Bayes is used. It helps in classification of bad words and rude comments from audio-to-text converted files as well as from student's community portal including explicit and vulgar remarks, text length as well as usage of upper-case letters in community forums etc. All such attributes are then listed for feature extraction by the model.
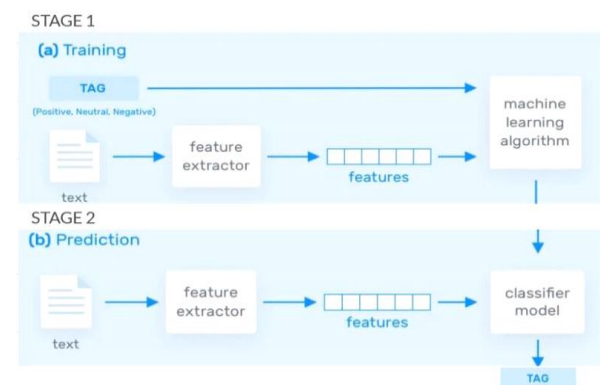


Fig. 4.  Sentiment Analysis

Even data like medical records, attendance, grades and teachers remark from student portals is used and processed for feature extraction. Logistic regression modelling is used for this analysis.

### Implementation Infrastructure

Proposed Solution can easily be integrated and implemented in schools, colleges, institutes and other places prone to bullying. A data-pipeline can be built  for fetching data at real-time with predictive analytical capabilities. Infrastructure for

such a solution is a one-time investment and will provide a great benefit to the coming future making schools bullying free, a dream most people never thought would become a reality.

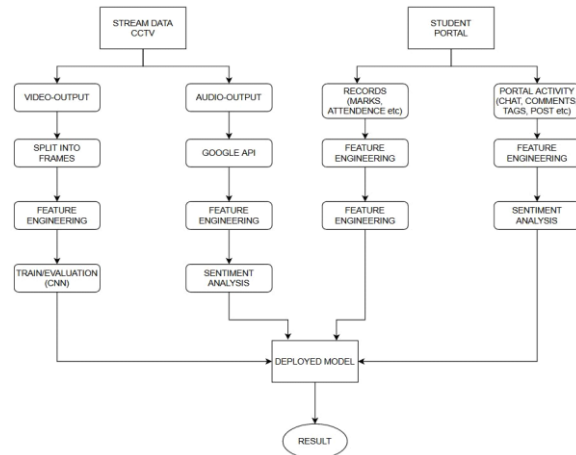Detailed information and diagrammatic representation of Recommended Architecture is provided below:



Fig. 5. Data Flow/Infrastructure Architecture Flow diagram

As shown in the diagram, data streams from CCTV is split as audio and video media streams. Audio is then sent to Google speech-to-text API which is used for sentiment analysis. In case of video, clips are sliced down into multiple frames based on time and frame rates. These images are analysed to detect any signs of bullying or any other sort of unethical action via the use of defined CNN model.

From student portals, web scraping is implemented to obtain information about posts, tags, comments, open-chats and for records which can be taken directly from the school database. All these attributes are fetched to be deployed into the model (on cloud premises) and results can be fetched remotely via internet.

***Table 3: Feature Scaled Values***



### Result

To analyse the accuracy in initial phases, labelled data was split into training and testing data in ratio of 7:3 respectively. At initial stages, with the gathered data our model is able to predict with an accuracy of 87 %.

The charts below depict cross-entropy loss which helps in measuring the performance of the classification model. The output represents a probability value between 0 and 1 [6] which shows our model's classification accuracy.
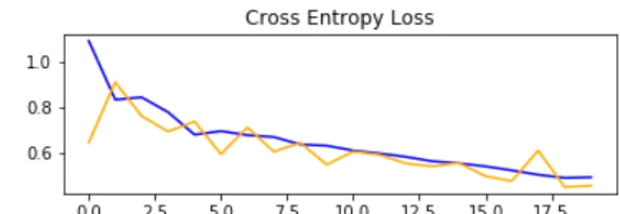


Fig. 6. Cross Entropy Loss (train dataset is represented by blue curve and test dataset is represented by orange curve)
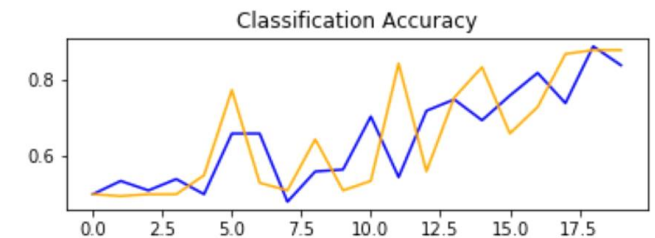


Fig. 7. Classification Accuracy (train dataset is represented by blue curve and test dataset is represented by orange curve)



Fig. 8. Confusion Matrix (This matrix is plotted on the labelled data having sample space of 2000 to check the accuracy of the deployed model)

### Improvements

Future Improvements can include real-time identification of bullies through facial recognition. Sending of alerts to management in case of any real time bullying incident. The model can also be used as a measure to check for depressed victims so as to avoid suicide and self-harm tendencies and help in countering it.

## Conclusion

The proof of concept solution developed by the authors was successful in identifying the bullying situation and hence the bully and the victim. The output was more accurate for audio-visual inputs. For the unstructured data, the accuracy of prediction will have to be improved by inter linking of the outputs with the audio visual parameters. Overall it was a successful proof of concept.

## Acknowledgment

We would like to thank Professor (HAG) Dr. Nidul Sinha of National Institute of Technology, Silchar for his guidance during the research.

## References:

[1]. National Academies of Sciences, Engineering, and Medicine. 2016. Preventing Bullying Through Science, Policy, and Practice. Washington, DC: The National Academies Press. https://doi.org/10.17226/23482.

[2]. D. Poeter. (2011) Study: A Quarter of Parents Say Their Child Involved in Cyberbullying. pcmag.com. [Online]. Available: http://www.pcmag.com/article2/0,2817,2388540,00.asp.

[3]. StopBullying.gov. (2020). What Is Bullying. [online] Available at: https://www.stopbullying.gov/bullying/what-is-bullying [Accessed 28 Jan. 2020].

[4]. Liu, B., 2020. Opinion Mining, Sentiment Analysis, Opinion Extraction. [online] Cs.uic.edu. Available at: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#lexicon> [Accessed 16 June 2020].

[5]. WildML. (2020). Understanding Convolutional Neural Networks for NLP. [online] Available at: .

http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/ [Accessed 28 Jan. 2020].

[6]. Brownlee, J. (2020). How to Classify Photos [online] Machine Learning Mastery. Available at: https://machinelearningmastery.com/how-to-develop-a-convolutional-neural-network-to-classify-photos-of-dogs-and-cats/ [Accessed 28 Jan. 2020].

[7]. J. Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification.

## Biographies:

Lalit Kumar is an associate consultant with Gazelle Information Technologies, New Delhi. He holds a bachelors in technology with special interests in AI and using AI for practical applications. He has built practical solutions for businesses to predict incidents in the supply chain using various algorithms.

Palash Goyal is currently a student of class XII at Mount Carmel School, Sector 23, Dwarka, New Delhi. He started coding at an early age and is currently working on Python, Automation and Image Processing. It was his personal experience at one of his previous schools, that led to the birth of this idea last year, which he converted into reality with the help of other authors.

Karan Malik is a 2nd year student of B. Tech Computer Science at University School of Information, Communication and Technology (USICT), Indraprastha University, New Delhi. His area of interests includes machine learning, image processing and natural Language Processing.

Rishav Kumar is currently a student of class XII at Mount Carmel School, Sector 23, Dwarka, New Delhi. His areas of interest include AI and internet of things.