

NLPIR 2024
8th International Conference on Natural Language Processing and Information Retrieval

JULIA HEINE, KRISTINA SCHAAFF AND TIM SCHLIPPE

INVESTIGATING TEXT AND IMAGE FEATURES FOR THE DETECTION OF CYBERBULLYING

Okayama, Japan
December 14, 2024

AGENDA

Introduction

1

Related Work

2

Experimental Setup

3

Experiments and Results

4

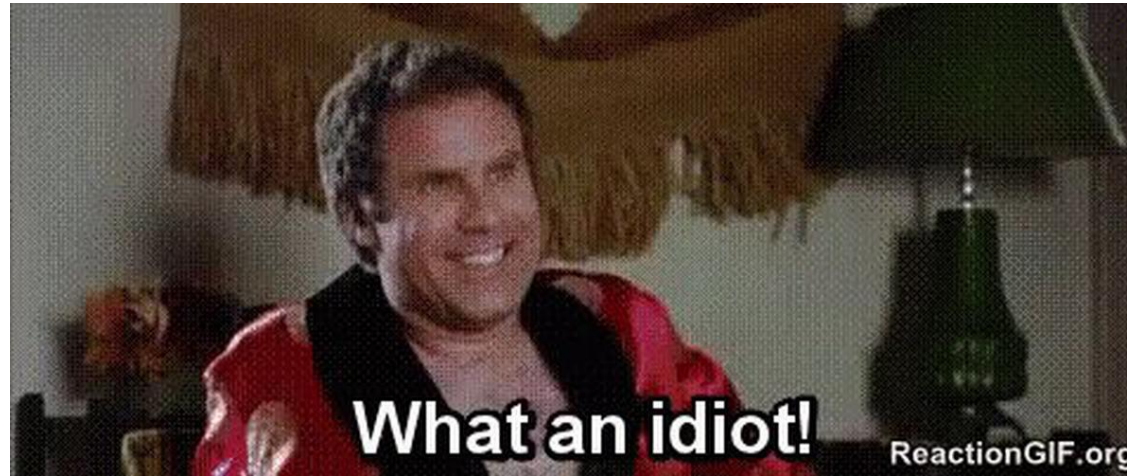
Conclusion and Future Work

5

1

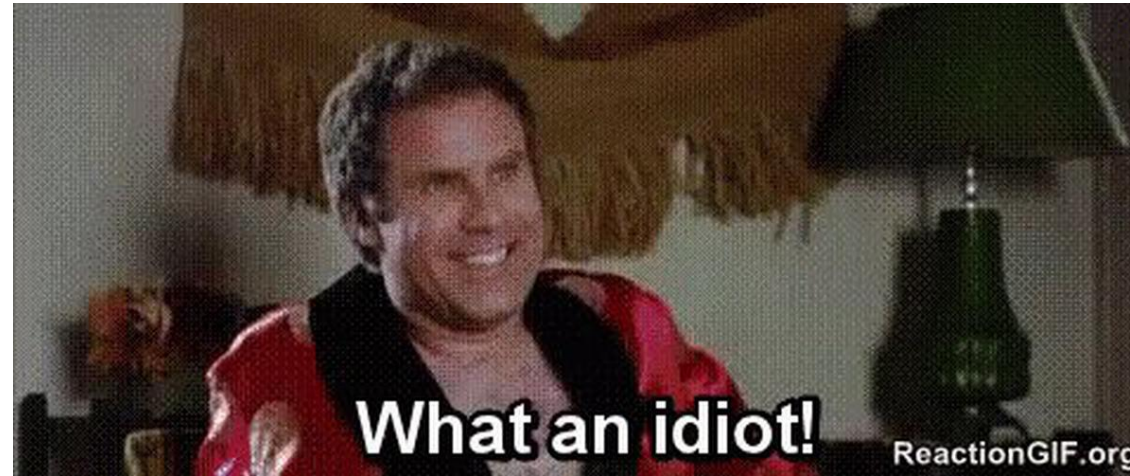
INTRODUCTION

MOTIVATION



~50% OF ALL TEENAGERS HAVE EXPERIENCE WITH CYBERBULLYING

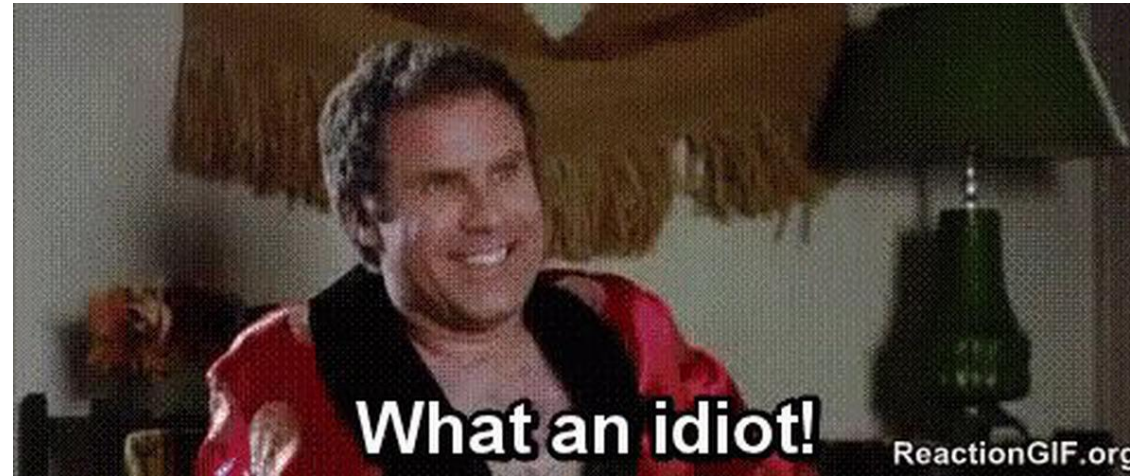
MOTIVATION



~50% OF ALL TEENAGERS HAVE EXPERIENCE WITH CYBERBULLYING

➔ trouble sleeping, difficulty concentrating, or feelings of anxiety

MOTIVATION

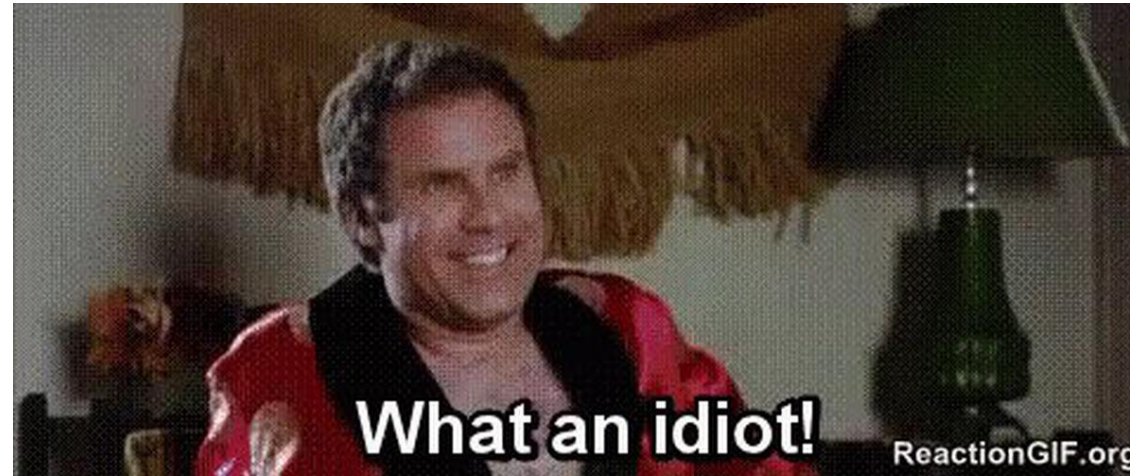


~50% OF ALL TEENAGERS HAVE EXPERIENCE WITH CYBERBULLYING

➔ trouble sleeping, difficulty concentrating, or feelings of anxiety

CYBERBULLYING CAN BE MULTIMODAL!

MOTIVATION



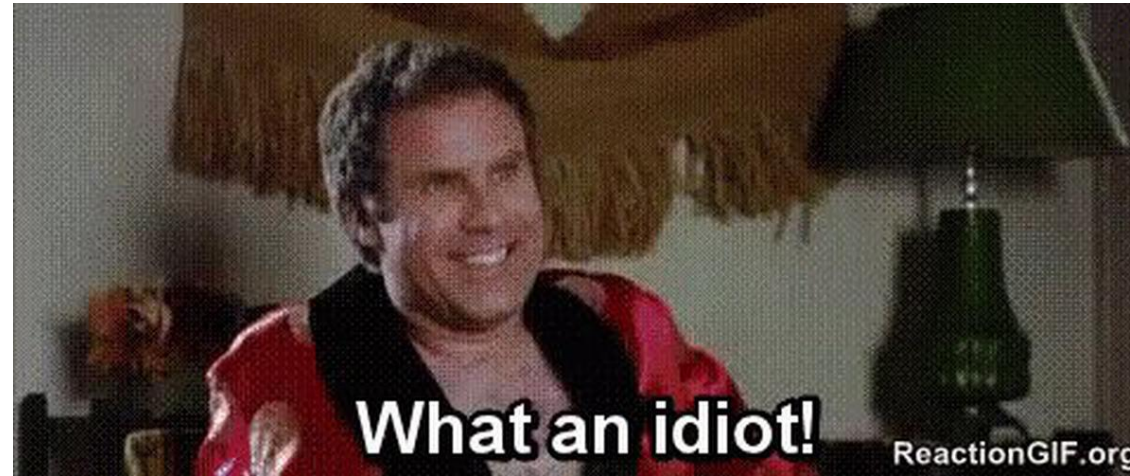
~50% OF ALL TEENAGERS HAVE EXPERIENCE WITH CYBERBULLYING

➔ trouble sleeping, difficulty concentrating, or feelings of anxiety

CYBERBULLYING CAN BE MULTIMODAL!

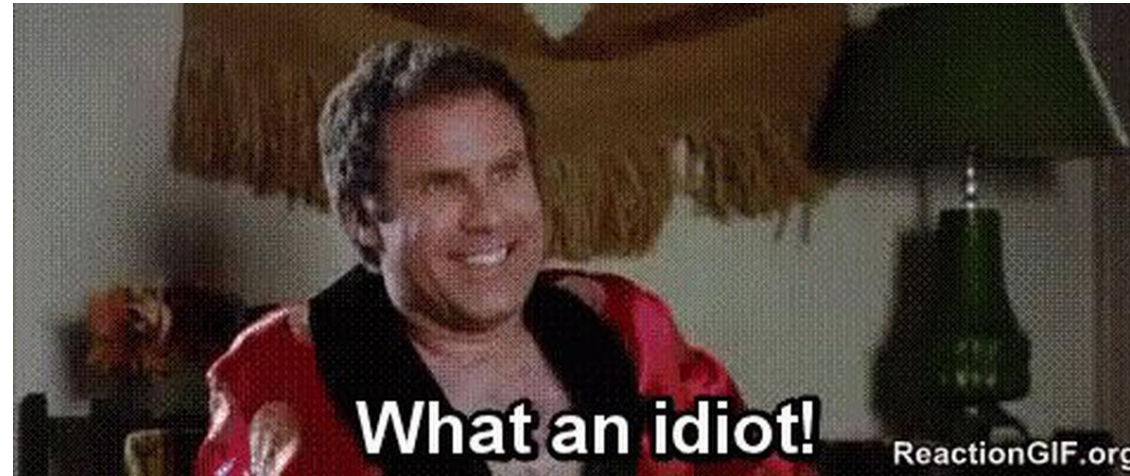
➔ text, images, videos

MOTIVATION



CHALLENGE: TOO MUCH CONTENT FOR MANUAL DETECTION

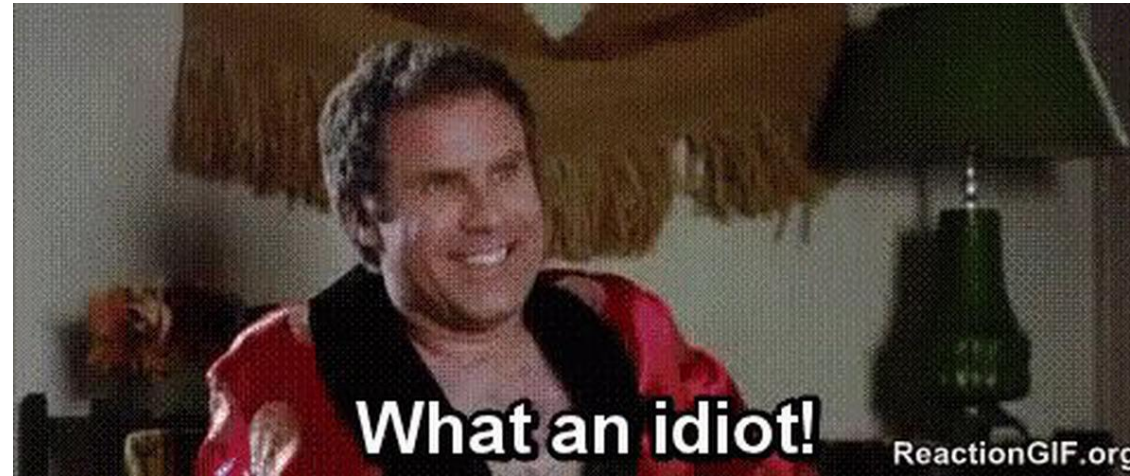
MOTIVATION



CHALLENGE: TOO MUCH CONTENT FOR MANUAL DETECTION

➔ volume of content shared online makes it impossible to manually monitor and address every violation

MOTIVATION

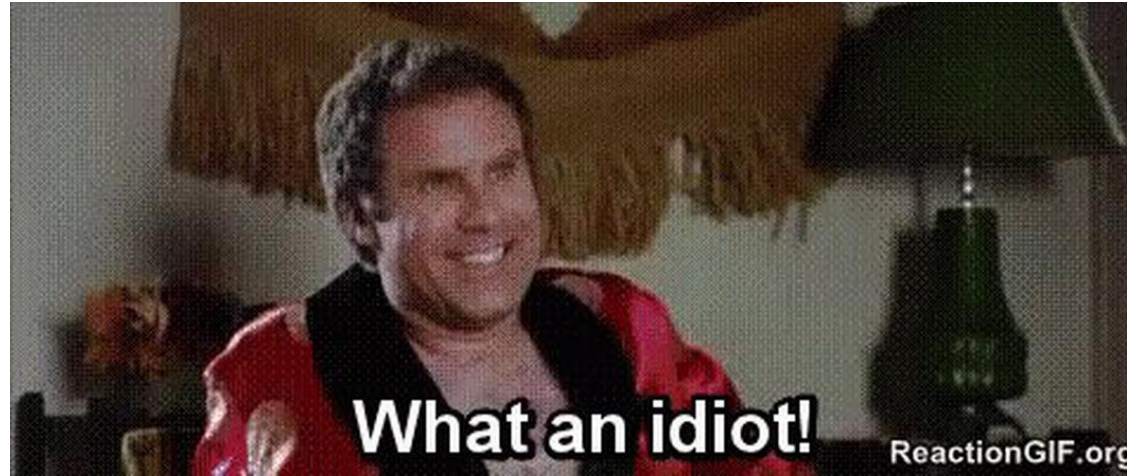


CHALLENGE: TOO MUCH CONTENT FOR MANUAL DETECTION

➔ volume of content shared online makes it impossible to manually monitor and address every violation

SOLUTION: AI TO DETECT CYBERBULLYING AUTOMATICALLY

MOTIVATION



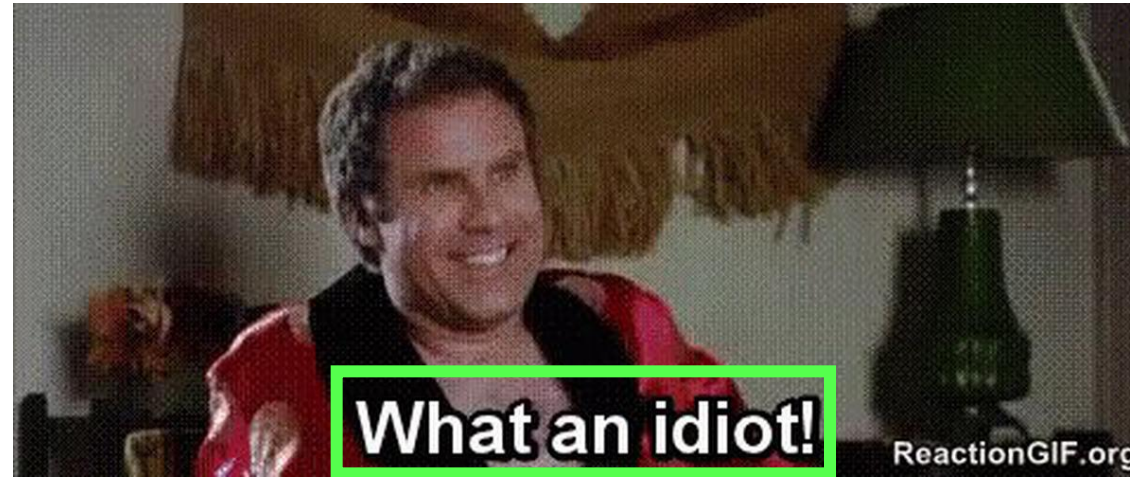
CHALLENGE: TOO MUCH CONTENT FOR MANUAL DETECTION

➔ volume of content shared online makes it impossible to manually monitor and address every violation

SOLUTION: AI TO DETECT CYBERBULLYING AUTOMATICALLY

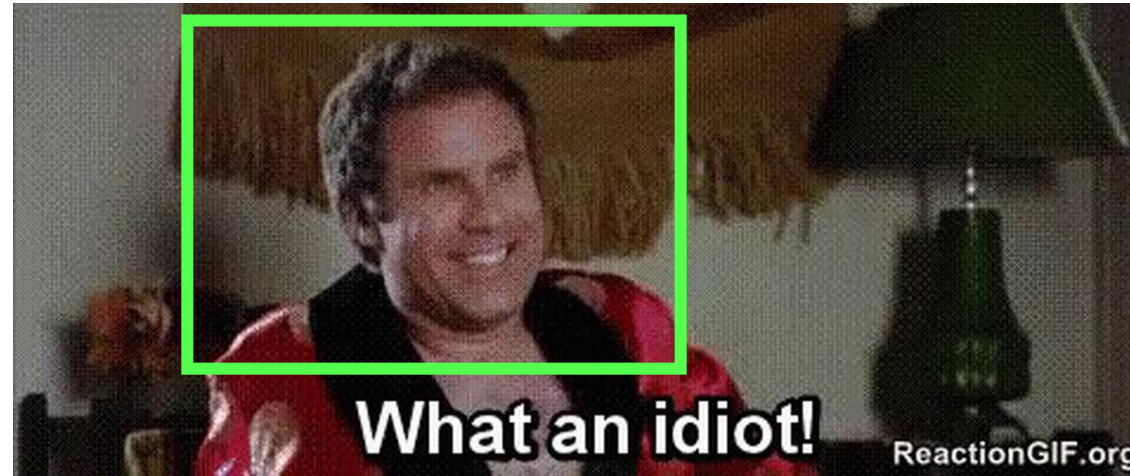
➔ automatically detect cyberbullying and address it more effectively and at scale

MOTIVATION: RESEARCH QUESTIONS



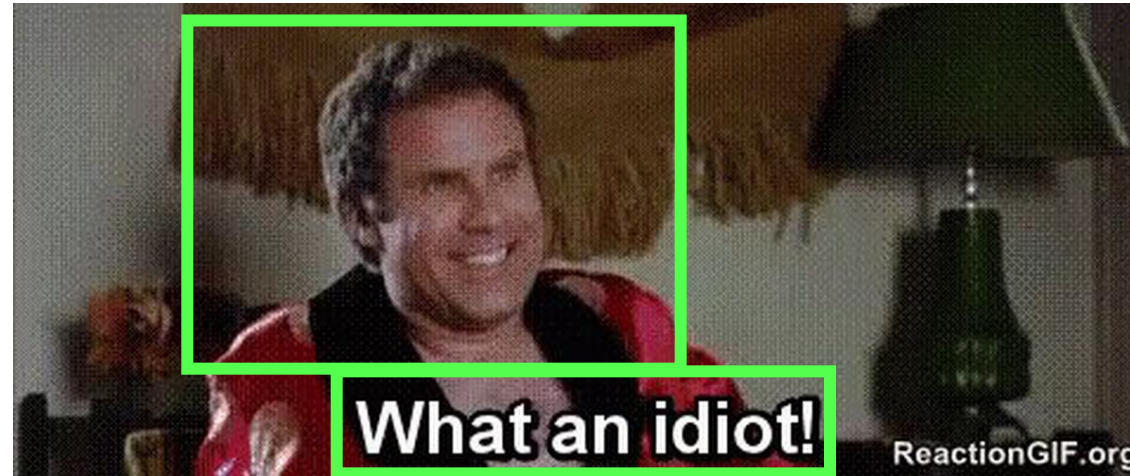
1. WHAT TEXT FEATURES HELP TO IMPROVE CYBERBULLYING DETECTION?

MOTIVATION: RESEARCH QUESTIONS



2. WHAT IMAGE FEATURES HELP TO IMPROVE CYBERBULLYING DETECTION?

MOTIVATION: RESEARCH QUESTIONS



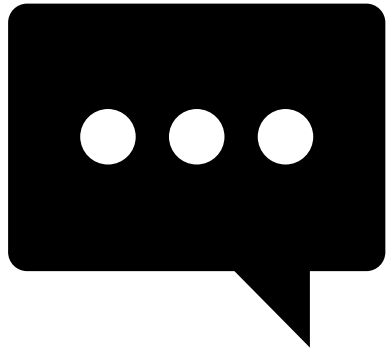
3. WHAT IS THE PERFORMANCE OF THE COMBINATION OF TEXT AND IMAGE FEATURES FOR CYBERBULLYING DETECTION?

2

RELATED WORK

RELATED WORK

Text features are more investigated



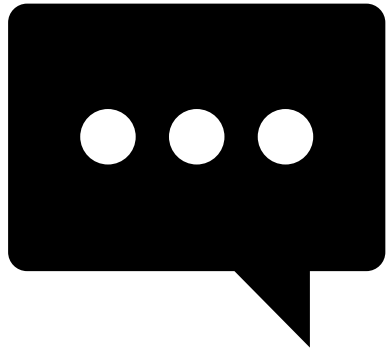
VS.



e.g., Batani et al. (2022), Islam et al. (2020), Raj et al. (2021), Alkomah et al. (2022), Atoum (2020), Van Bruwaene et al. (2020)

RELATED WORK

Text features are more investigated

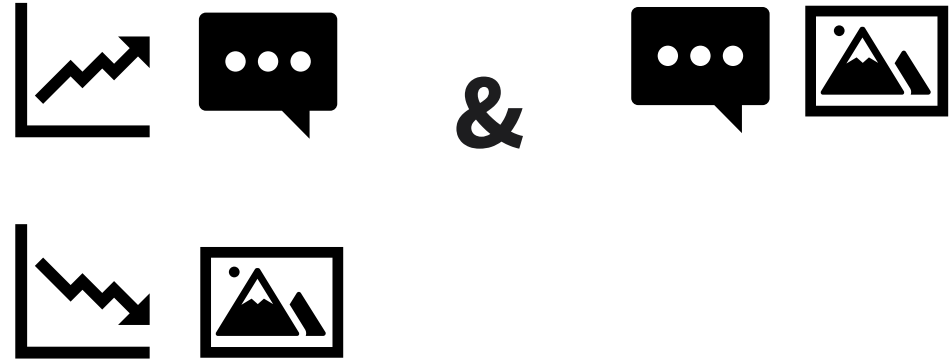


VS.



e.g., Batani et al. (2022), Islam et al. (2020), Raj et al. (2021), Alkomah et al. (2022), Atoum (2020), Van Bruwaene et al. (2020)

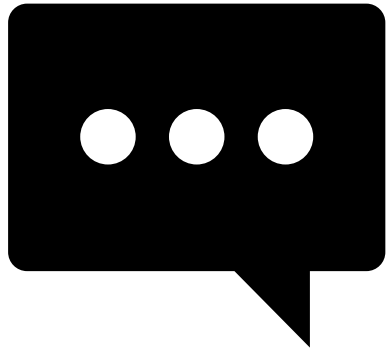
Image Models are not that accurate



e.g., Hosseinmardi et al. (2016), Gomez et al. (2019), Botelho et al. (2021), Sahu et al. (2021), Paul et al. (2022)

RELATED WORK

Text features are more investigated

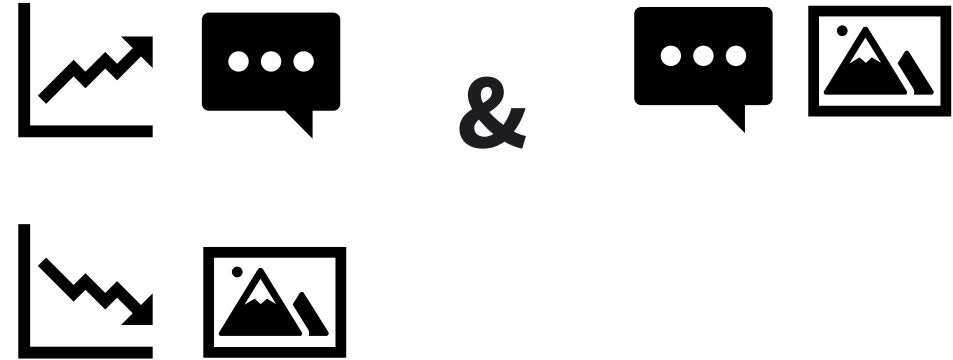


VS.



e.g., Batani et al. (2022), Islam et al. (2020), Raj et al. (2021), Alkomah et al. (2022), Atoum (2020), Van Bruwaene et al. (2020)

Image Models are not that accurate



e.g., Hosseinmardi et al. (2016), Gomez et al. (2019), Botelho et al. (2021), Sahu et al. (2021), Paul et al. (2022)

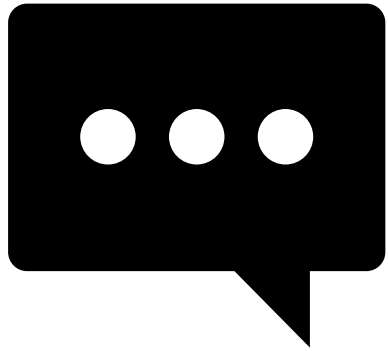
Cyberbullying = Hate Speech



e.g., Batani et al. (2022)

RELATED WORK

Text features are more investigated

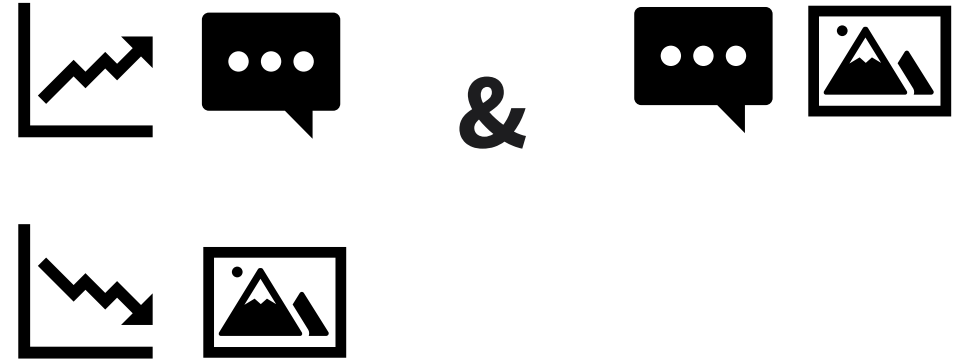


VS.



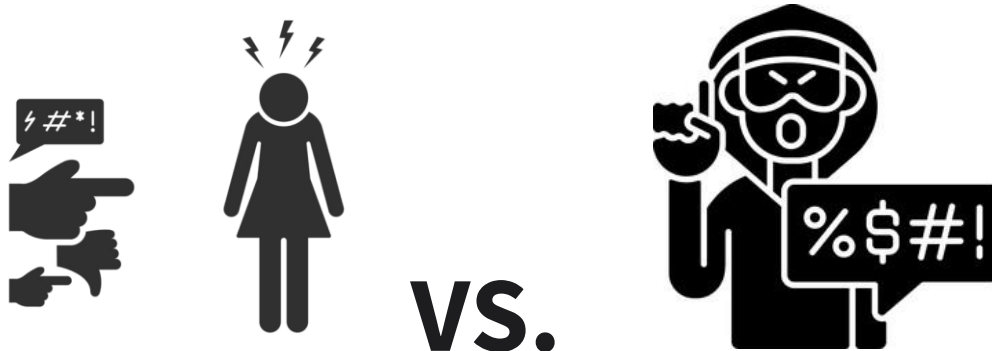
e.g., Batani et al. (2022), Islam et al. (2020), Raj et al. (2021), Alkomah et al. (2022), Atoum (2020), Van Bruwaene et al. (2020)

Image Models are not that accurate



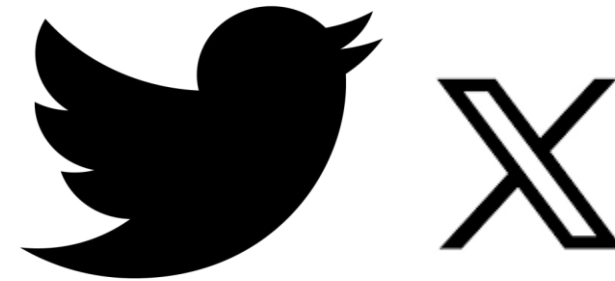
e.g., Hosseinmardi et al. (2016), Gomez et al. (2019), Botelho et al. (2021), Sahu et al. (2021), Paul et al. (2022)

Cyberbullying = Hate Speech



e.g., Batani et al. (2022)

Most datasets are based on Twitter



e.g., Gomez et al. (2019), Mahlangu & Tu (2019), Yadav et al. (2020), Mahat (2021), Roy et al. (2020), Bharti et al. (2021)

EXPERIMENTAL SETUP

EXPERIMENTAL SETUP: GOAL & APPROACH

GOAL

Develop a simple and structured approach to answer research questions

EXPERIMENTAL SETUP: GOAL & APPROACH

GOAL

Develop a simple and structured approach to answer research questions

APPROACH

Multimodel dataset

→ MMHS150K dataset (Gomez et al., 2019)

EXPERIMENTAL SETUP: GOAL & APPROACH

GOAL

Develop a simple and structured approach to answer research questions

APPROACH

Multimodel dataset

→ MMHS150K dataset (Gomez et al., 2019)

Binary Classification Problem

→ Convert 6-class dataset into binary labels

EXPERIMENTAL SETUP: GOAL & APPROACH

GOAL

Develop a simple and structured approach to answer research questions

APPROACH

Multimodal dataset

→ MMHS150K dataset (Gomez et al., 2019)

Binary Classification Problem

→ Convert 6-class dataset into binary labels

Balanced Dataset

→ Ensure equal distribution of *bullying* and *non-bullying* samples across train, validation, and test sets

EXPERIMENTAL SETUP: GOAL & APPROACH

GOAL

Develop a simple and structured approach to answer research questions

APPROACH

Multimodel dataset

→ MMHS150K dataset (Gomez et al., 2019)

Binary Classification Problem

→ Convert 6-class dataset into binary labels

Balanced Dataset

→ Ensure equal distribution of *bullying* and *non-bullying* samples across train, validation, and test sets

Simple Classifier

→ Multi-Layer Perceptron (MLP)

EXPERIMENTAL SETUP: FEATURE DESIGN & TESTING

FEATURE DESIGN

- Combine features from Hate Speech and Cyberbullying literature
- Create new features

EXPERIMENTAL SETUP: FEATURE DESIGN & TESTING

FEATURE DESIGN

- Combine features from Hate Speech and Cyberbullying literature
- Create new features

TESTING MODELS

Single-Feature Models

→ Evaluate performance of individual features

EXPERIMENTAL SETUP: FEATURE DESIGN & TESTING

FEATURE DESIGN

- Combine features from Hate Speech and Cyberbullying literature
- Create new features

TESTING MODELS

Single-Feature Models

→ Evaluate performance of individual features

Multi-Feature Models

→ Test combinations of features to improve performance

EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features



Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

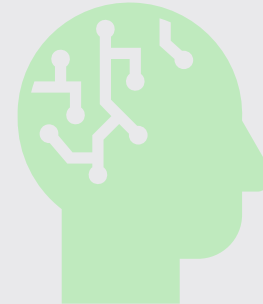
List Lookup Features

Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features



AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features

1010
Sentence Vector
1010

Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

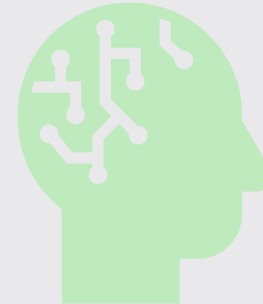
List Lookup Features

Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features

Negative
Neutral
Positive

AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features



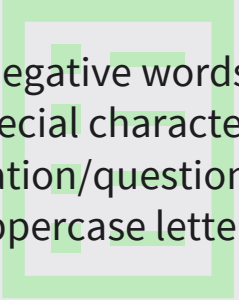
Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

List Lookup Features

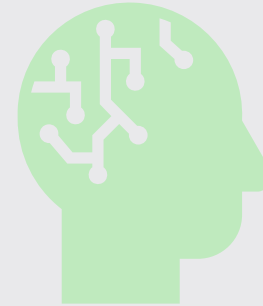


Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features



AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features



Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

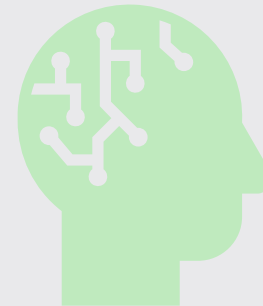
List Lookup Features

Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features



AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features

1010
Sentence Vector
1010

Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

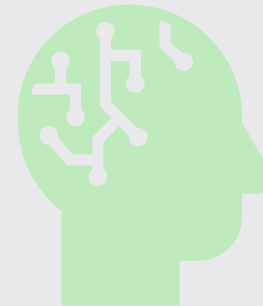
List Lookup Features

Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features

Negative
Neutral
Positive

AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features

1010
Sentence Vector
1010

Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

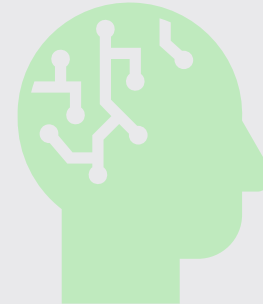
List Lookup Features

Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features

Negative
Neutral
Positive

AI Feedback Features



EXPERIMENTAL SETUP: TEXT FEATURES

Text Vector Features



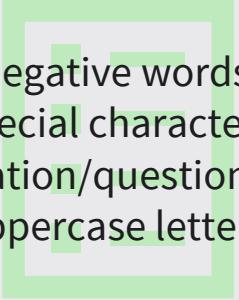
Error-Based Features



Semantic Features

Objectivity vs. Subjectivity
Sentiment
Sarcasm
Toxicity

List Lookup Features

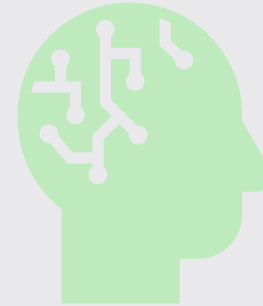


Negative words
Special characters
Exclamation/question marks
Uppercase letters

Emoji Features



AI Feedback Features



EXPERIMENTAL SETUP: IMAGE FEATURES

Image Vector Features

1010
Image Vector
1010

Image Text Features

1010
Sentence Vector
1010

Image Description Features

1010
1st best description
1st best description + Probability
5 best descriptions
1010

EXPERIMENTAL SETUP: IMAGE FEATURES

Image Vector Features

1010
Image Vector
1010

Image Text Features

1010
Sentence Vector
1010

Image Description Features

1010
1st best description
1st best description + Probability
5 best descriptions
1010

EXPERIMENTAL SETUP: IMAGE FEATURES

Image Vector Features

1010
Image Vector
1010

Image Text Features

1010
Sentence Vector
1010

Image Description Features

1010
1st best description
1st best description + Probability
5 best descriptions
1010

EXPERIMENTAL SETUP: IMAGE FEATURES

Image Vector Features

1010
Image Vector
1010

Image Text Features

1010
Sentence Vector
1010

Image Description Features

1010
1st best description
1st best description + Probability
5 best descriptions
1010

EXPERIMENTS AND RESULTS

RESULTS

Feature Category	F1
<i>Text Vector Features</i>	66.25%
<i>Semantic Features</i>	56.24%
- <i>sentimentscore</i>	39.92%
- <i>subjectivityscore</i> *	43.43%
- <i>sarcasmscore</i> *	53.21%
- <i>toxiciyscore</i> *	54.41%
<i>List Lookup Features</i>	51.90%
- <i>capitalizedCharsrel</i>	15.17%
- <i>negWordsrel</i>	42.10%
- <i>specialCharsrel</i> *	66.30%
<i>Emoji Features</i>	56.53%
- <i>negEmojiscore</i>	66.29%
- <i>neuEmojiscore</i>	66.39%
- <i>posEmojiscore</i>	65.54%
<i>Error-Based Features</i> *	56.76%
<i>AI Feedback Features</i> *	64.45%
<i>Text Features_{all}</i>	65.54%
<i>Text Features_{F1>50%}</i>	66.08%

* NEW FEATURES

RESULTS

Feature Category	F1
<i>Text Vector Features</i>	66.25%
<i>Semantic Features</i>	56.24%
- <i>sentimentscore</i>	39.92%
- <i>subjectivityscore</i> *	43.43%
- <i>sarcasm_score</i> *	53.21%
- <i>toxicityscore</i> *	54.41%
<i>List Lookup Features</i>	51.90%
- <i>capitalizedChars_{rel}</i>	15.17%
- <i>negWords_{rel}</i>	42.10%
- <i>specialChars_{rel}</i> *	66.30%
<i>Emoji Features</i>	56.53%
- <i>negEmoji_{score}</i>	66.29%
- <i>neuEmoji_{score}</i>	66.39%
- <i>posEmoji_{score}</i>	65.54%
<i>Error-Based Features</i> *	56.76%
<i>AI Feedback Features</i> *	64.45%
<i>Text Features_{all}</i>	65.54%
<i>Text Features_{F1>50%}</i>	66.08%

> 66% F1

* NEW FEATURES

RESULTS

Feature Category	F1
<i>Text Vector Features</i>	66.25%
<i>Semantic Features</i>	56.24%
- <i>sentimentscore</i>	39.92%
- <i>subjectivityscore</i> *	43.43%
- <i>sarcasm_score</i> *	53.21%
- <i>toxicityscore</i> *	54.41%
<i>List Lookup Features</i>	51.90%
- <i>capitalizedChars_{rel}</i>	15.17%
- <i>negWords_{rel}</i>	42.10%
- <i>specialChars_{rel}</i> *	66.30%
<i>Emoji Features</i>	56.53%
- <i>negEmoji_{score}</i>	66.29%
- <i>neuEmoji_{score}</i>	66.39%
- <i>posEmoji_{score}</i>	65.54%
<i>Error-Based Features</i> *	56.76%
<i>AI Feedback Features</i> *	64.45%
<i>Text Features_{all}</i>	65.54%
<i>Text Features_{F1>50%}</i>	66.08%

> 60% F1

* NEW FEATURES

Feature Category	F1
<i>Image Vector Features</i>	65.27%
<i>Image Text Features</i>	31.57%
<i>Image Description Features</i>	60.92%
- <i>firstBestDescription</i>	25.11%
- <i>firstBestDescriptionProb</i> *	62.37%
- <i>fiveBestDescription</i> *	28.72%
<i>Image Features_{all}</i>	30.39%
<i>Image Features_{F1>50%}</i>	56.10%

RESULTS

Feature Category	F1
<i>Text Vector Features</i>	66.25%
<i>Semantic Features</i>	56.24%
- <i>sentimentscore</i>	39.92%
- <i>subjectivityscore</i> *	43.43%
- <i>sarcasm_score</i> *	53.21%
- <i>toxicityscore</i> *	54.41%
<i>List Lookup Features</i>	51.90%
- <i>capitalizedChars_{rel}</i>	15.17%
- <i>negWords_{rel}</i>	42.10%
- <i>specialChars_{rel}</i> *	66.30%
<i>Emoji Features</i>	56.53%
- <i>negEmojiscore</i>	66.29%
- <i>neuEmojiscore</i>	66.39%
- <i>posEmojiscore</i>	65.54%
<i>Error-Based Features</i> *	56.76%
<i>AI Feedback Features</i> *	64.45%
<i>Text Features_{all}</i>	65.54%
<i>Text Features_{F1>50%}</i>	66.08%

Feature Category	F1
<i>Image Vector Features</i>	65.27%
<i>Image Text Features</i>	31.57%
<i>Image Description Features</i>	60.92%
- <i>firstBestDescription</i>	25.11%
- <i>firstBestDescriptionProb</i> *	62.37%
- <i>fiveBestDescription</i> *	28.72%
<i>Image Features_{all}</i>	30.39%
<i>Image Features_{F1>50%}</i>	56.10%

* NEW FEATURES

Feature Category	F1
<i>Text+Image Features_{all}</i>	64.42%
<i>Text+Image Features_{F1>50%}</i>	65.37%

5

CONCLUSION AND FUTURE WORK

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.
- Best image-based system achieved 65.27% F1 using image vector features.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.
- Best image-based system achieved 65.27% F1 using image vector features.
- Our innovative features provided valuable insights but limited improvement due to task difficulty.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.
- Best image-based system achieved 65.27% F1 using image vector features.
- Our innovative features provided valuable insights but limited improvement due to task difficulty.

Future Work

- Evaluate features with different classifiers.
- Test features on other datasets.
- Extend features to other domains.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.
- Best image-based system achieved 65.27% F1 using image vector features.
- Our innovative features provided valuable insights but limited improvement due to task difficulty.

Future Work

- Evaluate features with different classifiers.
- Test features on other datasets.
- Extend features to other domains.

CONCLUSION AND FUTURE WORK

Conclusion

- We have explored features for distinguishing between bullying and non-bullying messages with images.
- Best text-based systems achieved over 66% F1 on the MMSH150K test set.
- Successful text features: Text vector features, list lookup features, emoji features plus the combination of these features.
- Best image-based system achieved 65.27% F1 using image vector features.
- Our innovative features provided valuable insights but limited improvement due to task difficulty.

Future Work

- Evaluate features with different classifiers.
- Test features on other datasets.
- Extend features to other domains.

THANK YOU

Tim Schlippe

 tim.schlippe@iu.org