# A Multimodal Framework for the Detection of Hateful Memes

**Phillip Lippe**[*]
QUVA
University of Amsterdam
p.lippe@uva.nl

**Nithin Holla**[*]
University of Amsterdam
nithin.holla7@gmail.com

**Shantanu Chandra**[*]
ZS
University of Amsterdam
shantanu.chandra@zs.com

**Santhosh Rajamanickam**
Slimmer AI
rajamanickamsanthosh@gmail.com

**Georgios Antoniou**
King's College London
georgios.antoniou@kcl.ac.uk

**Ekaterina Shutova**
University of Amsterdam
e.shutova@uva.nl

**Helen Yannakoudakis**
King's College London
helen.yannakoudakis@kcl.ac.uk

## Abstract

An increasingly common expression of online hate speech is multimodal in nature and comes in the form of *memes*. Designing systems to automatically detect hateful content is of paramount importance if we are to mitigate its undesirable effects on the society at large. The detection of multimodal hate speech is an intrinsically difficult and open problem: memes convey a message using both images and text and, hence, require multimodal reasoning and joint visual and language understanding. In this work, we seek to advance this line of research and develop a multimodal framework for the detection of hateful memes. We improve the performance of existing multimodal approaches beyond simple fine-tuning and, among others, show the effectiveness of upsampling of contrastive examples to encourage multimodality, and ensemble learning based on cross-validation to improve robustness. We furthermore analyze model misclassifications and discuss a number of hypothesis-driven augmentations and their effects on performance, which we hope shall inform future research in the field. Our best approach comprises an ensemble of UNITER-based [3] models and achieves an AUROC score of $80.53$, placing us 4th on Phase 2 of the 2020 *Hateful Memes Challenge* organized by Facebook.

## 1 Introduction

Online abuse is an important societal problem of our time and one that is highly correlated with the rise of social media platforms such as Twitter and Facebook. A large number of Internet users have either encountered or witnessed some form of abusive behaviour online. For instance, it has been reported that $41\%$ of American adults have experienced online harassment Duggan [7]. Hate speech, a prevalent form of online abuse, has seen a dramatic rise in recent years [18]. It can be defined as any form of communication that attacks or uses discriminatory language with reference to a person or a group based on their race, gender, religion, etc. *Memes* – that have recently emerged as popular engagement tools and which, in their usual form, are image macros shared through social media platforms mainly for amusement – are also being increasingly used to spread hate and/or

---

[*]The authors contributed equally to this work.

instigate social unrest, and therefore seem to be a new form of expression of hate speech on online platforms. *Hateful memes* often target certain communities or individuals based for example on religion, gender, race, or physical attributes, by portraying them in a derogatory manner and/or by reinforcing stereotypes. Such memes promote racism [35] and sexism [6] among other forms of hate speech, threatening social peace and leading to damage at both the individual and societal level [2].

In this light, developing systems that can automatically detect hateful memes (and hateful content in general) is of paramount importance if we are to mitigate their undesirable effects. However, the detection of multimodal hate speech is an intrinsically difficult and open problem within the joint visual and language (V+L) understanding domain as it requires a holistic understanding of content, where reasoning about image and text is simultaneous. As the example in Figure 1 shows, it is not sufficient to rely on the text or the image modality individually for correct interpretation of content; rather both modalities should be jointly processed to infer the correct meaning of the meme. Due to the multimodal nature of the problem involving an interplay between the image and the text, existing and state-of-the-art multimodal systems perform rather poorly on the detection of hateful memes [15]. This highlights the need to further advance multimodal reasoning and understanding for systems to be able to more accurately flag hateful content online.

With that in mind, Facebook launched the *Hateful Memes Challenge* [15] as part of the NeurIPS 2020 competition track to encourage further research into multimodal reasoning and the design of systems that can detect hateful memes. Specifically, the task is formulated as a binary classification problem, where a meme can belong to one of two classes – *hateful* or *not hateful*. The challenge introduces a new Hateful Memes (HM) dataset [15], consisting of over $10,000$ memes. To enforce multimodality and reduce accidental biases in system classifications, the dataset includes non hateful, *benign confounders* (i.e., *contrastive* or *counterfactual* examples) for a subset of hateful memes. This makes the task particularly challenging for unimodal systems that only leverage signals from either of the two modalities and do not combine information from both the text and the image. This is corroborated by the performance



Figure 1: Example (non-hateful) multimodal meme from the HM dataset [15]. The image is a compilation of assets, including ©Getty Images.

of the baselines implemented by the organisers of the challenge [15], which clearly shows a substantial difference in performance between unimodal and multimodal systems; albeit the latter still performing rather poorly, and particularly compared to human performance, indicating that a more refined understanding of multiple signals is necessary.

From their thorough baseline experiments, the challenge organisers concluded that early-fusion multimodal models considerably outperform late-fusion architectures for the task. Building on their work, we experiment with a number of early-fusion multimodal approaches, namely LXMERT [34], UNITER [3], and Oscar [20], for the task of hateful memes detection. We propose various methods for improving model performance beyond simple fine-tuning and, among others, show the effectiveness of upsampling confounders to encourage multimodality, and ensemble learning based on cross-validation to improve model robustness. Our best system is an ensemble method based on UNITER that achieves an AUROC score of $80.53$, ranking our team, *Kingsterdam*, 4th on the final phase (Phase 2) of the Hateful Memes Challenge. We conclude with an error analysis of model misclassifications and discuss a number of hypothesis-driven model augmentations based on multi-task learning and their effects on performance, which we hope shall inform future research in the field. To facilitate further research, we make our code publicly available.[2]

## 2 Related work

### 2.1 Multimodal learning

Multimodal representation learning has recently gained traction due to the poor performance of existing (unimodal) models on multimodal tasks such as Visual Question Answering [1, 12] and Visual Reasoning [33]. These tasks involve V+L understanding and identifying the synergy be-

---

[2]`https://github.com/Nithin-Holla/meme_challenge`

tween the two modalities. Most existing multimodal systems adopt either a late-fusion (LF) or an early-fusion (EF) approach to process the two modalities. Late-fusion methods [14, 15] typically utilize unimodal models to process the two signals independently and then combine their features (usually via simple concatenation) before the final classification layer. Early-fusion methods such as MMBT [14], VisualBERT [19], and ViLBERT [23], on the other hand, employ more complex approaches to process the two modalities jointly within the model architecture. UNITER (UNiversal Image-TExt Representation) [3], LXMERT (Learning Cross-Modality Encoder Representations from Transformers) [34] and Oscar (Object-Semantics Aligned Pre-training) [20] are some of the recent and popular early-fusion multimodal representation learning models which are pretrained on various V, L and V+L tasks such as visual question answering and image captioning, and which we use in our work (discussed in detail in Section 4). Oscar, the more recent of the three, outperforms LXMERT and UNITER on the VQA [1], GQA [12], and NLVR2 [33] downstream tasks and achieves state-of-the-art performance on VQA and NLVR2, among others.

## 2.2 Multimodal hate speech detection

Existing work on hate speech detection has largely relied on text-based features. Mishra et al. [25] discuss emerging trends, resources and challenges, as well as outline the various approaches used in the domain of online abusive language detection. However, there is comparatively little work in the vision or multimodal domain, something which can be attributed to the scarcity of annotated datasets. Gomez et al. [9] introduce MMHS150K, a multimodal dataset of tweets consisting of both image and text and which are manually annotated for hate speech. The authors develop and evaluate three multimodal models – Feature Concatenation Model (FCM), Spatial Concatenation Model (SCM) and Textual Kernels Model (TKM) – and conclude that they fail to outperform unimodal counterparts. Hosseinmardi et al. [11] discuss the problem of cyberbullying and cyberaggression in Instagram posts and comments and evaluate the performance of a Naive Bayes classifier and a linear SVM classifier using a range of features, such as word n-grams and image categories, as well as metadata such as the number of followers and likes.

With regards to the *Hateful Memes Challenge*, the organizers provide several baseline models [15] and evaluate their performance using accuracy and AUROC. The baselines comprise unimodal and multimodal systems which are pretrained using diverse methods. Their experiments show that multimodal approaches outperform all the unimodal systems, reiterating the fact that this task requires a holistic understanding of the meme by processing the image and text signals jointly.

The multimodal approaches include several architectures based on BERT [5] that combine image and text features to get the final prediction. The best-performing multimodal baseline models are ViLBERT [23] and VisualBERT [19] which are pretrained on the CC (Conceptual Captions) [31] and the COCO (Common Objects in Context) [21] datasets respectively. ViLBERT extends the BERT architecture to a multimodal, two-stream model and incorporates separate transformers for the vision and language domain that interact through co-attentional transformer layers to learn joint representations of images and text. VisualBERT is a single-stream architecture which uses a self-attention mechanism within a layer of the transformer which includes both vision and language inputs.

## 3 Dataset

The challenge uses the Hateful Memes (HM) dataset [15] compiled by Facebook AI. It includes the memes (image with text) as well as the meme text separately to facilitate easier processing. The dataset consists of over 10k memes labeled as *hateful* or *non-hateful* using the definition of hatefulness presented by Kiela et al. [15] for this task – "a direct or indirect attack on people based on characteristics, including ethnicity, race, nationality, religion, caste, sex, etc." (see [15] for the full definition). The dataset also consists of an additional 2k unlabeled memes that form the test sets of the challenge and which were released in two phases of the competition: Phase 1 ("seen" test set) and Phase 2 (unseen and final test set used to determine system rankings).

The dataset is designed such that multimodal approaches are essential to perform well on the task. Specifically, for a set of hateful memes, *benign confounders* are devised and included in the dataset. These are defined as the result of minimal transformations to either the image or the text of the meme such that the corresponding label flips from hateful to non-hateful, or vice versa. Hence, there are

Table 1: Distribution of the different types of memes in the Hateful Memes dataset.

| Split | Multimodal hate | Unimodal hate | Benign confounders | Random benign | Dynamic adversarial benign confounders | Total |
|---|---|---|---|---|---|---|
| Train | 1300 | 1750 | 3200 | 2250 | – | 8500 |
| Dev-seen | 200 | 50 | 200 | 50 | – | 500 |
| Test-seen | 400 | 100 | 400 | 100 | – | 1000 |
| Dev-unseen | 200 | – | 200 | – | 140 | 540 |
| Test-unseen | 729 | – | 597 | – | 674 | 2000 |

two types of confounders – (a) *image confounders*, where the pair of memes have the same image but different text, and (b) *text confounders*, where the pair of memes have the same text but different images. Overall, the dataset comprises five different types of memes: *multimodal hate*, *unimodal hate*, *benign image* and *benign text* confounders, and finally *random non-hateful examples*. We refer the reader to Kiela et al. [15] for more details on the dataset and its construction.

There is a total of $12, 140$ memes, split into train, development (validation) and test sets. The training set contains $8, 500$ examples, with $36\%$ hateful and $64\%$ non-hateful memes. The development and test sets are composed of two sets each corresponding to Phase 1 (*dev-seen*, *test-seen*) and Phase 2 (*dev-unseen*, *test-unseen*) of the competition. The *dev-seen* and *test-seen* sets are class-balanced and contain $500$ and $1, 000$ examples respectively. In these, the distributions of the five different types of memes are as follows: $40\%$ multimodal hate, $10\%$ unimodal hate, $20\%$ benign text, $20\%$ benign image confounders and the remaining $10\%$ are random non-hateful. The *dev-unseen* set includes multimodal hateful and benign confounders from *dev-seen*, as well as a new set of $140$ dynamic adversarial benign confounders that are chosen such that the VisualBERT model fails to classify them correctly. The *test-unseen* set comprises several new multimodal hate and benign confounders as well a new set of dynamic adversarial benign confounders based on VisualBERT, resulting to a total of $2, 000$ examples. The composition of *test-unseen* makes this a challenging set and an appropriate test of multimodality. We provide the details of the distribution of the different types of memes on the various splits of the dataset in Table 1.

## 4   Models

Early-fusion architectures were shown to outperform the late-fusion ones for the task [15]; therefore, we focus on three early-fusion pretrained models, namely LXMERT [34], UNITER [3] and Oscar [20].

### 4.1   LXMERT

LXMERT is a large-scale transformer model which first processes the image and text by two independent, unimodal encoders. Another consecutive encoder combines the two unimodal representations via cross-attention modules. The model is pretrained on images from COCO and Visual Genome [17] as well as the image question answering datasets – VQA v2.0 [10], GQA [12], and VG-QA [38]. Specifically, it is pretrained on the tasks of masked language modeling, masked object prediction, cross modality matching, and image question answering. LXMERT has been shown to outperform ViLBERT and VisualBERT on downstream tasks such as visual reasoning tasks.

### 4.2   UNITER

UNITER is an early-fusion transformer model which utilizes self-attention on the joint image and textual input. It encodes visual features from bounding boxes using an image encoder and encodes word tokens using a text encoder into a common embedding space. The image input consists of features from Faster R-CNN [29] and 7-dimensional location features consisting of normalized top / left / bottom / right coordinates, width, height, and area for each of the bounding boxes. The text input consists of word embeddings as well as position embeddings. It is followed by several layers of self-attention.

The model is pretrained on four tasks, namely masked language/region modeling conditioned on image/text (MLM and MRM), image-text matching (ITM) and word-region alignment (WRA). The large-scale pretraining is performed over four V+L datasets – COCO, Visual Genome, CC, and SBU Captions [26]. UNITER has been shown to outperform baselines including LXMERT on six V+L tasks such as masked language modeling (MLM), ITM and MRFR (Masked Region Feature Regression) among others, across nine datasets.

### 4.3 Oscar

Oscar is another multimodal transformer model whose input consists of triples containing the word sequence, the set of object tags detected in the image, and the set of image region features. The object tags act as *anchor points* that make learning of semantic alignments between images and texts easier. Faster R-CNN is used to produce bounding box features as well as the set of object tags. Additionally, a 6-dimensional position vector is included for each of the bounding boxes.

The large-scale pretraining consists of two objectives – a masked token loss and a contrastive loss. Pretraining is performed on V+L datasets including COCO, CC, SBU, Flickr30k [36], VQA [1], GQA, and VG-QA. Oscar was shown to obtain state-of-the-art results on six downstream tasks such as VQA, VCR and NLVR, among others.

## 5 Approach

In this section, we present the various design aspects of our proposed framework for this task. We begin by describing how we extract image features for this dataset. We then describe how simple learning strategies can help us perform better on multimodal hate and benign confounders. We follow that with detailing the importance of creating an ensemble of models to help the model generalize better, and how this ensemble can be effectively optimized using an evolutionary algorithm (EA) for optimal weighting of model predictions. Finally, we discuss how supplementing the model with additional information from fine-grained object-detection classes may help the model identify the target groups in the image and the meme text.

### 5.1 Image feature extraction and base model selection

All our multimodal architectures – LXMERT, Oscar and UNITER (Section 4) – take both tokenized text as well as processed image features as input. To that effect, we first extract the image features for the HM dataset using a pretrained Faster R-CNN [29] backbone. For this, we made sure that we use the corresponding pretrained model checkpoint[3] as described in the original papers of these architectures rather than the features provided in the MMF library [32] for the dataset and released as part of the competition. This was done to ensure that the input features of the meme images are in the same latent space as those that these multimodal pretrained models were trained on, to avoid erroneous behavior. The tokenized text features were obtained using the standard pretrained BERT tokenizer [5] for the models, inline with their pretraining step.

Using the text and image features as input, we first fine-tune each of the models on the HM dataset (both the base and large versions of the models) for the supervised task of hateful memes detection using a binary classification objective. Specifically, we optimize for the binary cross-entropy (BCE) loss which is computed as:

$$\mathcal{L}_{BCE}(\boldsymbol{\theta}) = \sum_{i=1}^{N} -\left[y_i \log f_{\boldsymbol{\theta}}(x_i) + (1 - y_i) \log\left(1 - f_{\boldsymbol{\theta}}(x_i)\right)\right] \tag{1}$$

where $f_{\boldsymbol{\theta}}$ is the output of the model, $x_i$ is the input and $y_i$ is the gold label.

We find that UNITER outperforms LXMERT and Oscar by a large margin for both their base and large versions (see results for base models in Table 2), which could be attributed to its diverse set of pretraining tasks. We also note that the large versions of these architectures lead to poor generalization performance due to more severe overfitting on the training data as a result of their larger set of parameters. Given these observations, we proceed with UNITER-base (hereafter referred to as UNITER) as the base model of our framework and use it for all subsequent experiments.

---

[3]`https://github.com/MILVLG/bottom-up-attention.pytorch`

Using the pretrained UNITER (base), we also experimented with an additional pre-training / domain-adaptation step on the HM dataset in an attempt to align the model's weights to the new latent space of the memes domain. Specifically, we took the pretrained UNITER-base model and fine-tuned it further on the MLM, ITM and MRFR (within MRM) pretraining tasks using the HM dataset. This "warm-up" pretraining phase was then followed by the supervised fine-tuning step on the binary task of hateful memes detection. However, this did not yield any performance improvements and we therefore only apply supervised fine-tuning.

## 5.2 Confounder upsampling and loss re-weighting

A key characteristic of the HM dataset is the inclusion of benign confounders to counter the possibility of models exploiting unimodal priors rather than learning to reason multimodally. Thus, these sets of memes form an important source of truly multimodal instances that the model can directly leverage for effective learning.

During our model-benchmarking experiments (Section 5.1), we noticed that performance was quite poor on benign confounders and specifically text confounders. The models did not effectively exploit these instances during training such that they can multimodally infer their true, underlying meaning. As a simple strategy to alleviate this problem, we upsample the confounders in the batches during the supervised fine-tuning step of UNITER. As the model performed quite well on the image confounders, we only upsample text confounders. Intuitively, this helps the model to focus on input features of both modalities and subsequently improve its performance. We refer to the confounder upsampling approach using the abbreviation CFU.

Another important characteristic of the data that directly affects learning is the gold label distribution. The HM training set consists of 36% hateful and 64% non-hateful memes. To improve the detection of hateful memes, we employ a loss re-weighting strategy during training and weigh the loss of the hateful class higher, subsequently affecting the updates of the model parameters. Thus our new loss function is defined as follows:

$$\mathcal{L}_{HW}(\boldsymbol{\theta}) = \sum_{i=1}^{N} - \left[ \alpha_{pos} \cdot y_i \log f_{\boldsymbol{\theta}}(x_i) + \alpha_{neg} \cdot (1 - y_i) \log (1 - f_{\boldsymbol{\theta}}(x_i)) \right] \qquad (2)$$

where $\alpha_{pos}$ and $\alpha_{neg}$ are the weights for the hateful and not-hateful classes respectively such that $\alpha_{pos} > \alpha_{neg}$ and $\alpha_{pos} + \alpha_{neg} = 1$. We refer to this loss re-weighting approach on the hateful class using the abbreviation HW.

## 5.3 Cross-validation ensemble

Multimodal datasets are typically much smaller in size than other visual datasets used to train deep image classification models. Training on different subsets of the rather small training set of the HM dataset can lead to considerable variation in model predictions. To stabilize the predictions and tackle overfitting, we utilize an ensemble of UNITER models where we combine them using a weighted average of their predictions. The weight of each model is determined by an evolutionary algorithm (EA) that optimizes the weights with respect to the AUROC score of the ensemble on the development set. We denote models trained using cross-validation folds derived from the HM training set by the abbreviation CV.

The development set contains a good distribution of true multimodal examples which can be valuable during training. To fully utilize the data, exploit as many multimodal examples as possible and test generalization performance of fine-tuned parameters, we employ another cross-validation strategy such that we can also learn from the development set as well as use it for ensemble optimization. Specifically, we split the HM training set into CV folds (implementation details in Section 6) and, in each training fold we include half of the development set (dev-seen), while the other half is added to the test fold. While splitting the development set, we ensure that text confounder pairs (with different labels) remain together. This means that examples that form a confounder pair are never split across the CV training and test sets. The splitting of the development set is otherwise dynamic, i.e., random subsets of the development set are selected for inclusion in the CV folds. EA optimization of the ensemble weights is now performed on the augmented CV test folds. A visual representation of this cross-validation ensemble optimization process is presented in Figure 2. We denote models/ensembles trained using the augmented cross-validation folds with the subscript FINAL.
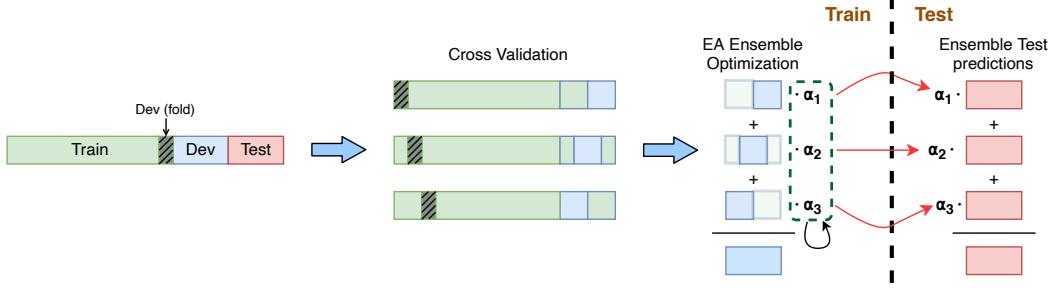
Figure 2: A representation of our cross-validation ensemble optimization process. The dev set (dev-seen; blue) is split into two parts (such that confounder pairs remain together) which are used to augment the CV training and test folds (green). The augmented CV test folds (green+blue) are used for the EA optimization on the right. The final ensemble weights $\alpha_1, ..., \alpha_M$, where $M$ is the number of CV folds, are used to combine the CV model predictions and evaluate performance on the final, unseen test set.

## 5.4 Training with margin ranking loss

In order to explicitly encourage learning from confounders, we also experiment with training using margin ranking loss as a means of contrastive learning. For every meme $x_i$ with label $y_i$, we sample another meme $\tilde{x}_i$ with label $\tilde{y}_i$ such that $\tilde{x}_i$ is the text confounder of $x_i$ if it exists, or a randomly sampled meme from the training set such that $y_i \neq \tilde{y}_i$. Training is performed on input pairs $\{(x_i, y_i, \tilde{x}_i, \tilde{y}_i)\}_{i=1}^N$ where $N$ is the total number of memes in the training set. The objective is to predict a higher probability score for the meme that is labeled as hateful. The new training loss is then calculated as the weighted sum of binary cross-entropy loss defined earlier and margin ranking (MR) loss:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}_{HW}(\boldsymbol{\theta}) + \gamma\mathcal{L}_{MR}(\boldsymbol{\theta}) \tag{3}$$

where $\boldsymbol{\theta}$ is the set of model parameters and $\gamma$ is a hyperparameter that controls the influence of the MR loss. The BCE loss $\mathcal{L}_{\mathcal{HW}}$ is computed as described in Section 5.2, while the MR loss is computed using the input pair as follows:

$$\mathcal{L}_{MR}(\boldsymbol{\theta}) = \sum_{i=1}^N \max\left[0, -y_{RANK}\left(f_{\boldsymbol{\theta}}(x_i) - f_{\boldsymbol{\theta}}(\tilde{x}_i)\right) + m\right] \tag{4}$$

where $m$ is the margin hyperparameter and $y_{RANK} = 1$ if $y_i = 1$ (hateful) and $\tilde{y}i = 0$ (not hateful), and $y_{RANK} = -1$ otherwise. At testing time, evaluation is performed in the standard setting, i.e., without any instance pairing. We refer to this training approach using the abbreviation MRL.

## 5.5 Fine-grained YOLO9000 object tags

During qualitative analyses of our models, we find misclassification errors on memes promoting hate towards certain communities or individuals based on religion, gender, and race, among others. We hypothesize that augmenting the model with information that would help identify the target group may lead to more accurate predictions. Specifically, we develop a variant of our model that uses YOLO9000 predictions [28]. YOLO9000 is a well-known image object detection model that can detect objects from a pre-defined set of 9000 fine-grained classes. Since not all of these classes are relevant to the HM dataset, we instead use a subset of 97 classes, such as "Nigerian", "Muslimah", "revolver", and "amputee" (see Section A.1 for more details). For every meme, we take the output classes predicted from the pretrained YOLO9000 model (multi-label classification), filter them to use only our selected set, and provide them as additional input to the model. Concretely, the input now consists of the meme text, YOLO9000 object tags, and bounding box features. We refer to this model using the subscript YOLO.

## 6 Experiments and Results

In this section, we provide our implementation details and the hyperparameters of our best-performing models. We then present and discuss the performance of a number of model variants.

7

Table 2: AUROC scores of our models on the development set (dev-seen) as well as the Phase 1 and 2 test sets of the challenge.

| Model | AUROC | | |
| --- | --- | --- | --- |
| | Development | Phase 1 | Phase 2 |
| ViLBERT CC | 70.07 | 70.03 | – |
| VisualBERT COCO | 73.97 | 71.41 | – |
| UNITER | 78.04 | 74.73 | – |
| LXMERT | 72.33 | – | – |
| Oscar | 72.00 | – | – |
| UNITER$_{\text{CV10}}$ | 79.81 | – | – |
| UNITER$_{\text{CV10 + CFU}}$ | 79.64 | – | – |
| UNITER$_{\text{CV10 + CFU + HW}}$ | 80.01 | 78.60 | – |
| UNITER$_{\text{CV15 + CFU + HW}}$ | 80.65 | 79.06 | – |
| UNITER$_{\text{CV30 + CFU + HW}}$ | 81.36 | 78.98 | – |
| UNITER$_{\text{CV15 + CFU + HW + MRL}}$ | 80.44 | 78.14 | – |
| UNITER$_{\text{CV15 + CFU + HW + YOLO}}$ | 80.67 | 78.21 | – |
| UNITER$_{\text{ENSEMBLE 1}}$ | 81.71 | 79.13 | 80.33 |
| UNITER$_{\text{ENSEMBLE 2}}$ | 81.76 | 79.10 | 80.40 |
| UNITER$_{\text{FINAL}}$ | 77.39 | 79.07 | **80.53** |

**Hyperparameters** For all our experiments, we use UNITER-base as our main model. We work with a batch size of 16 with gradient accumulation of 2 (making the effective batch size 32) due to memory constraints. We train the models with a learning rate of $3 \times 10^{-5}$ coupled with the cosine learning rate scheduler using 500 warmup steps. We optimize the binary cross entropy loss using the Adam optimizer [16] with a weight decay of $1 \times 10^{-3}$. During training, we optimize for the AUROC metric with early stopping patience of 5 for a maximum of 30 epochs. We limit the maximum text length to 60 tokens and image bounding box features to 100 per meme. Furthermore, we upsample the text confounders during training by a factor of 3 (which we empirically found to perform best), i.e., confounders are 3 times more likely to be sampled while forming a batch. Additionally, we scale up the loss of the hateful class by a factor of 1.8, attempting to mititgate the effect of the unbalanced training set distribution. We experiment with varying numbers of cross-validation folds; however, our best model is trained under a 15-fold cross-validation setting. We optimize the CV ensemble weights using an evolutionary algorithm with tournament selection, Gaussian noise mutation and uniform crossover [8, 24] strategies. Out of a population of 512 individuals and 100 generations, we pick the set of ensemble weights that achieve the highest AUROC score on the development set (dev-seen).

**Results** The results of our experiments are presented in Table 2, where we show the AUROC on the development set (dev seen), as well as the Phase 1 and Phase 2 test sets (seen and unseen respectively). The number of Phase 1 and Phase 2 submissions were limited to one per day and three in total respectively, and therefore some of the models have missing entries.

ViLBERT CC and VisualBERT COCO are the two best performing multimodal baselines provided by Kiela et al. [15]. We first note that simple supervised fine-tuning of UNITER using the HM task objective (denoted as UNITER in Table 2) already produces considerable improvement compared to these multimodal baselines. Similar fine-tuning with LXMERT performs worse compared to VisualBERT and UNITER. Oscar, despite having been shown to achieve state-of-the-art results on several V+L tasks, obtains lower scores compared to UNITER.

Continuing with UNITER as our best model, we evaluate the effectiveness of a number of variants. UNITER$_{\text{CV10}}$ is trained using 10-fold cross-validation using only the training set and the data split process described in Section 5.3 (first paragraph). The final predictions are obtained using an ensemble optimized using an evolutionary algorithm as described in Section 5.3. This achieves a slightly higher validation score on dev-seen, indicating that cross-validation training and ensembling can help boost performance compared to simple fine-tuning. When confounder upsampling (CFU) and hateful class loss re-weighting (HW) are used (UNITER$_{\text{CV10 + CFU + HW}}$), we observe further improvements on the development set. The Phase 1 results, however, show a larger increase in performance compared

to UNITER, which indicates that these modifications help the model to generalize better. When increasing the number of folds to 15 ($UNITER_{CV15 + CFU + HW}$), we observe additional improvements on the Phase 1 test set; while 30-fold cross-validation ($UNITER_{CV30 + CFU + HW}$) gives an even further increase in dev-seen performance, this does not translate to improvement on the Phase 1 test set. We also find that margin ranking loss ($UNITER_{CV15 + CFU + HW + MRL}$) and YOLO9000 [28] object tags ($UNITER_{CV15 + CFU + HW + YOLO}$) do not improve Phase 1 performance. We notice that, for memes targeting black people and Muslims, YOLO9000 predicts a range of different labels such as Nigerian, Ugandan, or Ethiopian, potentially adding noise to the model. We furthermore find that it fails to detect objects such as guns, missiles and wheelchairs in several images. Such noise may therefore be responsible for the effects on performance. Some approach of fine-tuning YOLO on the HM dataset might help to improve results.

$UNITER_{ENSEMBLE 1}$ consists of predictions obtained as an EA-optimized ensemble over three model variants – UNITER trained using 30-fold cross-validation, UNITER with margin ranking loss, and UNITER with YOLO9000 object tags. Each of these versions of UNITER are themselves results of ensembles of cross-validation style training on different folds (based on the HM training set) and optimized using EA as detailed in Section 5.3. Since these models employ different strategies, we hypothesize their errors will be different and therefore ensembling may lead to complementary effects and help to improve performance. This ensemble achieves the highest Phase 1 score – 79.13. The second ensemble we evaluate is $UNITER_{ENSEMBLE 2}$ which includes the aforementioned three UNITER models as well as three UNITER models trained using 15-fold cross-validation on the training set and using varying sets of seeds and batch sizes. This model leads to an improvement in performance albeit a particularly small one. Finally, $UNITER_{FINAL}$ uses 15-fold cross-validation where part of the development set (dev-seen) is included in the per-fold training set, as described in Section 5.3. This set of three different ensemble models are the ones that we submitted to Phase 2 of the challenge (the final phase of the competition), with the last one, $UNITER_{FINAL}$, achieving an AUROC of 80.53 and ranking us 4th on Phase 2's leaderboard. This improvement in performance could be attributed to the fact that *dev-seen* has more truly multimodal instances, and incorporating these during training enables the model to reason better using both the modalities.

# 7   Error analysis and further model variants

In this section, we present qualitative analyses we conducted for one of our best models by looking at its misclassification patterns on dev-seen, and detail additional experiments we designed in an attempt to improve performance further but which were not successful. We hope that these observations can further inform the basis of future research in the field. Specifically, we designed two more experiments to supplement the model with additional, complementary information, as well as experimented with independent attention blocks for different modalities to mitigate overfitting on the text features. The details of the experiments are presented below.

## 7.1   Analyzing misclassifications

We investigated the misclassifications of $UNITER_{CV15 + CFU + HW}$ (one of the best models of Phase 1) to understand the errors the model makes and identify possible ways in which the model could be improved. We present some of the model's false negatives and false positives in Figure 3. It can be seen that the false negative predictions are on memes that are truly multimodal and hateful in nature when image and text are put together. We find the model cannot reason about target groups; for example, Muslims, wheelchair users, and the Ku Klux Klan. Furthermore, it does not recognize real-life persons such as Anne Frank, or symbols.

The false positives are also mainly multimodal as some of them could become hateful if a different text or image is used. Here too, the ability to identify the target groups might help the model to better identify hate and benign memes.

## 7.2   Capturing target group information

In light of the observations above, we hypothesized that identifying the target groups in the memes may help improve classification performance. Below we present another method aimed at including this information in the model.
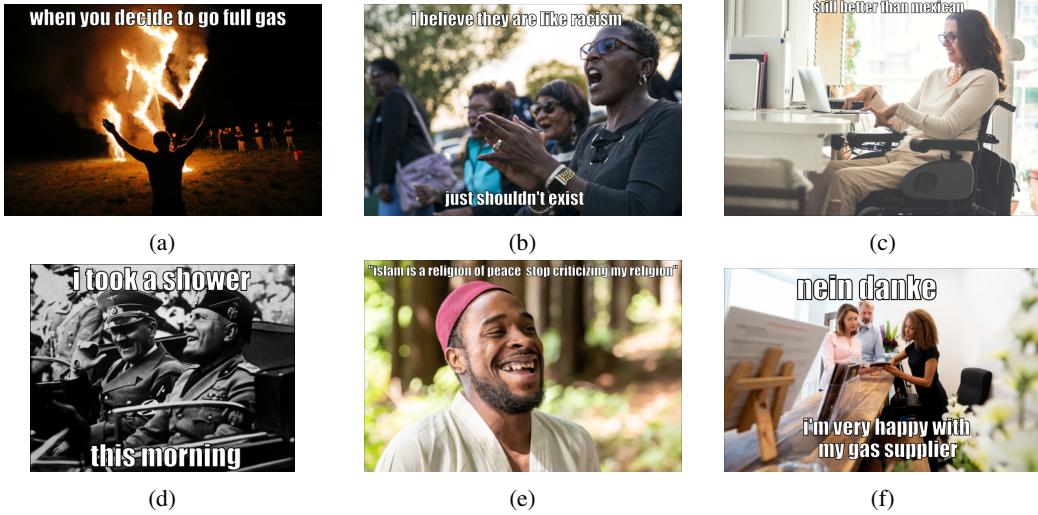
Figure 3: (a)-(c) False negatives: memes from the HM dataset that are labeled as hateful but predicted as being non-hateful. (d)-(f) False positives: memes that are labeled as non-hateful but predicted as being hateful. These are obtained from the UNITER$_{CV15 + CFU + HW}$ variant. Images above are a compilation of assets, including ©Getty Images. The hateful content does not represent the authors' views or opinions in any way whatsoever and it is only used for the purposes of demonstrating model misclassifications.

**Social Bias Frames**    We also performed experiments using the Social Bias Inference Corpus (SBIC) [30]. This dataset contains around 150k structured annotations regarding lewdness, offensiveness, intent to offend, targeted group and more across 44, 761 social media posts. Our first set of experiments focused on multi-task learning (MTL), i.e., simultaneously learning to classify memes as hateful/non-hateful and predicting SBIC classes as an auxiliary task. Specifically, we experimented with both offensiveness classification (binary) and target group prediction (multi-class) as the auxiliary task in separate experiments using the SBIC training data. During MTL training, we trained on alternate batches of the primary and auxiliary tasks with different classification heads for each. However, these setups yielded lower scores on dev-seen compared to simple fine-tuning of UNITER. Subsequently, we experimented with an alternative implementation where we first fine-tune RoBERTa [22] solely on SBIC's target group classification task and then used the model to generate target group labels for each of the memes (based on the meme text as input). We then fed this as input to our UNITER model (similarly to the process used with the YOLO9000 dataset; Section 5.5). However, this too failed to improve results, possibly because of noise in the generated labels and the fact that the meme text alone is unlikely to be sufficient for accurate identification of the target group.

### 7.3    Emotion detection

Rajamanickam et al. [27] showed that multi-task learning with emotion detection can lead to significant improvements in abusive language detection. Taking inspiration from this work, we implemented a multi-task learning variant with emotion detection as the auxiliary task using the GoEmotions dataset [4]. The dataset consists of annotations for 58k Reddit comments over 27 emotion categories including a neutral class (multi-label). Similarly to SBIC, we trained on batches of the primary and auxiliary tasks sampled in different ratios, and with different classification heads for each task. However, the model lead to no improvements in performance on dev-seen.

### 7.4    Overfitting text features

During our experiments, we found that the model performs well on image confounders but struggles on text confounders. This suggests that it relies more on textual features and not as much on image features. In order to mitigate this, we split the internal attention layers of the UNITER transformer architecture into four independent chunks to enable varying dropout rates on different modalities. Specifically, the single attention mechanism was split into *text-to-text, text-to-image, image-to-image*

and *image-to-text* blocks, allowing us to apply different dropout rates on each section. We also ensured that the dropout is consistent across all the multiple-attention heads of each layer, i.e., the dropout across the independent chunks was broadcasted along all the heads. We aimed for a higher *text-to-text* dropout and lower *image-to-text* and *text-to-image* dropouts. We found that this approach also fails to improve results further.

## 8   Conclusion

We proposed a multimodal hateful memes detection framework that ranked 4th in Phase 2 of the 2020 Hateful Memes Challenge launched by Facebook. The proposed solution showcases a number of effective techniques: how to utilize the truly multimodal memes (confounders) in the dataset during training to strengthen the model's reasoning capability; how to effectively optimize an ensemble of models based on cross-validation and an evolutionary algorithm for weight tuning; re-weighting the loss of the minority class; using image features in the correct feature space, i.e., using the exact same checkpoint of the object detector backbone as used by these multimodal architectures during their pretraining phase. We furthermore conducted an error analysis based on model misclassifications and detailed additional experiments that were designed in an attempt to improve performance further but that were not successful. We hope that these observations can further inform future research in the field.

We note that for both our proposed solution and other state-of-the-art architectures such as Oscar there is a lot of room for improvement in terms of their ability to perform "true" multimodal reasoning. Fine-tuning the image extractor during training can help to improve image understanding with respect to the task. However, multimodal reasoning ability is hindered by sub-par image understanding largely due to lack of world knowledge in the current image feature extractor architectures. Identifying target groups as well as public figures and symbolisms are important aspects that can facilitate accurate interpretation of multimodal content. We leave the implementation of alternative approaches to capturing such information for future work. Recent developments in the field, such as ERNIE-ViL [37], a multimodal model that incorporates structured knowledge from scene graphs [13] in its pretraining tasks, is another interesting avenue for future work.

## References

[1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

[2] Prithvi Bhattacharya. 2019. Social degeneration through social media: A study of the adverse impact of 'memes'. In *2019 Sixth HCT Information Technology Trends (ITT)*, pages 44–46.

[3] 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 *European Conference on Computer Vision*, pages 104–120. Springer.

[4] Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A dataset of fine-grained emotions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054, Online. Association for Computational Linguistics.

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

[6] Jessica Drakett, Bridgette Rickett, Katy Day, and Kate Milnes. 2018. Old jokes, new media–online sexism and constructions of gender in internet memes. *Feminism & Psychology*, 28(1):109–127.

[7] Maeve Duggan. 2017. Men, women experience and view online harassment differently. *Pew Research Center*.

[8] A. E. Eiben and James E. Smith. 2015. *Introduction to Evolutionary Computing*, 2nd edition. Springer Publishing Company, Incorporated.

[9] Raul Gomez, Jaume Gibert, Lluis Gomez, and Dimosthenis Karatzas. 2020. Exploring hate speech detection in multimodal publications. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1470–1478.

[10] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and D. Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–6334.

[11] Homa Hosseinmardi, Sabrina Arredondo Mattson, Rahat Ibn Rafiq, Richard Han, Qin Lv, and Shivakant Mishra. 2015. Detection of cyberbullying incidents on the instagram social network. *arXiv preprint arXiv:1503.03909*.

[12] Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6700–6709.

[13] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. 2015. Image retrieval using scene graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[14] Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, and Davide Testuggine. 2019. Supervised multimodal bitransformers for classifying images and text. *arXiv preprint arXiv:1909.02950*.

[15] 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. *arXiv preprint arXiv:2005.04790*.

[16] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

[17] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73.

[18] Zachary Laub. 2019. Hate speech on social media: Global comparisons. `https://www.cfr.org/backgrounder/hate-speech-social-media-global-comparisons`.

[19] 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. *arXiv preprint arXiv:1908.03557*.

[20] Xiujun Li, Xi Yin, C. Li, X. Hu, Pengchuan Zhang, Lei Zhang, Longguang Wang, H. Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*.

[21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.

[22] 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. *arXiv preprint arXiv:1907.11692*.

[23] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23.

[24] B. Miller and D. Goldberg. 1995. Genetic algorithms, tournament selection, and the effects of noise. *Complex Syst.*, 9.

[25] Pushkar Mishra, Helen Yannakoudakis, and Ekaterina Shutova. 2019. Tackling online abuse: A survey of automated abuse detection methods. *arXiv preprint arXiv:1908.06024*.

[26] Vicente Ordonez, G. Kulkarni, and T. Berg. 2011. Im2text: Describing images using 1 million captioned photographs. In *NIPS*.

[27] Santhosh Rajamanickam, Pushkar Mishra, Helen Yannakoudakis, and Ekaterina Shutova. 2020. Joint modelling of emotion and abusive language detection. *arXiv preprint arXiv:2005.14028*.

[28] Joseph Redmon and Ali Farhadi. 2017. YOLO9000: Better, Faster, Stronger. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7263–7271.

[29] 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*, pages 91–99.

[30] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490, Online. Association for Computational Linguistics.

[31] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565.

[32] Amanpreet Singh, Vedanuj Goswami, Vivek Natarajan, Yu Jiang, Xinlei Chen, Meet Shah, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. 2020. Mmf: A multimodal framework for vision and language research. `https://github.com/facebookresearch/mmf`.

[33] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. 2018. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*.

[34] Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.

[35] Amanda Williams, Clio Oliver, Katherine Aumer, and Chanel Meyers. 2016. Racial microaggressions and perceptions of internet memes. *Computers in Human Behavior*, 63:424 – 432.

[36] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.

[37] Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. ERNIE-ViL: Knowledge enhanced vision-language representations through scene graph. *arXiv preprint arXiv:2006.16934*.

[38] Yuke Zhu, O. Groth, Michael S. Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4995–5004.

# A  Appendix

## A.1  YOLO9000 relevant classes

Table 3 presents the list of the 97 YOLO9000 classes we use in our experiments:

Table 3: List of YOLO9000 classes used in our models.

| | | | | |
|---|---|---|---|---|
| military soldier | marcher | domestic goat | domestic sheep | Nigerian |
| Muslimah | Zimbabwean | Sudanese | Eritrean | Punjabi |
| Yemeni | Ugandan | niqab | wild goat | sniper rifle |
| Tommy gun | revolver | machine gun | Black African | mountain gorilla |
| pygmy chimpanzee | gal | Mongol | Tibetan | farm boy |
| cover girl | homeless | amputee | Guyanese | Iraqi |
| heavyweight | Albanian | guy | Nicaraguan | Abyssinian |
| South African | Cameroonian | Haitian | Jordanian | Afghan |
| lady | old man | Ethiopian | Kenyan | lass |
| Senegalese | clown | general | captain | minister |
| ambassador | wheelchair | ram | missile | bomber |
| goat herder | diocesan | eparchy | great grandson | Labrador retriever |
| Korean | schoolchild | Lithuanian | Bolivian | Japanese |
| Arabian | stallion | trotting horse | Omani | Ugandan |
| Bornean | orphan | ape | sweetheart | waiter |
| freight train | German shepherd | Siberian husky | Eskimo dog | hydrogen bomb |
| khakis | Father | hijab | Guinean | Papuan |
| monk | native | Kalashnikov | Mexican | stylist |
| rabbi | beard | kitten | kitty | fireman |
| man | Yugoslav | | | |