# Mapping Memes to Words for Multimodal Hateful Meme Classification

*Giovanni Burbi[1]    *Alberto Baldrati[1,2]    Lorenzo Agnolucci[1]    Marco Bertini[1]

Alberto Del Bimbo[1]

[1] University of Florence - Media Integration and Communication Center (MICC)
[2] University of Pisa
Florence, Italy - Pisa, Italy

`[name.surname]@unifi.it`

## Abstract

*Multimodal image-text memes are prevalent on the internet, serving as a unique form of communication that combines visual and textual elements to convey humor, ideas, or emotions. However, some memes take a malicious turn, promoting hateful content and perpetuating discrimination. Detecting hateful memes within this multimodal context is a challenging task that requires understanding the intertwined meaning of text and images. In this work, we address this issue by proposing a novel approach named ISSUES for multimodal hateful meme classification. ISSUES leverages a pre-trained CLIP vision-language model and the textual inversion technique to effectively capture the multimodal semantic content of the memes. The experiments show that our method achieves state-of-the-art results on the Hateful Memes Challenge and HarMeme datasets. The code and the pre-trained models are publicly available at* `https://github.com/miccunifi/ISSUES`

***Disclaimer: This paper contains hateful content that may be disturbing to some readers.***

## 1. Introduction

Multimodal image-text memes are a unique form of memes that combine visual elements and textual content to convey humor, ideas, or emotions. Unfortunately, some of them are used to perpetuate discrimination against individuals or groups based on their identity [9]. Classifying hateful memes is particularly challenging, as their true intent subtly emerges only when text and images are integrated. Figure 1 shows some examples of multimodal memes where innocuous images and texts turn hateful when combined together.

The Hateful Memes Challenge [9] (HMC) played a pivotal role in advancing research on automated hateful meme classification. The challenge presented curated memes where visual and textual information were tightly intertwined, making the use of multimodal approaches essential. In particular, the organizers crafted non-hateful "confounder" memes by altering only the image or text in the hateful memes while preserving their overall context. These confounder memes show that a seemingly harmless image or text could turn hateful depending on the contextual cues present in the other modality.

To tackle the hateful meme classification task we present a novel approach named ISSUES (mappIng memeS to wordS for mUltimodal mEme claSsification) that leverages a pre-trained CLIP vision-language model and the recently introduced textual inversion technique [1, 7]. Following the terminology introduced in [7], we refer to *textual inversion* as the process of mapping an image into a pseudo-word token residing in the CLIP token embedding space. ISSUES introduces a powerful framework based on three key concepts. First, by exploiting the textual inversion technique, we enhance the multimodal capabilities of the model, allowing the creation of a multimodal representation within the textual embedding space. Second, we disentangle image and text features and we adapt both embedding spaces to the specific downstream task. Finally, inspired by [3], we design an effective multimodal fusion network. Experiments show that our approach achieves SotA results on the challenging HMC [9] dataset and on the HarMeme [18] dataset.

We summarize our contributions as follows:

- We introduce ISSUES, a novel approach for multimodal meme classification which leverages textual inversion in conjunction with a frozen pre-trained CLIP vision-language model;

- To the best of our knowledge, this is the first work that shows that textual inversion can be effectively used to enrich the textual features in a classification task;

- The proposed method achieves SotA results on two datasets: Hateful Memes Challenge and HarMeme.
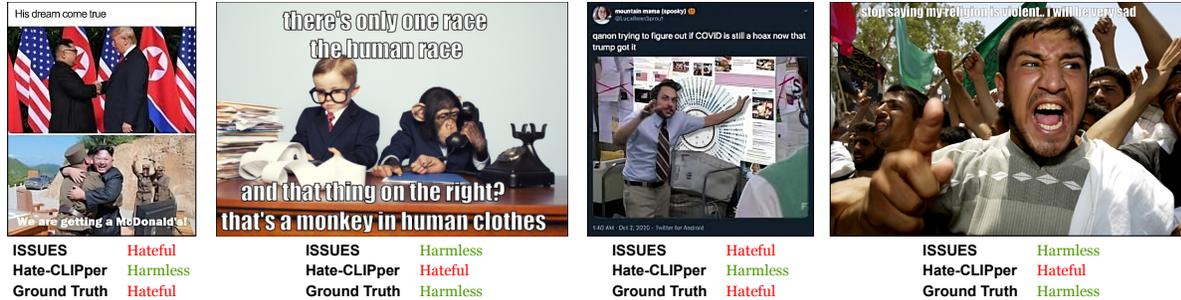
---

* Equal contribution

Figure 1. Examples of multimodal image-text memes. Given a meme, we want to classify whether its content conveys hate. The proposed ISSUES approach is more effective at evaluating the hatefulness of the memes than the state-of-the-art method Hate-CLIPper.

## 2. Related Work

**Hateful Meme Classification** Hateful meme classification has gained prominence as an emerging multimodal task, especially thanks to Facebook's organization of the Hateful Memes Challenge (HMC) [9], which established a benchmark dataset. HMC provided some baselines as a comparison for the competitors, such as VilBERT [15] and VisualBERT [13]. The challenge report [8] showed how all the top five submissions [14, 17, 22, 23, 24] outperformed the baselines, but mainly thanks to the use of external data, additional input features and/or ensemble models.

After the end of the competition, other methods have been presented. For instance, [5] proposes PromptHate, a prompt-based model that prompts pre-trained language models to perform hateful meme classification. Most similar to our work is Hate-CLIPper [10], which explicitly models the cross-modal interactions between the CLIP image and text features through intermediate fusion, thanks to a feature interaction matrix. Similarly to Hate-CLIPper, we rely on CLIP, but we introduce a 2-stage training approach that involves a Combiner network to fuse the features and the use of textual inversion to obtain a multimodal representation within the textual embedding space.

**Textual Inversion** Textual inversion has emerged as a powerful approach for the personalized image generation task in the domain of text-to-image synthesis [7, 11, 20]. For instance, [7] employs the reconstruction loss of a latent diffusion model to carry out textual inversion.

More recently, textual inversion has also been employed in virtual try-on [16] and personalized [6] and composed image retrieval [1, 21]. In particular, [1] proposes SEARLE, an approach that involves training a textual inversion network $\phi$ with a distillation-based loss.

## 3. Proposed Approach

### 3.1. Preliminaries

**CLIP** The CLIP [19] vision-language model is designed to align visual and textual data within a common embedding space. It comprises a visual encoder denoted as $V_E$ and a text encoder denoted as $T_E$. These encoders extracts feature representations $V_E(I) \in \mathbb{R}^d$ and $T_E(E_L(Y)) \in \mathbb{R}^d$ from an input image $I$ and its corresponding text caption $Y$, respectively. Here, $d$ represents the dimension of the CLIP embedding space, and $E_L$ signifies the embedding lookup layer that maps each tokenized word in $Y$ to the CLIP token embedding space $\mathcal{W}$.

**SEARLE** SEARLE [1] is an approach that involves training a textual inversion network, denoted as $\phi$, by distilling knowledge from an optimization-based method. This network exhibits the remarkable capability of efficiently performing textual inversion in a single forward pass. Given an image $I$, the $\phi$ network maps its CLIP visual features, represented by $V_E(I) \in \mathbb{R}^d$, into a pseudo-word token $v_*$ within the CLIP token embedding space $\mathcal{W}$, such that $v_* = \phi(V_E(I)) \in \mathcal{W}$. The pre-training of the textual inversion network $\phi$ aims to ensure that the pseudo-word token $v_*$ not only captures the visual information of $I$ but also enables effective interactions with actual words.

### 3.2. ISSUES

Figure 2 shows an overview of the approach. ISSUES focuses on three main areas: (1) enhancing the textual representation of the meme through textual inversion; (2) adapting the embedding spaces of the pre-trained model using linear projections; (3) employing an expressive multimodal fusion function.

**Enhancing textual representation** As reported in [19], the CLIP model exhibits strong performance in vision-language semantic tasks by utilizing only an image overlaid with some text. This underscores the efficacy of the CLIP visual encoder capabilities in directly extracting meaningful textual semantic information from raw pixel data, thus creating a powerful multimodal representation. To further enhance the expressiveness of the approach, ISSUES aims to integrate the visual information of a meme within the textual embedding space, thereby generating a multimodal representation also within this domain.

To this end, we employ the textual inversion technique [7], which involves mapping the image of a meme $I_M$ to a pseudo-word token $v_*$ residing in the CLIP token
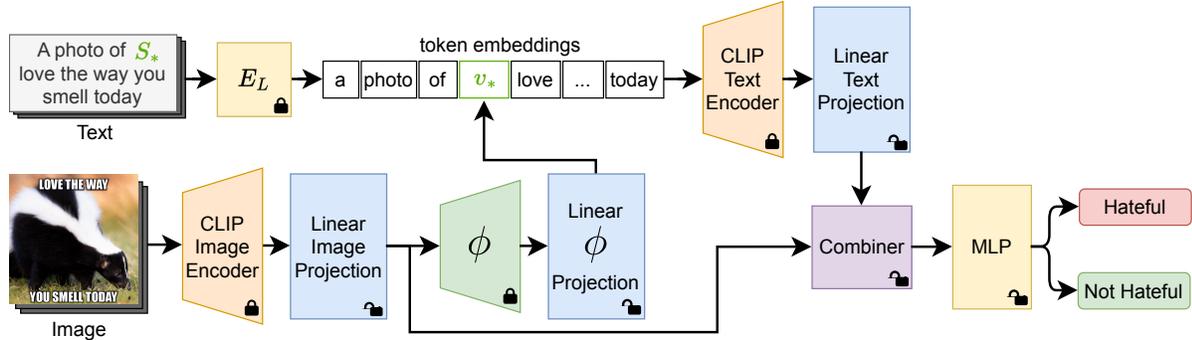
Figure 2. Overview of the proposed approach. We disentangle CLIP common embedding space via linear projections. We employ textual inversion to make the textual representation multimodal. We fuse the textual and visual features with a Combiner architecture. $E_L$ represents the CLIP embedding lookup layer. $\phi$ indicates the SEARLE textual inversion network [1].

embedding space $\mathcal{W}$. To achieve this mapping efficiently and effectively, we leverage the pre-trained textual inversion network $\phi$ proposed in SEARLE [1]. Let $S_*$ represent the pseudo-word associated with $v_*$. Instead of relying solely on the text of the meme {*meme text*}, we compute the textual features from the following prompt: "a photo of $S_*$, {*meme text*}". Remarkably, thanks to the textual inversion process, these text features encompass both textual and visual information, offering a powerful multimodal representation of the meme.

**Adapting the embedding spaces** The pre-training tasks of both $\phi$ and CLIP aim to align text and image representations in a shared multimodal latent space. However, this alignment is not suitable for our task, as memes often contain text and images with different meanings. Thus, it is crucial to disentangle image and text representations in the shared multimodal latent space. To achieve this, following [10] we train linear projections after freezing the $\phi$ and CLIP encoders to adapt the embedding spaces. This way, we effectively disentangle the image and text modalities, enabling the learning of two distinct embedding spaces that are better suited for our specific task.

To effectively adapt the embedding spaces to our task, we follow [2] by adopting a two-stage training strategy. In the first stage, we pre-train the linear projection of the CLIP visual encoder using the same approach and architecture described in [10]. This step allows us to achieve the desired adaptation while retaining the valuable prior knowledge captured in the pre-trained encoders. However, pre-training the textual encoder projection layer using the same approach is not optimal due to the presence of textual inversion, since the pseudo-word tokens generated by $\phi$ would not be incorporated in the pre-training phase.

To address this issue we introduce a second stage of training, where we train the textual and $\phi$ projection layers in conjunction with the multimodal fusion function.

**Multimodal fusion function** Inspired by [3], we adopt the Combiner network as the multimodal feature fusion function in ISSUES. Combiner is designed to take two mul-

timodal representations as input. The first representation originates from the pre-trained projection of the CLIP visual encoder. The second representation stems from the projection of the output of the CLIP text encoder augmented with the textual inversion, resulting in a multimodal embedding within the projected CLIP textual embedding space. We argue that utilizing two multimodal representations empowers the Combiner and enhances its capability to capture the nuanced semantics of memes. The primary objective of the Combiner is to produce a meaningful combined representation, which then serves as input to an MLP for performing the final classification.

## 4. Experimental Results

For all the experiments, we employ the ViT-L/14 as the backbone for CLIP. The whole system is trained with a standard binary cross-entropy loss.

### 4.1. Datasets

We employ two datasets for the experiments: HMC [9] and HarMeme [18]. HMC was proposed by Facebook for the Hateful Memes Challenge and contains 8500, 500, and 2000 synthetic memes in the training, development seen and test unseen sets, respectively. The hatred in HMC is aimed mainly at religion, race, disability, and sex. HarMeme is composed of COVID-19-related real-world memes labeled with three classes: *very harmful*, *partially harmful*, and *harmless*. Similarly to [5], we combine the first two classes in a single *hateful* class. HarMeme contains 3013 and 354 memes for the training and test sets, respectively.

### 4.2. Quantitative Results

We compare ISSUES with baselines and competing methods in Tabs. 1 and 2 for the HMC test unseen and HarMeme test sets, respectively. We do not consider Dis-MultiHate [12] and PromptHate [5] for the HMC dataset because the original papers report only the results on the dev seen split and the pre-trained models are not public.

| Method | Acc. | AUROC |
|---|---|---|
| CLIP Text-Only[†] | 63.50 | 63.43 |
| CLIP Image-Only[†] | 74.65 | 81.35 |
| CLIP Text + textual inversion[†] | 76.65 | 83.31 |
| CLIP Sum[†] | 76.80 | 82.92 |
| VisualBERT COCO [13] | 69.95 | 74.59 |
| ViLBERT CC [15] | 70.03 | 72.78 |
| HMC $2^{nd}$ prize [17] | 69.50 | 83.10 |
| HMC $1^{st}$ prize [24] | 73.20 | 84.49 |
| Hate-CLIPper[†] [10] | 77.15 | 84.36 |
| **ISSUES** | **77.70** | **85.51** |

Table 1. Quantitative results for the HMC test unseen set. Best scores are highlighted in bold. Methods marked with [†] were evaluated by us since the corresponding results are not available.

| Method | Acc. | AUROC |
|---|---|---|
| CLIP Text-Only[†] | 79.38 | 87.00 |
| CLIP Image-Only[†] | 83.90 | 91.59 |
| CLIP Sum[†] | 81.92 | 91.62 |
| DisMultiHate [12] | 81.24 | 86.39 |
| PromptHate [5] | **84.47** | 90.96 |
| Hate-CLIPper[†] [10] | 83.90 | 91.87 |
| **ISSUES** | 81.64 | **92.83** |

Table 2. Quantitative results for the HarMeme test set. Best scores are highlighted in bold. Methods marked with [†] were evaluated by us since the corresponding results are not available.

We evaluate the performance of the models using the AU-ROC [4] – which was the primary evaluation metric for Facebook's Hateful Memes Challenge – and the accuracy.

The presence of confounders renders the HMC dataset particularly challenging for unimodal methods and necessitates robust multimodal reasoning, as evidenced by the poor results obtained by CLIP Text-Only. However, employing only the CLIP image encoder yields significantly better results, demonstrating its capability to extract representations that capture the semantic meaning of both the image and the overlaid text. When we enhance the CLIP textual encoder using the textual inversion (see CLIP Text + textual inversion in Tab. 1), we successfully extract semantic visual features from memes and transfer them into CLIP textual embedding space. This process yields a meaningful multimodal representation, improving the results obtained by the visual encoder and the sum of visual and textual features.

ISSUES outperforms all the baselines and prize winners of the challenge. When compared to the current SotA approach Hate-CLIPper [10], the proposed method shows a significant improvement. Since both models employ the same CLIP backbone, the increase in the metrics is attributable to the effectiveness and impact of the novel components that constitute our approach.

Table 2 shows the results for the HarMeme dataset. Differently from the HMC dataset, utilizing a uni-modal architecture that considers only the text of the meme obtains

| Method | Combiner | 2-stage | TI | Acc. | AUROC |
|---|---|---|---|---|---|
| Base | | | | 77.15 | 84.36 |
| | ✓ | | | **77.80** | 84.68 |
| | ✓ | ✓ | | 77.55 | 85.05 |
| **ISSUES** | ✓ | ✓ | ✓ | 77.70 | **85.51** |

Table 3. Ablation study on the test unseen set of the HMC dataset. TI stands for Textual Inversion.

good results. Consequently, this dataset does not challenge multimodal reasoning to the same extent as the HMC one. However, since the memes are directly collected "in the wild" from social networks, they may exhibit lower quality, but they provide valuable insights into real-world behaviors and interactions. Therefore, the experiments on this dataset offer an excellent opportunity to assess the generalization capabilities of ISSUES. Remarkably, even in this scenario, our model outperforms all existing SotA architectures.

### 4.3. Ablation Study

We present an ablation study on the HMC dataset in Tab. 3. We start with the Hate-CLIPper architecture [10] as the base model and gradually integrate the main parts of IS-SUES: the Combiner network, 2-stage training, and textual inversion. Finally, we measure their relative importance.

Replacing the fusion method in Hate-CLIPper with the Combiner network results in a performance increase. The Combiner enhances the expressiveness of the fusion function, thus improving its ability to model the semantic interaction between text and image representations. The improvement observed with the two-stage training process underscores the importance of carefully adapting the embedding spaces to our specific task by disentangling the image and text representation. Finally, by using textual inversion, we provide two meaningful multimodal representations to the Combiner network in both the visual and textual embedding spaces, which enables better modeling of the semantic meaning of the meme.

### 5. Conclusion

In this paper, we introduce ISSUES, a novel method for classifying hateful memes. Our approach leverages CLIP and textual inversion to generate meaningful multimodal representations in both textual and visual domains. To adapt the pre-trained models, we integrate trainable linear projections into our architecture and adopt a two-stage training strategy. Also, we enhance the expressiveness of the fusion function through a Combiner network, thus improving the modeling of semantic interactions between the meme representations. Experiments on the HMC and HarMeme datasets show that ISSUES achieves SotA results.

# References

[1] Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Alberto Del Bimbo. Zero-shot composed image retrieval with textual inversion. *arXiv preprint arXiv:2303.15247*, 2023. 1, 2, 3

[2] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Conditioned and composed image retrieval combining and partially fine-tuning clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4959–4968, 2022. 3

[3] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Effective conditioned and composed image retrieval combining clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21466–21474, 2022. 1, 3

[4] Andrew P Bradley. The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern recognition*, 30(7):1145–1159, 1997. 4

[5] Rui Cao, Roy Ka-Wei Lee, Wen-Haw Chong, and Jing Jiang. Prompting for multimodal hateful meme classification. *arXiv preprint arXiv:2302.04156*, 2023. 2, 3, 4

[6] Niv Cohen, Rinon Gal, Eli A Meirom, Gal Chechik, and Yuval Atzmon. "this is my unicorn, fluffy": Personalizing frozen vision-language representations. In *European Conference on Computer Vision*, pages 558–577. Springer, 2022. 2

[7] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 1, 2

[8] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Casey A Fitzpatrick, Peter Bull, Greg Lipstein, Tony Nelli, Ron Zhu, et al. The hateful memes challenge: Competition report. In *NeurIPS 2020 Competition and Demonstration Track*, pages 344–360. PMLR, 2021. 2

[9] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *Advances in Neural Information Processing Systems*, 33:2611–2624, 2020. 1, 2, 3

[10] Gokul Karthik Kumar and Karthik Nanadakumar. Hateclipper: Multimodal hateful meme classification based on cross-modal interaction of clip features. *arXiv preprint arXiv:2210.05916*, 2022. 2, 3, 4

[11] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023. 2

[12] Roy Ka-Wei Lee, Rui Cao, Ziqing Fan, Jing Jiang, and Wen-Haw Chong. Disentangling hate in online memes. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 5138–5147, 2021. 3, 4

[13] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 2, 4

[14] Phillip Lippe, Nithin Holla, Shantanu Chandra, Santhosh Rajamanickam, Georgios Antoniou, Ekaterina Shutova, and Helen Yannakoudakis. A multimodal framework for the detection of hateful memes. *arXiv preprint arXiv:2012.12871*, 2020. 2

[15] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019. 2, 4

[16] Davide Morelli, Alberto Baldrati, Giuseppe Cartella, Marcella Cornia, Marco Bertini, and Rita Cucchiara. Ladi-vton: Latent diffusion textual-inversion enhanced virtual try-on. *arXiv preprint arXiv:2305.13501*, 2023. 2

[17] Niklas Muennighoff. Vilio: State-of-the-art visio-linguistic models applied to hateful memes. *arXiv preprint arXiv:2012.07788*, 2020. 2, 4

[18] Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md Akhtar, Preslav Nakov, Tanmoy Chakraborty, et al. Detecting harmful memes and their targets. *arXiv preprint arXiv:2110.00413*, 2021. 1, 3

[19] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2

[20] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2

[21] Kuniaki Saito, Kihyuk Sohn, Xiang Zhang, Chun-Liang Li, Chen-Yu Lee, Kate Saenko, and Tomas Pfister. Pic2word: Mapping pictures to words for zero-shot composed image retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19305–19314, 2023. 2

[22] Vlad Sandulescu. Detecting hateful memes using a multimodal deep ensemble. *arXiv preprint arXiv:2012.13235*, 2020. 2

[23] Riza Velioglu and Jewgeni Rose. Detecting hate speech in memes using multimodal deep learning approaches: Prize-winning solution to hateful memes challenge. *arXiv preprint arXiv:2012.12975*, 2020. 2

[24] Ron Zhu. Enhance multimodal transformer with external label and in-domain pretrain: Hateful meme challenge winning solution. *arXiv preprint arXiv:2012.08290*, 2020. 2, 4