

HateProof: Are Hateful Meme Detection Systems really Robust?

Piush Aggarwal*

piush.aggarwal@fernuni-hagen.de
 CATALPA, FernUniversität in Hagen
 Hagen, Germany

Pranit Chawla*

Mithun Das*
 pranitchawla98@iitkgp.ac.in
 mithundas@iitkgp.ac.in
 Indian Institute of Technology
 Kharagpur
 Kharagpur, West Bengal, India

Punyajoy Saha

punyajoy@iitkgp.ac.in
 Indian Institute of Technology
 Kharagpur
 Kharagpur, West Bengal, India

Binny Mathew

binnymathew@iitkgp.ac.in
 Indian Institute of Technology
 Kharagpur
 Kharagpur, West Bengal, India

Torsten Zesch

torsten.zesch@fernuni-hagen.de
 CATALPA, FernUniversität in Hagen
 Hagen, Germany

Animesh Mukherjee

animeshm@cse.iitkgp.ac.in
 Indian Institute of Technology
 Kharagpur
 Kharagpur, West Bengal, India

ABSTRACT

Exploiting social media to spread hate has tremendously increased over the years. Lately, multi-modal hateful content such as memes has drawn relatively more traction than uni-modal content. Moreover, the availability of implicit content payloads makes them fairly challenging to be detected by existing hateful meme detection systems. In this paper, we present a use case study to analyze such systems' vulnerabilities against external adversarial attacks. We find that even very simple perturbations in uni-modal and multi-modal settings performed by humans with little knowledge about the model can make the existing detection models highly vulnerable. Empirically, we find a noticeable performance drop of as high as 10% in the macro-F1 score for certain attacks. As a remedy, we attempt to boost the model's robustness using contrastive learning as well as an adversarial training-based method - *VILLA*. Using an ensemble of the above two approaches, in two of our high resolution datasets, we are able to (re)gain back the performance to a large extent for certain attacks. We believe that ours is a first step toward addressing this crucial problem in an adversarial setting and would inspire more such investigations in the future.

CCS CONCEPTS

• Computing methodologies; • Natural language processing;
 • Social and professional topics; • Censorship;

KEYWORDS

Hateful memes, robustness, multi-modal, social media, ethics, accountability

1 INTRODUCTION

Leveraging the doctrine of freedom of speech [3, 23, 67], misanthropists are spreading hatred in the society using social media platforms. Memes have been widely used in the online world to spread harmful content at an alarming rate. To detect such malicious content, social media companies employ moderators to manually screen posts publicized on their platforms. However, due to the high volume of content dissemination, it has become challenging to

*Equal Contribution.

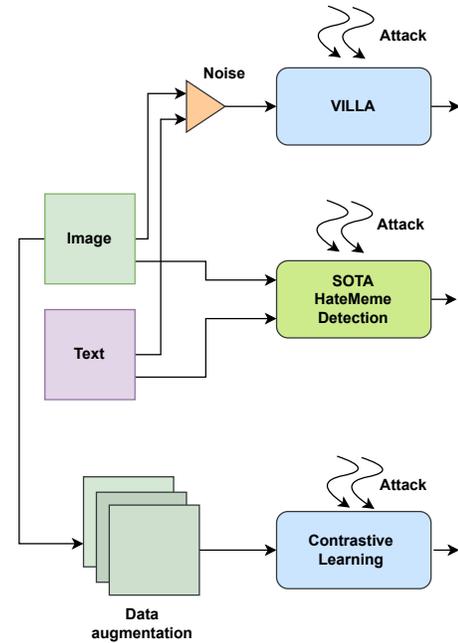


Figure 1: A schematic showing the overall setup of our experiments.

label such contents manually. Hence, automatic moderation techniques are required. With the advancements in NLP and vision technologies machines are now able to interpret the semantics of hate, thereby, facilitating deceleration of the spread of such content [2, 35, 41, 51, 59, 70].

However, systems like these that rely more on machine learning than empirical principles are oftentimes susceptible to unexpected outcomes. A critical difference between multimodal systems and their unimodal counterparts is the fusion mechanism [60]. This fusion mechanism merges multiple input modalities to learn their common representation, which is then processed by numerous fully linked layers (assuming any popular deep learning setup) to predict

classification results depending on the nature of the corresponding downstream tasks. The system adopts multiple strategies to learn strong fusion embeddings of the input modalities. This fusion mechanism presents a new challenge for studying the robustness of these models to adversity. Hateful memes comprise implicit properties, which any of their modalities can acquire. As a result, the possibility of experiencing unexpected events is more. Consequently, performance could be potentially localized and the model may get biased [25] thus posing significant challenges.

Misuse of hate meme detection systems is a common problem similar to deliberate misconduct in AI-based biometric authentication systems [57, 62]. Evtimov et. al. [14] compares the model vulnerabilities toward attacks that are attempted when full knowledge of the model is available to the one having partial or no access to model properties. In this paper, we compose adversarial attacks based on partial knowledge of the model [14] where attackers are assumed to know that the model takes clues from text tokens and image pixels to learn the hate classifier, and propose nine different attacks that can be easily integrated with any existing meme while maintaining the message perception. The attacks are applied to the individual as well as both of the modalities of the meme. One of the key aspects of our analysis is to consider hateful meme detection as an end-to-end task. Therefore, we employ OCR technology to extract meme texts rather than directly using the text provided in the datasets. We find that the state-of-the-art models are highly vulnerable to the attacks designed by us. We, therefore, propose *adversarial training* and *contrastive learning* based countermeasures to tackle such vulnerabilities. Figure 1 illustrates our overall experimental pipeline.

The major contributions of this paper include the following.

- We extend the study in [14] to examine the model’s vulnerability to partial model knowledge-based adversarial attacks. In particular, we propose nine different attacks to investigate the vulnerabilities of the existing hate meme detection systems to human-induced adversarial attacks. We find that all of these models are highly vulnerable with a drop of as high as 10% in macro-F1 performance in certain cases.
- We develop two different countermeasures to tackle such adversarial inputs. Precisely, we implement a variant of SIMCLR [10] which is based on contrastive learning as well as the popular *VILLA* [17] approach based on adversarial training. In the former case, a function learns additional signals from simple augmentations of each data point in the training set and that is blended with hate meme detection model’s objectives through a suitable loss function thus increasing the overall generalization capabilities of the model. In the latter case, in addition to available data points, a model is trained on adversarial examples that are generated by perturbing one modality of the input data while keeping another one unchanged at a time. For image modality, unlike [55] where pixels are perturbed, we consider image-region features. Similarly, the final embedding layer is used for textual tokens, where a minute amount of noise is introduced.
- In addition, we also propose an *ENSEMBLE* of the above two approaches to salvage the best of both of them. We observe that for two of our high resolution datasets, the *ENSEMBLE* based countermeasure is able to successfully tackle various

different forms of adversarial attacks. The *VILLA* approach is a close second after the *ENSEMBLE* approach.

In summary, our analysis¹ reveals the plausible loopholes that potential haters can adopt in order to escape from the detection systems. Given this, we play proactively and propose generalized countermeasures that are efficient under many attack conditions.

2 RELATED WORK

Hate meme detection: Detecting hate memes has become an increasingly popular research topic owing to the alarming growth in hateful memes across different social media platforms. The models often perform multimodal pre-training on huge unsupervised corpus [1] followed by a fine-tuning on a relatively smaller set of supervised hate data. In the case of memes, the early and late fusion of features belonging to each modality is done before generating the final predictions. In the Hateful Memes challenge [28] most effective systems featuring at the top of the leaderboard were heavily relying on large multimodal transformer models such as VL-BERT [2, 70], OSCAR [34], UNITER [35], VISUALBERT [59], and LXMERT [56]. Recently many works have also proposed an ensemble of different visual-linguistic models for obtaining better detection performance [2, 41, 70].

Datasets: Multiple datasets [16, 28, 46, 54] have been built to bring in more diversity in the detection task. In addition, there are other similar datasets that emphasize sentiment, humor, offensiveness, and motive of memes [52]. Apart from English, memes have also been analyzed in other languages datasets [20, 40].

Adversarial attacks: Previous studies raise a concern about model biasing and domain-specific responses [25]. Among images, tampering has been used for the malicious purpose for quite sometime now. Bayar et.al [5] embed noise using image processing functions namely scaling, blurring and introduction of white noise. For automatic image alteration, GAN-based systems are commonly used [39]. Color filter array (CFA) based spectroscopy has been performed to detect pixel level abnormalities [15]. Further, several perturbation-based algorithms have been proposed to attack image modality [7]. Among text-based attacks, [22] added *LOVE* to every text input and found that the nature of the data and the labeling standards are more crucial than the model architecture. The study showed that character-level feature training helped to uplift model robustness. Content code blurring has also been used to both emphasize [21] and forge [69] the textual quality. Attacks also have been developed while considering both modalities. For example, [61] observed that decoupling of image and text components directly defeats the sole purpose of the visual-linguistic models.

In this paper, unlike [14], we specifically focus on the model-independent attacks which humans can introduce with little technical expertise². Attack for instance adding noise to the image and texts in the memes can be easily implemented without prior knowledge of the classification models. Unlike previous work which provides a very abstract understanding of this noise category, we attempt to broaden the attack’s horizon at a fine-grained level and analyze their impact on a variety of hate meme detection model paradigms.

¹<https://github.com/aggarwalpiush/Robust-hatememe-detection>

²They are aware of using internet facilities.

Multi-modal robustness Earlier work in model robustness were generally limited to uni-modal settings with end-to-end adversarial training (AT) [22, 38, 43, 58, 65, 68] where training data was augmented with synthetic data to facilitate provably robust training. For multi-modal cases, Yang et. al. [66] emphasized the robust fusion of modalities rather than end-to-end parameter training. Recent studies show that, by injecting adversarial perturbations into feature space, AT can further improve model generalization on language understanding [27, 70] as well as visual-linguistic tasks [17]. Exploiting contrastive learning to improvise model robustness exhibits impressive performance in text classification tasks [42, 47]. In addition, there are approaches [17, 18, 32] which insert visual-linguistic perturbation into the embedding space of the models to increase the adversarial examples in the input. In this work, we intend to exploit contrastive learning proposed by Chen et.al. [9] to augment the robustness of hate meme classification. In another countermeasure that we develop, unlike [17], which introduces uniformly distributed perturbation in the embedding space of both modalities, we use Gaussian noise based perturbations to generate the adversarial examples to perform the training.

3 DATASETS

In order to analyze the vulnerability of available hate meme classifiers as well as evaluate our countermeasure proposal, we have used three benchmark datasets (see Table 1). Note that we apply the proposed adversarial attacks only on the test set in order to analyze the model robustness.

- Kiela et.al. [28] (FBHM) contains 10,000 memes collected from Getty images and are semi-artificially annotated using benign confounders. The dataset consists of five varieties of memes. These include *multimodal hate* where both modalities possess benign confounders, *unimodal hate* where at least one of the modalities is already hateful, *benign image* as well as *benign text* confounders and *random not-hateful* examples. First four categories are annotated with *hateful* label and rest with *non hateful* label. Since the meme text is also annotated here, no error is induced during the OCR process. The dataset is split into 85% training, 5% development, and 10% test sets. The development and test set are fully balanced with fixed proportions of each variety discussed above.
- Pramanick et.al. [46] (HARMEME) contains COVID-related memes posted in US social media. Some of the query keywords include *Wuhan virus*, *US election*, *COVID vaccine*, *work from home* and *Trump not wearing mask*. Unlike [28], all of the memes are original and shared across social media. Also, the associated textual content is the output from Google Vision API rather than the manual. The resolution from the memes is also preserved. We group the whole dataset into two categories – (i) *hateful* which consists of both *harmful* and *partially harmful* labelled memes provided in the dataset, and (ii) *non hateful* which consists of memes tagged as *non harmful* in the dataset. Out of a total of 3,544 data points, we again use 85%, 5%, and 10% for training, validation and test respectively.

	FBHM [28]	HARMEME [46]	MAMI [16]
Train/Dev/Test	8500/500/1000	3013/177/354	8000/1000/1000
Hate %	37.56	26.21	50
Domain	<i>in the wild</i>	Covid-19/US Election	Misogynistic
Bit depth (Avg.)	9.54	43.90	4.30

Table 1: Dataset statistics.

- Fersini et. al. [16] (MAMI) We focus only on SUBTASK-A of the dataset where memes are either labelled as *misogynist* or *non misogynist*. We relabel the former as *hateful* and the latter as *non-hateful* making it consistent with our setup. The memes have been collected from social media websites and pertain to threads with women as subjects, or having antifeminist content. Such memes having famous women personalities such as Scarlett Johansson, Emilia Clarke, etc. as well as hashtags such as #girl, #girlfriend, #women, #feminist are also considered. Likewise [46], meme text is extracted with Google Vision API. In total, we use a balanced set of 10,000 instances. Out of this, we use 10% posts each for development and test which are randomly stratified.

4 ADVERSARIAL ATTACKS

Szegedy et.al [55] pointed out that neural networks are highly susceptible to minor changes in the input aka adversarial attacks unless they are properly trained to handle such changes. The human eye is almost insensitive to these minor changes; however, they are capable of hindering the prediction probabilities of the neural classifiers. Adversarial attacks can be *whitebox* where internal parameter and gradients of the models are known. However, in the case of *blackbox* attacks, the model is considered as an oracle and attacks are developed based on model confidence and output prediction probabilities. Therefore, in this case, attacks can be implemented in unlimited settings. As corporate models are not disclosed and updated time to time, therefore attacks in even blackbox settings can vary. In this work, we consider attacks based on partial knowledge of the models. Therefore, we compose model independent attacks and categorize them in nine different types including simple attacks such as SaltPepper noise which has shown noticeable effects in the past [13, 26]. Figure 2 illustrates example of attacks chosen for our study.

4.1 Attacks on text

Unlike previous work [14] where edit distance-based attacks were chosen, in our case, we apply textual attacks in such a way that directly affects the recognition quality of OCRs applied to the memes to extract text. Such attacks study the importance of text modality in real-time settings. The different attack forms include

- Add: Adding high polarity token to the input text often misleads the classifier [22]. Taking cue from this paper, we consider the token *LOVE* to have a small font of size five pixels and embed it at random locations of the input image. To establish the generalizability we also repeat the attack for a (semantically) related token *CARESS*.
- Blur: In addition to visual-linguistic models, OCRs in the pre-processing stage also play important role in the overall

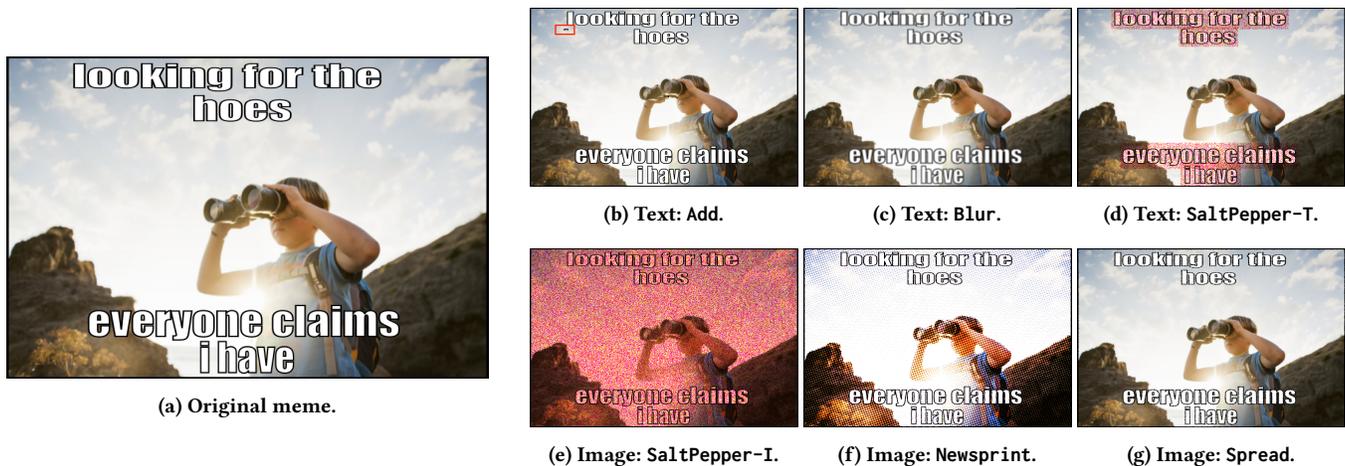


Figure 2: Different forms of adversarial attacks that could be presented to a hate meme detection model.

performance of the hateful meme detection problem. Diminishing the quality of the textual part can hamper the recognition quality of the OCR models. We blur the text using the openCV tool³.

- SaltPepper-T: Likewise blurring, we add salt-and-pepper noise to only the text area of the image. To do so, we extract the bounding box where text is available using easyOCR⁴ python library. Then we randomly add white and black pixels only to the area contained in this box.

4.2 Attacks on image

Cases where the internal behavior of the threat model is not predetermined, the attacker or the adversary cannot use any gradient information to generate appropriate attacks. However, they can introduce random modifications to the image while preserving its message. Therefore we apply image distortions using the GIMP software⁵ in different settings.

- SaltPepper-I: White and black pixels are injected into the whole image in such a way that neither the image perception nor the message conveyed is compromised. We inject the pixels in two different settings (*high* and *low*).
- Newsprint: In this case, the image is poisoned with clustered-dot dithering to get a feel of newspaper printing.
- Spread: Swapping of each pixel with another randomly chosen pixel is done in order to get a slightly jittery output. Here again, both *high* and *low* amounts of spreading have been done.

In addition, we consider combining text as well as image-based attacks. To do so, we choose the common attack scenarios in each modality and aggregate SaltPepper-I (commonly used in photographs [31]) and Add (token *LOVE* in our case [22]) in order to analyze the model vulnerability against joint modality attacks.

5 COUNTERMEASURES

In this section, we discuss two different countermeasures that we employ to tackle such adversarial attacks one of which is based on contrastive learning while the other on adversarial training. We keep the model agnostic of the attack type since otherwise it would learn the properties of certain types only and would lose generalizability. All our attack types are thus unseen to the model.

5.1 Contrastive learning through data augmentation

Data augmentation [48, 49] is the one among those promising methods that are prominently used to improve model generalization and robustness. The main idea is to introduce a small mutation in the original training data and synthetically generate new samples to virtually increase the amount of training data [53]. The points are generally represented as interpolated output of actual training samples. This resolves under-determined and poor generalization problems that occur in the model when they are trained on scarce dataset [36]. During the formalization of augmented samples, human knowledge is needed to describe a neighborhood around each instance in the training data [8]. For instance, during the image classification process, a variety of image transformation operations are performed such as horizontal reflections, slight rotations, and mild scalings. However, sharing the same class by each neighbor examples restricts the models to learn the vicinity relationships across different classes [36]. Therefore, in this work, we consider contrastive learning (CL), which represents a function learned on simple augmentations of each data point in the training set guided via a suitable loss function. The loss function is derived in such a way that it maximizes the agreement among the augmentations of a single data point while spreading apart the augmentations from different data points. The loss function can later be integrated into any classifier to enhance its generalization capabilities [9]. In this work, we apply the *SimCLR* [10] implementation of contrastive learning with crucial revisions to adapt it to the multi-modal settings (memes in our case).

³<https://github.com/openCV/openCV-python>

⁴<https://github.com/JaidedAI/EasyOCR>

⁵<https://www.gimp.org/>



Figure 3: Representative examples of image and text variations used for synthetic augmentation in contrastive learning.

Algorithm 1 Computation of the contrastive loss function \mathcal{L}

Input: batch size N , temperature constant τ , pretrained encoder network f , projection generator g , Augly function A , instance x

```

1: for all  $i \in 1, \dots, N$  do
2:   # Learn representation
3:    $z_i = g(f(x_i))$ 
4:   apply Augly  $x'_i \leftarrow A(x_i)$ 
5:    $h_j = f(x'_i)$ 
6:    $z_j = g(h_j)$ 
7:   # Compute pairwise cosine similarity
8:    $s_{(i,j)} = z_i^T z_j / (||z_i|| ||z_j||)$ 
9:    $s_{(j,i)} = z_j^T z_i / (||z_i|| ||z_j||)$ 
10:  # Define loss function  $\ell(i, j)$ 
11:   $\ell(i, j) = -\log \frac{\exp(s_{(i,j)}/\tau)}{\sum_{n=1}^N \mathbb{1}_{[k \neq i]} \exp(s_{(i,k)}/\tau)}$ 
12: # Contrastive loss function  $\mathcal{L}$ 
return  $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
    
```

To implement CL, the first step would be to obtain synthetic augmentations. Therefore we use the *Augly* [44] python library to obtain such augmentations. In case of the image component, for each input sample we choose one of the noisy versions among blurring, random noising, color jittering, horizontal flipping or gray scaling. Figure 3 illustrates the examples of augmented samples. Similarly to generate textual samples, we use *Augly text*. Three different types of text augmentation that we use are casing alteration, similar character replacement as well as typo induction. In the end, we randomly combine augmented versions from different modalities to produce the related sample for each input instance. We redefine the loss function of the (best performing) state-of-the-art hate meme classification model to incorporate an additional contrastive loss estimated from the augmented data. Precisely, in addition to the cross entropy loss used during the training of the classification model, we add the contrastive loss function computed in Algorithm 1. We calculate the encoder network’s projection g (aka the embedding) for each input instance x and its augmented version $A(x)$. Next we calculate the cosine

similarity between the projections. Subsequently, contrastive loss function is calculated as shown in line 11 of the algorithm where the numerator encapsulates the similarity of a given instance (aka positive example) and the denominator encapsulates the similarity between pairs of remaining instances of the batch which act as negative examples. $\mathbb{1}_{[k \neq i]}$ is an indicator function which is 1 for negative examples and τ denotes the temperature parameter.

5.2 Adversarial training with VILLA

Adversarial training has shown a promising role to counter adversarial attacks in machine learning models. Earlier analysis of model robustness has been limited to performance on surprising altercations appearing in relatively clean data [6, 55]. However, robustness against deliberate obfuscation activities, for instance human induced adversarial attacks on hate meme, has recently emerged as a critical problem. A robust classifier should be able to accurately identify adversarially poisoned images. Training the model with augmented samples as discussed in the previous section can be helpful in order to surpass the vulnerabilities of the existing hate meme detection models. However, due to limited flexibility in explicit augmented example generation, the robustness can only be increased up to a certain extent. Therefore, here, we attempt to generate adversarial examples by performing perturbation in the embedding space [70] of both the image and the text for each training input instance. For textual and image parts, we use BERT and regional embeddings respectively [19]. To do so, we implement the *VILLA* framework [17]⁶ where one modality is perturbed keeping the other one unchanged at a time. For image modality, unlike [55] where pixels are perturbed, we consider image-region features. Similarly for textual tokens, final embedding layer is used where a small amount of noise is introduced. Small perturbations in the feature space keep the instances within the classification boundary and act as a form of latent augmentation. Following prior work [70], we use uniform distribution based perturbation at embedding level to produce adversarial examples during pre-training (task agnostic) and fine-tuning (task specific) stages. In addition, for our work, we use Gaussian noise based perturbation which actually results in a better performance.

⁶<https://github.com/zhegan27/VILLA>

MAMI			FBHM			HARMEME		
UNITER	VISUALBERT	ROB+RESNET	UNITER	VISUALBERT	ROB+RESNET	UNITER	VISUALBERT	ROB+RESNET
84.99	85.00	85.50	63.53	60.88	63.10	78.51	75.57	73.67

Table 2: Performance (macro-F1 scores) of the vanilla hate meme detection models.

6 EXPERIMENTAL SETUP

In this section we discuss the overall experimental setup comprising a description of the hateful meme detection models, the training framework and the hyperparameters.

6.1 Hate meme detection models

We select a series of state-of-the-art hate meme detection models in order to test their vulnerability against adversarial attacks. To analyse the model vulnerability, we start with reproducing results for hate meme detection task. These models are discussed in brief for better readability.

VISUALBERT: This [33] builds up on the BERT architecture [12] and exploits the attention layer to align the input meme text with its image regions which are generated from variety of object detectors. Task specific MLM is used for pretraining. Unlike the original implementation where token ids are persistent during fine-tuning, we use the Vilio⁷ framework where weights for token type are trained from scratch which has been shown to exhibit promising performance in hate meme detection task [41]. Likewise we also use an additional classification head with very high learning rate (500 times the original) and multi-sample dropout.

UNITER: Like VISUALBERT, the UNiversal Image-TExt Representation (UNITER) model [11] builds on an early fusion approach and is pre-trained on large text-image datasets. Visual features extracted from faster R-CNN [50] and textual ones from word piece encodings [64] are combined through a transformer based architecture to align them in a shared embedding space. The model emphasizes on fine grained alignment between the modalities with conditional masking and optimal transport [45]. We update the activation functions and embedding calculations as noted in [41].

ROB+RESNET: Our proposal is to adopt a late fusion approach so that the individual contributions of two modalities are better captured. To this purpose we suggest separate extraction pipelines for the image and the text features. To extract the image features we use RESNET - a very deep residual learning framework for image feature generation [24]. To obtain the text representation, we use the popular RoBERTa [37] model. Finally, we concatenate the features from both the modalities and pass it through a feed-forward network having 128 hidden layers with ReLU activation and dropout = 0.2 to generate the predictions. We train the model for 30 epochs on Adam optimizer [29] with learning rate = 10^{-5} and weight decay = 0.1.

COUNTERMEASURES: In addition to the vanilla implementation of the above-mentioned models, we repeat our experiments in presence of the two countermeasures discussed in the previous section to combat the effect of attacks while keeping the performance. For VILLA, we use three different variants. These include VILLA-PT-UN where adversarial training is only performed during the

pre-training phase with uniform noise, VILLA-FT-UN and VILLA-FT-GN where additional training is done at the fine-tuning stage with uniform and Gaussian noise respectively. In all the variants, the default number of adversarial steps is kept at 3. Due to the requirement of huge computation, we limit ourselves to not doing Gaussian noise-based perturbation at the pre-training stage. For contrastive learning, the loss is calculated for each modality and later on added to the loss function of the hate meme classification model. We test this for $\tau = 0.5$ and call it CL-IND-0.5. We ensemble both the countermeasures (call it ENSEMBLE) by taking averages of their prediction probabilities.

6.2 Model training

To study the robustness of hate meme detection models, performance must be evaluated before and after the application of adversarial attacks. Therefore we start with introducing the attacks on the test instances for each dataset. Consequently for each dataset, we generate 10 different poisoned test sets including the original one.

Text pre-processing: First, we extract the text from the memes. To this purpose, we use the paid service API from Google Vision called GoogleOCR⁸. We choose a meme and apply the OCR model directly on it without performing any preprocessing as the OCR is already equipped with its own preprocessor component. For tokenization, we use BERT-base [12] as well as RoBERTa-base tokenizers [37].

Image pre-processing Following [41], we extract the image features using the detectron2 framework [63] because it makes the training process faster without compromising the relevant content. To be precise, we use Faster R-CNN with ResNet-101, using object and attribute annotations from Visual Genome [4]. The pretrained model generates bottom-up attention features corresponding to salient image regions. Minimum and Maximum boxes are set to 36 for such region of interests.

Both text and image features are then fed to the hate meme detection models in accordance to their architecture requirements.

Hyperparameters: Apart from the COUNTERMEASURES, for rest of the detection models, we keep the default hyperparameters as provided in [41] implementation. Since VILLA need to be trained on adversarial examples generated from each input steps we tune the variable which is responsible for setting the limits of the number examples. We choose this variable as 1, 2, 3, 4, 5, 6, 10, and 100. Generation of adversarial examples also depends upon the type of noise, i.e., uniform or Gaussian noise. In case of contrastive learning, we tune the temperature coefficient τ which is responsible for distinguishing positive and negative samples. We perform our experimentation with $\tau = 0.5$ as has been noted earlier.

⁷<https://github.com/Muennighoff/vilio>

⁸<https://cloud.google.com/vision/docs/ocr>

	MAMI			FBHM			HARMEME		
	UNITER	VISUALBERT	ROB+RESNET	UNITER	VISUALBERT	ROB+RESNET	UNITER	VISUALBERT	ROB+RESNET
Add	0.19	0.46	0.11	0.13	-0.60	0.35	-1.48	-2.90	-0.17
Blur	1.29	1.54	0.52	-1.44	0.69	-0.19	9.99	9.42	5.09
SaltPepper-T	0.09	0.77	1.24	-1.64	0.43	-1.10	3.89	2.24	-2.84
SaltPepper-I-Low	2.29	3.62	6.10	0.34	1.19	1.61	4.89	3.70	4.52
SaltPepper-I-High	4.42	5.50	5.11	2.04	2.65	3.92	16.62	12.02	8.27
Newsprint	6.81	8.57	14.47	3.29	3.80	0.54	7.86	7.52	4.58
Spread-Low	1.19	1.46	1.21	0.52	1.82	2.34	13.39	11.54	16.56
Spread-High	5.14	4.20	8.78	2.70	6.65	3.49	27.66	28.20	25.44
Add+SaltPepper-I	5.30	6.41	16.47	1.40	3.09	3.15	6.61	5.36	4.98
Average	2.97	3.61	6.00	0.82	2.19	1.57	9.94	8.57	7.38

Table 3: % Δ in macro-F1 hate meme detection systems are exposed to different adversarially modified test inputs. The $+ve$ values (highlighted with gradient **RED**) indicate the amount of performance drop when an attack has been applied while $-ve$ values (highlighted with **GREEN**) indicate performance improvements over the baselines.

7 RESULTS

This section is divided into three parts. In the first part, we report the performance of vanilla hate meme detection models. In the second part, we report the change in the performance when the data points are adversarially altered. In the third part, we report the performance (re)gained due to the introduction of the countermeasures. Since two out of the three datasets are unbalanced, we use macro-F1 to measure performance all through.

7.1 Performance of the hate meme detection models

In Table 2 we report the performance of the vanilla hate meme models in terms of macro-F1. All the models perform almost similarly with ROB+RESNET doing slightly better for the MAMI dataset and the UNITER model doing slightly better for the other two datasets. Explicit significant tests (M-W U test) show that none of the models significantly outperform one another.

7.2 Adversarial attacks

To detect vulnerability of the hate meme detection models against adversarial attacks, we replace test set instances with their attacked versions. Finally, each test set have nine different noisy variants which are used for the analysis below as well as for evaluating the proposed countermeasures in the next subsection. Table 3 illustrates the vulnerability of the existing state-of-the-art hate meme detection models with GoogleOCR text recognition model. Each value in the table represents the percentage difference between the macro F1 score before and after application of the attacks. Therefore, greater the positive value, the more severe is the effect of the attack on the model. Almost all models suffer a drop in performance across the different datasets and for different attack schemes. For the MAMI dataset the Newsprint noise strongly affects all the models with ROB+RESNET seeing a drop of as high as $\sim 14.5\%$. The Spread-High noise is next in line among the unimodal attacks. Naturally, the multimodal attack Add+SaltPepper-I also adversely affects all the models with ROB+RESNET again seeing the highest drop of $\sim 16.5\%$. Compared to the other two datasets, for the FBHM dataset all the models suffer less from the attacks. Here, the VISUALBERT

model suffers the most for the Spread-High noise, the ROB+RESNET suffers the most for the SaltPepper-I-High noise and UNITER model suffers the most for the Newsprint noise. For some textual noise, the models also show an improvement in performance which might be attributed to the fact that the training dataset possibly already had some similar adversarial examples by default which allowed the models to learn their characteristics. For the HARMEME dataset we observe that all the models suffer the most. Out of all, the Spread-High noise deteriorates the performance of all the models the most with 27.7%, 28.2% and 25.4% reduction in macro-F1 scores for UNITER, VISUALBERT and ROB+RESNET respectively. On average all models suffer a reduction in the macro-F1 scores due to the introduction of the adversarial attacks with image based attacks being always more severe than text based attacks. Among the image based attacks, SaltPepper-I-High, Spread-High and Newsprint noises affect the model performances the most. For all the Add attacks we show the results when LOVE is embedded in the memes. The results when CARESS is embedded are the same, as is expected, and hence not shown.

7.3 Countermeasures

In this section we present the results of our countermeasure approaches. We choose UNITER as the model since it showed slightly better performance in two out of the three datasets (observations from the other models show similar trends). We compare the adversarially attacked model performance in presence of the countermeasure techniques with those in the absence of these techniques. Precisely, we call the percentage difference in the macro-F1 score of the adversarially attacked model in absence and presence of countermeasures as x . Therefore a higher negative value of x indicates that a particular countermeasure is more effective in reducing the adverse effect of the attack on the model. We also measure how close the countermeasure brings the model performance to the original unattacked model. Thus, we compute the percentage difference in the macro-F1 score of the unattacked model and the attacked model with the countermeasure and call it y . The expected value of y is 0, i.e., after the application of the countermeasure the model performance should be equivalent to the performance of the unattacked model (which was also the main purpose of introducing

	MAMI			FBHM			HARMEME		
	CL-IND-0.5	VILLA-FT-GN	ENSEMBLE	CL-IND-0.5	VILLA-FT-GN	ENSEMBLE	CL-IND-0.5	VILLA-FT-GN	ENSEMBLE
Add	(4.61,-4.21)↓	(2.08,-1.68)↓	(1.43,-1.03)↓	(-1.08,0)↑	(-5.19,0)↑	(-6.39,0)↑	(-1.04,0)↑	(4.42,-4.1)↓	(3.3,-2.98)↓
Blur	(5.37,-6.07)↓	(1.47,-2.17)↓	(1.12,-1.82)↓	(-1.6,0)↑	(-3.14,0)↑	(-4.75,0)↑	(-0.34,-10.81)↑	(-2.91,-8.24)↑	(-3.5,-7.65)↑
SaltPepper-T	(7.6,-7.1)↓	(3.18,-2.68)↓	(2.05,-1.55)↓	(-2.25,0)↑	(-4.33,0)↑	(-4.35,0)↑	(0.53,-5.58)↓	(-2.03,-3.02)↑	(-2.56,-2.49)↑
SaltPepper-I-Low	(7.8,-9.5)↓	(5.7,-7.4)↓	(2.23,-3.93)↓	(-2,0)↑	(-2.44,0)↑	(-2.07,0)↑	(3.66,-9.71)↓	(1.47,-7.52)↓	(1.59,-7.64)↓
SaltPepper-I-High	(9.32,-13.15)↓	(3.73,-7.56)↓	(2.38,-6.21)↓	(-0.08,-1.89)↑	(-1.9,-0.07)↑	(-1.91,-0.06)↑	(-1.65,-16.13)↑	(-1.81,-15.97)↑	(-3.66,-14.12)↑
Newsprint	(11.12,-17.34)↓	(4.81,-11.03)↓	(4.93,-11.15)↓	(4.64,-7.86)↓	(0.26,-3.48)↓	(0.25,-3.47)↓	(2.56,-11.58)↓	(1.06,-10.08)↓	(2.05,-11.07)↓
Spread-Low	(3.4,-4)↓	(2.05,-2.65)↓	(1.02,-1.62)↓	(-2.32,0)↑	(-3.67,0)↑	(-4.16,0)↑	(1.58,-16.13)↓	(-1.83,-12.72)↑	(-1.79,-12.76)↑
Spread-High	(4.72,-9.27)↓	(-0.25,-4.3)↑	(0.48,-5.03)↓	(1.92,-4.55)↓	(-1.76,-0.87)↑	(-0.22,-2.41)↑	(3.32,-32.14)↓	(-5.03,-23.79)↑	(-5.08,-23.74)↑
Add+SaltPepper-I	(16.06,-20.77)↓	(6.21,-10.92)↓	(4.82,-9.53)↓	(3.36,-4.69)↓	(1.8,-3.13)↓	(0.99,-2.32)↓	(3.7,-11.47)↓	(2.41,-10.18)↓	(2.03,-9.8)↓
Average	(7.78,-10.75)↓	(3.22,-6.19)↓	(2.28,-5.25)↓	(0.06,-0.88)↓	(-2.27,0)↑	(-2.52,0)↑	(1.37,-11.31)↓	(-0.48,-9.46)↑	(-0.85,-9.09)↑

Table 4: Performance of best version of each countermeasure over the UNITER model. Each cell in the table is a (x, y) tuple as defined in the text. ↓ and ↑ respectively indicate whether a countermeasure worsens or improves the attacked model’s performance in terms of the macro-F1 score. We bold out the best countermeasures (only improvement) on each attack across the datasets.

the countermeasure in the first place). In general, a lower negative value of y should be better.

Key results and observations: Table 4 presents the results of this experiment. Each entry in the table shows the (x, y) values for a particular attack and a particular countermeasure. We show the performance of CL-IND-0.5 and VILLA-FT-GN since we find them to be the best among all the variants. Further, to unfold the best of these two countermeasures we also present an ENSEMBLE of both of them by considering average of their prediction probabilities. Overall, we make the following observations.

- For the FBHM and the HARMEME datasets, the countermeasures are effective for multiple of the attack strategies. For both these datasets, on average, the ENSEMBLE countermeasure performs the best followed by VILLA-FT-GN and CL-IND-0.5 countermeasures in that order.
- For the FBHM dataset the highest improvement in case of text based attacks is obtained for Add using the ENSEMBLE based countermeasure. The performance in presence of the countermeasure improves by 6.39% over the performance in absence of the countermeasure. In fact, in this case, the difference between the performance of the attacked model with countermeasure and the unattacked model is 0 as expected. The two other text based attacks – Blur and SaltPepper-T – also see large benefits in presence of the ENSEMBLE based countermeasure. Among the image based attacks, the VILLA-FT-GN and the ENSEMBLE based countermeasures are most effective for the SaltPepper-I-Low and the SaltPepper-I-High attacks respectively. For Spread-Low attack, ENSEMBLE based countermeasure is most effective and for Spread-High, VILLA-FT-GN is most effective.
- For the HARMEME dataset, the CL-IND-0.5 countermeasure does best for the Add attack. For the other two text based attacks the ENSEMBLE based countermeasure work the best. In the image based attacks, the ENSEMBLE based countermeasure once again works the best for the SaltPepper-I-High and the Spread-High attacks. VILLA-FT-GN is most effective for the Spread-Low attack.

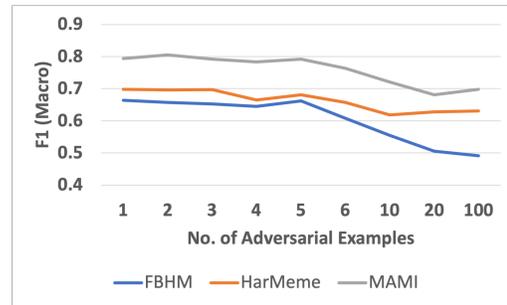


Figure 4: Average macro-F1 over all the attacks with increase in the adversarial examples used during fine-tuning process of VILLA model.

- For both the FBHM and HARMEME datasets none of the countermeasures work for certain attacks like Newsprint and Add+SaltPepper-I⁹ and is a scope for future improvement.
- Finally, none of the countermeasures seem to work for the MAMI dataset. On a deeper investigation we find that the resolution of the original images for this dataset is already very bad; the average *bit depth* for this dataset is 4.30 compared to 43.90 and 9.54 for the HARMEME and the FBHM datasets. This is possibly the reason why the introduction of the adversarial examples in the training phase does not benefit the models.

Impact of the number of adversarial examples: We further observe that the number of adversarial examples added in the training set has a positive impact only up to a certain extent. During the fine-tuning process of VILLA, for each training input sample we vary the number of adversarial examples that we add and estimate the performance after each addition. is used to generate the multiple adversarial examples. All other parameter settings remain the same. Figure 4 illustrates that the inclusion of up to five adversarial examples is tolerable by the system after which it steadily deteriorates.

⁹ x becomes positive in these cases which means the introduction of the countermeasure worsens the performance

8 CONCLUSION AND OUTLOOK

In this study, we systematically analyse vulnerability of multi-modal hate meme detection systems when they are attacked with partial model knowledge based adversaries. We audited widely popular visual linguistic models for a number of text and image based attacks and observed that there is a steep performance decrease in performance. Overall the image based attacks are found to be more severe. In order to counter such attacks we proposed two different methods – contrastive learning and adversarial training and found that an ensemble of these two methods work well for a large majority of attacks for two of the three datasets.

Future directions: A natural follow up of this work is to develop countermeasures that would work for very severe adversarial attacks like Newsprint and Add+SaltPepper-I. Another direction would be to identify countermeasure approaches for datasets where the original images are themselves of very low quality (e.g., the MAMI dataset in our study). Last but not the least, the set of our attack schemes are far from exhaustive and more variants need to be tried in future.

ACKNOWLEDGMENTS

This work was supported by CATALPA - Center of Advanced Technology for Assisted Learning and Predictive Analytics of the FernUniversität in Hagen, Germany.

REFERENCES

- [1] Tariq Habib Afridi, Aftab Alam, Muhammad Numan Khan, Jawad Khan, and Young-Koo Lee. 2020. A Multimodal Memes Classification: A Survey and Open Research Issues. <https://doi.org/10.48550/ARXIV.2009.08395>
- [2] Piush Aggarwal, Michelle Espranita Liman, Darina Gold, and Torsten Zesch. 2021. VL-BERT+: Detecting Protected Groups in Hateful Multimodal Memes. In *Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)*. Association for Computational Linguistics, Online, 207–214. <https://doi.org/10.18653/v1/2021.woaah-1.22>
- [3] Natalie Alkiviadou. 2019. Hate speech on social media networks: towards a regulatory framework? *Information & Communications Technology Law* 28, 1 (2019), 19–35. <https://doi.org/10.1080/13600834.2018.1494417>
- [4] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [5] Belhassen Bayar and Matthew C. Stamm. 2016. A Deep Learning Approach to Universal Image Manipulation Detection Using a New Convolutional Layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security (Vigo, Galicia, Spain) (IH&MMSec'16)*. Association for Computing Machinery, New York, NY, USA, 5–10. <https://doi.org/10.1145/2909827.2930786>
- [6] Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. 2013. Evasion Attacks against Machine Learning at Test Time. In *Machine Learning and Knowledge Discovery in Databases*, Hendrik Blockeel, Kristian Kersting, Siegfried Nijssen, and Filip Železný (Eds.), Springer Berlin Heidelberg, Berlin, Heidelberg, 387–402.
- [7] N. Carlini and D. Wagner. 2017. Towards Evaluating the Robustness of Neural Networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, Los Alamitos, CA, USA, 39–57. <https://doi.org/10.1109/SP.2017.49>
- [8] Olivier Chapelle, Jason Weston, Léon Bottou, and Vladimir Vapnik. 2000. Vicinal Risk Minimization. In *Advances in Neural Information Processing Systems*, T. Leen, T. Dietterich, and V. Tresp (Eds.), Vol. 13. MIT Press. <https://proceedings.neurips.cc/paper/2000/file/ba9a56ce0a9bfa26e8ed9e10b2cc8f46-Paper.pdf>
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv preprint arXiv:2002.05709* (2020).
- [10] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. 2020. Big Self-Supervised Models are Strong Semi-Supervised Learners. *arXiv preprint arXiv:2006.10029* (2020).
- [11] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *ECCV*.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [13] Samuel Dooley, Tom Goldstein, and John P. Dickerson. 2021. Robustness Disparities in Commercial Face Detection. *ArXiv abs/2108.12508* (2021).
- [14] Ivan Evtimov, Russel Howes, Brian Dolhansky, Hamed Firooz, and Cristian Canton Ferrer. 2020. Adversarial Evaluation of Multimodal Models under Realistic Gray Box Assumption. <https://doi.org/10.48550/ARXIV.2011.12902>
- [15] Edgar González Fernández, Ana Sandoval Orozco, Luis García Villalba, and Julio Hernandez-Castro. 2018. Digital Image Tamper Detection Technique Based on Spectrum Analysis of CFA Artifacts. *Sensors* 18, 9 (Aug. 2018), 2804. <https://doi.org/10.3390/s18092804>
- [16] Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. 2022. SemEval-2022 Task 5: Multimedia automatic misogyny identification. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*. Association for Computational Linguistics.
- [17] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. 2020. Large-Scale Adversarial Training for Vision-and-Language Representation Learning. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 6616–6628. <https://proceedings.neurips.cc/paper/2020/file/49562478de4c54fafd4ec46fd2b72de5-Paper.pdf>
- [18] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. 2020. Large-Scale Adversarial Training for Vision-and-Language Representation Learning. *arXiv:2006.06195 [cs.CV]*
- [19] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. 2020. Large-Scale Adversarial Training for Vision-and-Language Representation Learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (Vancouver, BC, Canada) (NIPS'20)*. Curran Associates Inc., Red Hook, NY, USA, Article 555, 13 pages.
- [20] Darina Gold, Piush Aggarwal, and Torsten Zesch. 2021. GerMemeHate: A Parallel Dataset of German Hateful Memes Translated from English. <https://www.inf.uni-hamburg.de/en/inst/ab/lt/publications/2021-alacamyimam-konvens-mmhs21.pdf#page=9>
- [21] Thomas Gottron. 2008. Content Code Blurring: A New Approach to Content Extraction. In *2008 19th International Workshop on Database and Expert Systems Applications*. 29–33. <https://doi.org/10.1109/DEXA.2008.43>
- [22] Tommi Gröndahl, Luca Pajola, Mika Juuti, Mauro Conti, and N. Asokan. 2018. All You Need is "Love": Evading Hate-speech Detection. <https://doi.org/10.48550/ARXIV.1808.09115>
- [23] Amos Guiora and Elizabeth A. Park. 2017. Hate Speech on Social Media. *Philosophia* 45, 3 (July 2017), 957–971. <https://doi.org/10.1007/s11406-017-9858-4>
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. <https://doi.org/10.48550/ARXIV.1512.03385>
- [25] Ming Shan Hee, Roy Ka-Wei Lee, and Wen-Haw Chong. 2022. On Explaining Multimodal Hateful Meme Detection Models. <https://doi.org/10.48550/ARXIV.2204.01734>
- [26] Siddharth Jaiswal, Karthikeya Duggirala, Abhisek Dash, and Animesh Mukherjee. 2022. Two-Face: Adversarial Audit of Commercial Face Recognition Systems. *Proceedings of the International AAAI Conference on Web and Social Media* 16, 1 (May 2022), 381–392. <https://doi.org/10.1609/icwsm.v16i1.19300>
- [27] Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. 2020. SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 2177–2190. <https://doi.org/10.18653/v1/2020.acl-main.197>
- [28] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes. <https://doi.org/10.48550/ARXIV.2005.04790>
- [29] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. <https://doi.org/10.48550/ARXIV.1412.6980>
- [30] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the Carbon Emissions of Machine Learning. *arXiv preprint arXiv:1910.09700* (2019).

- [31] Vijaysinh Lendave. 2021. A guide to different types of noises and image denoising methods. <https://analyticsindiamag.com/a-guide-to-different-types-of-noises-and-image-denoising-methods/>
- [32] Linjie Li, Zhe Gan, and Jingjing Liu. 2021. A Closer Look at the Robustness of Vision-and-Language Pre-trained Models. arXiv:2012.08673 [cs.CV]
- [33] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. VisualBERT: A Simple and Performant Baseline for Vision and Language. <https://doi.org/10.48550/ARXIV.1908.03557>
- [34] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*. Springer, 121–137.
- [35] Phillip Lippe, Nithin Holla, Shantanu Chandra, Santhosh Rajamanickam, Georgios Antoniou, Ekaterina Shutova, and Helen Yannakoudakis. 2020. A Multimodal Framework for the Detection of Hateful Memes. (2020). <https://doi.org/10.48550/ARXIV.2012.12871>
- [36] Xiaofeng Liu, Yang Zou, Lingsheng Kong, Zhihui Diao, Junliang Yan, Jun Wang, Site Li, Ping Jia, and Jane You. 2018. Data Augmentation via Latent Space Interpolation for Image Classification. In *2018 24th International Conference on Pattern Recognition (ICPR)*. 728–733. <https://doi.org/10.1109/ICPR.2018.8545506>
- [37] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. <https://doi.org/10.48550/ARXIV.1907.11692>
- [38] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards Deep Learning Models Resistant to Adversarial Attacks. <https://doi.org/10.48550/ARXIV.1706.06083>
- [39] Francesco Marra, Diego Gragnaniello, Davide Cozzolino, and Luisa Verdoliva. 2018. Detection of GAN-Generated Fake Images over Social Networks. In *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. 384–389. <https://doi.org/10.1109/MIPR.2018.00084>
- [40] Martina Miliani, Giulia Giorgi, Ilir Rama, Guido Anselmi, and Gianluca E. Lebani. 2020. DANKMEMES @ EVALITA 2020: The Memeing of Life: Memes, Multimodality and Politics. In *EVALITA Evaluation of NLP and Speech Tools for Italian - December 17th, 2020*. Accademia University Press, 275–283. <https://doi.org/10.4000/books.aaccademia.7330>
- [41] Niklas Muennighoff. 2020. Vili: State-of-the-art Visio-Linguistic Models applied to Hateful Memes. <https://doi.org/10.48550/ARXIV.2012.07788>
- [42] Lin Pan, Chung-Wei Hang, Avirup Sil, and Saloni Potdar. 2021. Improved Text Classification via Contrastive Adversarial Training. <https://doi.org/10.48550/ARXIV.2107.10137>
- [43] Tianyu Pang, Xiao Yang, Yinpeng Dong, Kun Xu, Jun Zhu, and Hang Su. 2020. Boosting Adversarial Training with Hypersphere Embedding. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 7779–7792. <https://proceedings.neurips.cc/paper/2020/file/5898d8095428ee310bf7fa3da1864ff7-Paper.pdf>
- [44] Zoe Papakipos and Joanna Bitton. 2022. AugLy: Data Augmentations for Robustness. arXiv:2201.06494 [cs.AI]
- [45] Gabriel Peyré and Marco Cuturi. 2018. Computational Optimal Transport. (2018). <https://doi.org/10.48550/ARXIV.1803.00567>
- [46] Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021. Detecting Harmful Memes and Their Targets. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 2783–2796. <https://doi.org/10.18653/v1/2021.findings-acl.246>
- [47] Yao Qiu, Jinchao Zhang, and Jie Zhou. 2021. Improving Gradient-based Adversarial Training for Text Classification by Contrastive Learning and Auto-Encoder. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 1698–1707. <https://doi.org/10.18653/v1/2021.findings-acl.148>
- [48] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A. Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann. 2021. Data Augmentation Can Improve Robustness. <https://doi.org/10.48550/ARXIV.2111.05328>
- [49] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A. Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann. 2021. Fixing Data Augmentation to Improve Adversarial Robustness. <https://doi.org/10.48550/ARXIV.2103.01946>
- [50] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (Eds.), Vol. 28. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2015/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf>
- [51] Benet Oriol Sabat, Cristian Canton Ferrer, and Xavier Giro i Nieto. 2019. Hate Speech in Pixels: Detection of Offensive Memes towards Automatic Moderation. arXiv:1910.02334 [cs.MM]
- [52] Chhavi Sharma, Deepesh Bhageria, William Scott, Srinivas PYKL, Amitava Das, Tanmoy Chakraborty, Viswanath Pulabaigari, and Björn Gambäck. 2020. SemEval-2020 Task 8: Memotion Analysis- the Visuo-Lingual Metaphor!. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*. International Committee for Computational Linguistics, Barcelona (online), 759–773. <https://doi.org/10.18653/v1/2020.semeval-1.99>
- [53] Patrice Y. Simard, Yann A. LeCun, John S. Denker, and Bernard Victorri. 2012. *Transformation Invariance in Pattern Recognition - Tangent Distance and Tangent Propagation*. Springer Berlin Heidelberg, Berlin, Heidelberg, 235–269. https://doi.org/10.1007/978-3-642-35289-8_17
- [54] Shardul Suryawanshi, Bharathi Raja Chakravarthi, Mihael Arcan, and Paul Buitelaar. 2020. Multimodal Meme Dataset (MultiOFF) for Identifying Offensive Content in Image and Text. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*. European Language Resources Association (ELRA), Marseille, France, 32–41. <https://aclanthology.org/2020.trac-1.6>
- [55] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. <https://doi.org/10.48550/ARXIV.1312.6199>
- [56] Hao Tan and Mohit Bansal. 2019. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. <https://doi.org/10.48550/ARXIV.1908.07490>
- [57] Önsen Toygar, Felix O Babalola, and Yiltan Bitirim. 2020. FYO: a novel multimodal vein database with palmar, dorsal and wrist biometrics. *IEEE Access* 8 (2020), 82461–82470.
- [58] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. 2017. Ensemble Adversarial Training: Attacks and Defenses. <https://doi.org/10.48550/ARXIV.1705.07204>
- [59] Riza Velioglu and Jewgeni Rose. 2020. Detecting Hate Speech in Memes Using Multimodal Deep Learning Approaches: Prize-winning solution to Hateful Memes Challenge. <https://doi.org/10.48550/ARXIV.2012.12975>
- [60] Nishant Vishwamitra, Hongxin Hu, Ziming Zhao, Long Cheng, and Feng Luo. 2021. Understanding and Measuring Robustness of Multimodal Learning. <https://doi.org/10.48550/ARXIV.2112.12792>
- [61] Nishant Vishwamitra, Hongxin Hu, Ziming Zhao, Long Cheng, and Feng Luo. 2021. Understanding and Measuring Robustness of Multimodal Learning. arXiv:2112.12792 [cs.LG]
- [62] Guoqing Wang, Chuanxin Lan, Hu Han, Shiguang Shan, and Xilin Chen. 2019. Multi-modal face presentation attack detection via spatial and channel attentions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.
- [63] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. 2019. Detectron2. <https://github.com/facebookresearch/detectron2>.
- [64] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. <https://doi.org/10.48550/ARXIV.1609.08144>
- [65] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He. 2019. Feature Denoising for Improving Adversarial Robustness. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [66] Karren Yang, Wan-Yi Lin, Manash Barman, Filipe Condessa, and Zico Kolter. 2021. Defending Multimodal Fusion Models Against Single-Source Adversaries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 3340–3349.
- [67] Savvas Zannettou, Barry Bradlyn, Emiliano De Cristofaro, Haewoon Kwak, Michael Sirivianos, Gianluca Stringini, and Jeremy Blackburn. 2018. What is Gab: A Bastion of Free Speech or an Alt-Right Echo Chamber. In *Companion Proceedings of the The Web Conference 2018 (Lyon, France) (WWW ’18)*. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 1007–1014. <https://doi.org/10.1145/3184558.3191531>
- [68] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. 2019. Theoretically Principled Trade-off between Robustness and Accuracy. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), PMLR, 7472–7482. <https://proceedings.mlr.press/v97/zhang19p.html>
- [69] Lin Zhao, Changsheng Chen, and Jiwu Huang. 2021. Deep Learning-Based Forgery Attack on Document Images. *IEEE Transactions on Image Processing* 30 (2021), 7964–7979. <https://doi.org/10.1109/TIP.2021.3112048>
- [70] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. 2020. FreeLB: Enhanced Adversarial Training for Natural Language Understanding. In *International Conference on Learning Representations*. <https://openreview>.

net/forum?id=BygzbyHFvB

A ETHICAL CONSIDERATION

A.1 Unintentional bias

Any biases found in our study are unintentional, and we do not intend to do harm to any group or individual. We note that determining whether a meme is harmful can be subjective, and thus it is inevitable that there would be biases in model training process and its analysis.

A.2 Code misuse

We intend to publish our code base in order to replicate our study for further research. However, it can also be exploited to generate

attacks in order to confuse the algorithms deployed by social media. Therefore, human moderation in addition to algorithmic detection is needed in order to ensure that this does not occur.

A.3 Carbon emission

We carried out most of our experiments on GPUs to analyse their robustness and to generate counter-measure models. Our experiments were conducted using a private infrastructure, which has a carbon efficiency of 0.432 kgCO₂eq/kWh. A cumulative of 300 hours of computation was performed on hardware of type RTX 2080 Ti (TDP of 250W). Total emissions are estimated to be 32.4 kgCO₂eq of which 0% was directly offset. Estimations were conducted using the Machine Learning Impact calculator presented in [30].