# Emotion-Aware Multimodal Fusion for Meme Emotion Detection

Shivam Sharma[1,4], Ramaneswaran S[3], Md. Shad Akhtar[2] and Tanmoy Chakraborty[1]

[1]*IIT Delhi, India*; [2]*IIIT Delhi, India*; [3]*VIT, Vellore, India*; [4]*Wipro AI Labs, India*

{shivam.sharma,tanchak}@ee.iitd.ac.in, shad.akhtar@iiitd.ac.in, s.ramaneswaran2000@gmail.com

**Abstract**—The ever-evolving social media discourse has witnessed an overwhelming use of memes to express opinions or dissent. Besides being misused for spreading malcontent, they are mined by corporations and political parties to glean the public's opinion. Therefore, memes predominantly offer affect-enriched insights towards ascertaining the societal psyche. However, the current approaches are yet to model the affective dimensions expressed in memes effectively. They rely extensively on large multimodal datasets for pre-training and do not generalize well due to constrained visual-linguistic grounding. In this paper, we introduce MOOD (Meme emOtiOns Dataset), which embodies six basic emotions. We then present ALFRED (emotion-Aware muLtimodal Fusion foR Emotion Detection), a novel multimodal neural framework that (i) explicitly models emotion-enriched visual cues, and (ii) employs an efficient cross-modal fusion via a gating mechanism. Our investigation establishes ALFRED's superiority over existing baselines by 4.94% F1. Additionally, ALFRED competes strongly with previous best approaches on the challenging Memotion task. We then discuss ALFRED's domain-agnostic generalizability by demonstrating its dominance on two recently-released datasets – HarMeme and Dank Memes, over other baselines. Further, we analyze ALFRED's interpretability using attention maps. Finally, we highlight the inherent challenges posed by the complex interplay of disparate modality-specific cues toward meme analysis.

**Index Terms**—Memes, multimodality, emotion analysis, social media, information fusion.

✦

## 1 INTRODUCTION

MEMES on social media represent a new digital artifact genre that has become ubiquitous with time. Internet memes contain a short piece of text embedded over an image, often expressing information sarcastically or humorously, but sometimes in illicit ways. Internet memes also offer rich potential to understand or sway the sentiment and opinion of communities on social media, essentially facilitating a systematic study of their affective characteristics.

The increase in the amount of multimodal content being disseminated over the web has spurred innovation in allied areas involving multimodality in general; however, fewer efforts have been made toward analyzing memes, especially to the extent its inherent complexity solicits. Memes, having multimodal constructs, require a joint interpretation of both the embedded text and the visuals to assimilate the intended meaning comprehensively. Although existing approaches perform much better on conventional multimodal tasks, particularly the ones involving visual-linguistic grounding [1], they are yet to deliver for other scenarios like memes [2], [3], [4], [5], [6].

One of the critical challenges within the studies related to social media content analysis is the *subjective perception problem* [7], which leads to ambiguous data labeling. Consequently, several recent participatory efforts involving the detection of emotion from memes consider the categories that capture higher levels of affective abstraction like *humor, sarcasm, offense, and motivation* [2], [8]. Although these efforts effectively capture affective phenomena from memes that pertain to macro societal aspects like political, socio-cultural, demographics, etc., they constrain fine-grained analysis of meme emotions. Efforts are needed towards studies that target more fine-grained analysis of user behavior, decisions, and perceptions [9], encompassing a broader spectrum of emotions that help address the multimodal affective characterization at the fundamental level.



Fig. 1: Example of memes: (a,b) text modality is misleading; (c,d) same image but emotion being differentiated by the text.

Additionally, the inherent multimodality involving image-text configurations presents an additional challenge in detecting emotions due to the complexity posed by the implicit background context abstracted by *memetic visuals*. For instance, memes shown in Figs. 1(a) and 1(b) depict faces with expressions conveying *anger* and *sadness*, respectively; whereas the corresponding embedded text suggests something different, thereby vesting the required onus of meme emotion detection, primarily upon the visual cues.

Also, examples shown in Figs. 1(c) and 1(d) carry the same background image, but their corresponding embedded texts convey contrasting emotions – *anger* and *joy*, respectively, indicating triviality of the visual cues present. These aspects pose challenges like visual-linguistic dissociation and cross-modal noise, and induce reasoning complexity by having dependencies on implicit contextual cues, effectively inhibiting the overall progress towards addressing such non-trivial tasks.

In this work, we propose a novel task of classifying the emotions expressed within memes amongst *six* basic Ekman [10] emotions. To this end, we also introduce a manually curated large-scale multimodal dataset, MOOD (Meme emOtiOns Dataset). MOOD constitutes real multimodal memes (images with overlaid text) expressing various emotions that are discretely mapped to the six fundamental Ekman [10] emotions – *fear, anger, joy, sadness, surprise,* and *disgust* via manual annotation. We benchmark this dataset against several unimodal and multimodal systems emulating competitive baselines for meme emotion detection and report class-wise, along with macro-averaged performances. Further, we investigate the design of an effective approach to detect emotions from memes and propose ALFRED, a multimodal approach that employs systematic modality-specific interactions via gating. ALFRED constitutes (a) gated multimodal fusion (GMF) towards explicitly incorporating emotion-enriched visual features, followed by (b) gated cross-attention (GCA) to fuse emotion-enriched image and text representations, conditioned upon the visual cues learned. We observe significant gains by ALFRED over other strong baselines on the MOOD dataset, along with competitive performance on Memotion tasks [2], followed by distinct indications of ALFRED's strong generalizability over other related tasks such as HarMeme [11] and Dank Memes [12].

We also examine the interpretability of the predictions by analyzing visual attention maps corresponding to the encoding mechanisms that ALFRED employs and highlight both affective affinity and limitations exhibited therein. Besides discussing and analyzing the overall and emotion-specific performances of various systems and ALFRED in detail, we delineate the contributory aspects of MOOD and ALFRED. Finally, we perform an extensive error analysis elucidating some imminent challenges and limitations like *modality-specific obscurity* and *thematic overlaps* rendered by the complex memetic dynamics.

Through this work, we intend to address the imperative necessity of characterizing multimodal content like memes and their fundamental affective spectrum by emphasizing their esoteric visual semiotics. In particular, we make the following contributions:

- We introduce MOOD (Meme emOtiOns Dataset) that captures six basic *Ekman's* emotions for memes.
- We propose ALFRED (emotion-Aware muLtimodal Fusion foR Emotion Detection): a multimodal neural framework that uses affect-enriched features from memes and fuses them via a gated cross-attention mechanism.
- We benchmark the dataset via several unimodal and multimodal baselines and discuss their limitations.
- Further, we empirically demonstrate the efficacy of ALFRED over strong baselines on the MOOD and Memotion datasets.
- Finally, we perform interpretability analysis and establish the generalizability of ALFRED on meme datasets capturing distributions beyond six basic Ekman emotions.

**Reproducibility:** The source codes and the sample dataset are uploaded at: `https://github.com/LCS2-IIITD/ALFRED_Meme EmotionDetection`.

## 2 RELATED WORK

Several studies contributed to the understanding of detecting emotions from multimodal content, encompassing modalities like images and text, along with signals like audio, video, EEG, eye movement, etc. There have also been numerous recent efforts toward analyzing memes and detecting various allied harmful aspects. Since there have been limited exploratory studies on detecting emotions from memes, the field can benefit significantly from the findings of the aforementioned applications. We systematically review these areas to set the necessary background.

**Multimodal emotion detection:** Emotion detection is a well-studied area explored for various modalities such as text, speech, and audio [13], [14]. Significant emphasis has been laid on multimodal emotion detection as well. Unlike traditional emotion detection tasks involving single modalities, multimodal approaches require a mechanism to effectively learn features from multiple correlated modalities. An initial effort in this domain [15] proposed multi-kernel learning based deep-CNN towards emotion and sentiment recognition on different multimodal datasets. This was followed by proffering a pooling-based fusion mechanism in [16] along with introducing a multimodal social media dataset (and metadata) from Reddit towards the domain of emotion classification. Efforts have been made toward multimodal feature characterization [17], wherein authors introduced a Tumblr-based multimodal dataset and demonstrated the efficacy of multimodal approaches when contrasted with their unimodal counterparts. Another notable finding [18] explored an approach advocating co-ordinated representation learning for multimodal emotion recognition. The authors used a recurrent neural network to emulate correlated attention and calculate the correlation between EEG and eye movement signals.

Recent efforts include a Transformer-based inter-modality attention mechanism [19] with self-supervision [20], while materializing the design for fusing features from different modalities for multimodal emotion recognition. While there has been significant progress from a computational standpoint toward emotion recognition, Mittal et al. [21] attempted to investigate an approach based on Frege's Context principle, which provides different interpretations of context for emotion recognition. The authors studied different interpretations using modality-specific features, semantic content, and depth-map to develop their algorithm. Although these efforts pave the way for addressing critical challenges prevalent within the affect-oriented applications for multimodal content, there is still scope for further exploring multimodal content representing dynamic cross-modal semiotics, like memes.

**Meme analysis:** A significant influx of memes from online fringe communities, such as Gab, Reddit, and 4chan, to mainstream platforms, such as Twitter and Instagram, resulted in a massive epidemic of intended harm [22]. This has imminently solicited addressing the prevailing challenges from the computational social science point of view. Towards this end, several datasets capturing offensiveness [23], hatefulness [24], [25], and harmfulness [26] in memes have been curated. Besides the tasks corresponding to these resources, there are a variety of other tasks, such as detecting sexism [27], racism [28], and harmful propaganda [29] from memes, that have been explored from the perspective of critical discourse analysis.

Participatory events like the Facebook Hateful Meme Challenge [24] and shared-task on detecting *hero, villain* and *victim* from memes [30] have laid a strong foundation for community-level initiatives for detecting hate speech [31], [32], [33], [34], [35] and connotative role-labels in memes. As part of these challenges, several interesting approaches besides ensembling large language models, utilizing meta information, attentive interactions, and adaptive loss are attempted in the multimodal setting [36], [37], [38], [39]. Other notable insights from meme analyses suggest the utility of commonsense knowledge [40], web entities, racial aspects [26], [41], and other external cues for detecting offense, harm, and hate speech in memes.

Most of these efforts either address the detection tasks at various levels for harmfulness (see a recent survey [42]) or design ensemble techniques lacking cost-optimality. However, they tend to ignore the primary spectrum of emotions that plays a crucial role in cascading any adverse effect over social media. The current study aims to address this fundamental aspect of detecting the basic emotions of memes.

**Emotion detection from memes:** Although several studies have analyzed emotions of social media content, fewer efforts have been made toward characterizing the emotions of Internet memes. French [43] studied the correlation between the semantics of a meme and the textual discussions in the thread related to a multimodal post. Their study indicates the effectiveness of memes as a sentiment predictor over social media. `Memotion` [2], [8], a series of participatory shared tasks on meme emotion classification, initiated the task of detecting meme emotions at different levels of granularity. Their dataset was curated towards three sub-tasks, emulating various combinations of *multi-class/label* formulations. The three sub-tasks were – (a) *sentiment classification*: multi-class classification amongst positive and negative categories; (b) *emotion classification*: multi-class classification amongst categories *sarcastic, humorous, offensive and motivational*; and (c) *quantification*: multi-class/label classification amongst emotion intensities, represented by *slightly, mildly* and *very* and across emotion categories. Participants of this task explored several unimodal and multimodal approaches. For unimodal feature extraction, participants used a variety of models such as BERT [44], GloVe [45] for text modality and pre-trained image models such as EfficientNet [46] and ResNet [47] for image modality. Singh et al. [48] and Vlad et al. [49] used multi-task learning to jointly predict emotion and sentiment. While these solutions, as discussed in the Introduction section (c.f. Sec. 1), address the detection of affect categories at a higher level of abstraction, they do not consider the multimodal emotion characterization from a fundamental perspective. This effectively renders the investigation of basic emotions from memes obscure.

## 3  MEME EMOTIONS DATASET (MOOD)

Memotion dataset [2], [8] includes affective categories like *motivation, offense, sarcasm*, and *humor*, which represent high-level emotion abstraction within memes. Although these categories are critical for studying the imminent implications of memetic discourse over social media, they are insufficient for characterizing their impact on an individual's psyche. Therefore, as part of this work, we aim to set up a framework for addressing emotion recognition from memes w.r.t. the basic emotions. Ekman and Cordaro [10] empirically suggested that human beings exhibit six basic

TABLE 1: Summary of MOOD & Extended AffectNet.

| Dataset | Split | # memes | FER | AGR | JOY | SAD | SPR | DGT |
|---|---|---|---|---|---|---|---|---|
| MOOD | Train | 7004 | 612 | 1413 | 1920 | 1822 | 855 | 382 |
| | Val | 1500 | 131 | 292 | 394 | 416 | 168 | 99 |
| | Test | 1500 | 128 | 312 | 392 | 398 | 185 | 85 |
| | Total | 10004 | 871 | 2017 | 2706 | 2636 | 1208 | 566 |
| Ext. AffectNet | Subset | 50389 | 6540 | 10000 | 10000 | 10000 | 10000 | 3849 |
| | Add-On | 1447 | 162 | 198 | 400 | 472 | 169 | 46 |
| | Total | 51836 | 6702 | 10198 | 10400 | 10472 | 10169 | 3895 |

emotions, namely *fear* (FER), *anger* (AGR), *joy* (JOY), *sadness* (SDN), *surprise* (SPR), and *disgust* (DGT). Besides constituting the primary spectrum for studying the human affective response, basic emotions decide the overall affective tone of memes towards analyzing their disseminative outcomes. To this end, we manually curate MOOD, a multimodal dataset for detecting basic emotions from memes.

### 3.1  Dataset collection and de-duplication

The memes were collected primarily from two sources – Google image search[1] and imgflip[2]. We used keywords like 'happy memes', 'depression memes', 'sad cat memes', etc., to crawl a diverse set of memes, capturing the six basic Ekman emotions. Many duplicates in the collected set were removed using an off-the-shelf API called imagededup[3], followed by manually filtering out the low-quality memes. We used a set of filtering criteria for memes to ensure that the memes collected were of high quality. A meme is discarded if any of the following criteria are met: (1) The resolution of the meme is bad, such that the meme image is unclear or the readability of the meme text is affected; (2) If the meme did not exhibit any of the 6 Ekman emotions; (3) If the meme had content that induces hate towards or is harmful to individuals or communities; (4) Contain any personal information of a user; (5) Contains text that is not in English or is code-switched. This resulted in a total of $10,004$ memes (c.f. Table 1), that constituted our primary proposed meme emotion dataset MOOD. The dataset constitutes *generic* memes corresponding to the six basic Ekman emotions, on general topics like *birthdays, relationships, family, friends* (to name a few), with an appropriate mix of human subjects, pop-culture references, animated characters, and animals.[4]

There is a variation in the relative proportions of the memes collected for different emotion categories. Many memes were collected from Google image search results for categories *anger, joy, sadness, and disgust.* Interestingly, almost half of the *fear and surprise* memes we collected are obtained from imgflip. This suggests the categorical diversity that the platform can provide. Also, very few memes for categories *disgust and sadness* could be sourced from imgflip, suggesting the positive sentiment-based genre that dominates the platform. This distribution also reflects the realistic availability of such multimodal data over different platforms.

MOOD consists of a total of 2017, 2706, and 2636 memes from categories *anger, joy* and *sadness*, respectively, constituting the majority share within the dataset. These are followed by the 1208 from *surprise*, 871 from *fear* and the least from disgust, with 576 memes. The summary of MOOD can be observed from Table 1.

---

1. Google Images ⧉
2. Imgflip ⧉
3. imagededup ⧉
4. See Appendix A.3 for more details on thematic-distribution in MOOD

(a) Fear  (b) Anger  (c) Joy  (d) Sadness  (e) Surprise  (f) Disgust
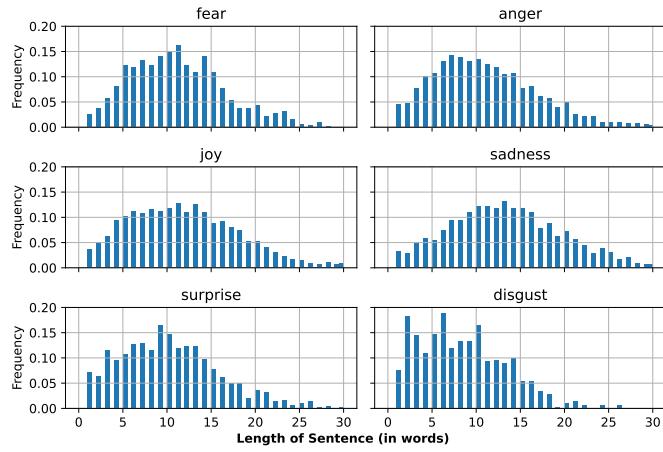
Fig. 2: Examples depicting memes for six basic *Ekman* emotions from our dataset, MOOD.



Fig. 3: Normalized histograms of meme text length per class.



(a) Fear  (b) Joy  (c) Sadness

Fig. 4: Examples of the images added to *extend* AffectNet dataset, especially towards capturing the *non-human* subjects within the meme visuals like, (a) *fearful* 'Spongebob', (b) *joyous* cat and (c) *sad* puppy, depicted in subfig. (a), (b) and (c), respectively.

TABLE 2: Prescribed guidelines for MOOD's annotation.

| Annotation Guidelines | |
| --- | --- |
| 1 | Emotion labeling should consider the meme author's perspective. |
| 2 | Emotion labels should emphasize upon meme's content only. |
| 3 | Emotion label should be one of the six basic Ekman emotions. |
| 4 | Both textual and visual cues should be factored-in while annotating. |
| 5 | Pop-cultural or vernacular-oriented references must be queried from additional information sources wherever needed. |
| 6 | Sample leading to multiple possible interpretations, or lack of mapping to any of Ekman's six emotions, should be skipped. |

Further, towards capturing the emotional signals expressed via visual modality, we leverage AffectNet [14], a large-scale human facial expression dataset that consists of around $400K$ manually annotated facial expression images capturing *neutral* and a total of 7 emotions, namely *happy, sad, surprise, fear, disgust, anger* and *contempt*. We randomly sample around $50K$ images corresponding to the Ekman emotion categories (excluding neutral and contempt from AffectNet) for our framework. AffectNet contains facial expressions only for *humans*, rendering the animated and graphical emotion depiction and comprehension obscure, as required for our scenario. In order to robustly identify the emotions in realistic meme images within MOOD, especially the ones containing cartoons, animated figures, and animals, we needed to emphasize our multimodal framework's capability to handle them. To this end, we collected a dataset of $1,447$ images (c.f. Table 1) with animated characters, cartoons, and animals and manually labeled them with Ekman emotion classes. We queried the web for using a simple formatted query string as 'reactionary templates for' + <emotion category> + <animal/cartoon>, and followed similar filtering/annotation process as for MOOD. These images represent various standard reactionary templates disseminated over social media. Fig. 4 depicts some examples of the images added to *extend* the AffectNet dataset, especially towards capturing the *non-human* subjects within the memes with query prefix as 'reactionary template for', and suffixes as: *fearful Spongebob*, *joyous cat* and *sad puppy*, depicted in Figs. 4(a), 4(b) and 4(c), respectively. We call this extended set **Ext. AffectNet**.

### 3.2 Text Length Analysis

The meme text length analysis indicates the complexity that could be posed within a corpus. The more diverse the text-length distribution is for each category, the more difficult it could be to model the sequence and underlying association. It can be visualized from Fig. 3 that the distributions for *anger, sadness, and joy* are relatively more close to being *normal*, also suggesting that there is a consistent pattern for the creation of content for memes from these categories. In contrast, *disgust*, along with *surprise and fear* (with slight variability), have relatively more variation regarding the text lengths used. This suggests the challenge it poses to the language models and the diversity with which such content is created online.

### 3.3 Annotation

Two annotators annotated the dataset, while a consolidator oversaw the entire annotation exercise. One of the annotators is male, while the other female, and their ages range from 24-35 years. Moreover, both of them were professional lexicographers and social media savvy. On the other hand, the consolidator was an expert working in the fields of Computational Social Sciences and NLP. Before starting the annotation process, they were briefed on the task using detailed guidelines. They were requested to assess memes' textual and visual content towards the final adjudication of the meme emotion. In particular, they were asked to identify the emotion that the meme's author is trying to express via a *multi-class* labeling setup. A tabulated list of prescribed guidelines adopted towards the annotation is shown in Table 2. We conducted

(a) Fear    (b) Anger    (c) Joy    (d) Sadness    (e) Surprise    (f) Disgust

Fig. 5: Word clouds depicting the category-wise lexicon comprising the embedded texts for memes in the MOOD dataset.



(a) I+T (FER)    (b) I (AGR)    (c) T (SPR)

Fig. 6: Example memes depicting modality-specific influence for emotion recognition. I: Image, T: Text.

the annotation process in two stages – a dry run and a final annotation stage. The Cohen's Kappa [50] was computed to assess the inter-annotator agreement prior to the commencement of the final annotating process and was found to be nearly perfect with a score of 0.86.

Based on the annotation, MOOD dataset can be exemplified via Fig. 2, which depicts example memes for each of the six *Ekman* emotions from the MOOD dataset. Although most memes in MOOD are designed by their authors to disseminate some form of humor via sarcasm, satire, or benign limericks, their primary objective is to resonate with the consumer's emotional appeal. Typically, this leads to the memes exhibiting a *primary emotion* by design. Samples depicted in Figs. 2(a)-2(f) are specially hand-picked to ensure a clear understanding of the annotation strategy. As can be observed from the samples, the primary emotions conveyed in the form of *disgust, anger, surprise, fear, sadness* and *joy* within the respective samples are expressed by both text and visual cues within the memes. The modality-specific expressivity varies extensively across MOOD, which constitutes the key multimodal challenge posed while performing analysis over realistic memes. The demonstration via Fig. 6 depicts the independent and joint influence of both image and text modalities via memes. Fig. 6(a) depicts *fear* through both text and image; Fig. 6(b) shows *anger* mainly via the facial expressions; and Fig. 6(c) expresses *surprise* only via text.

### 3.4 Lexical Analysis of MOOD

The lexical summary exhibited by the textual cues present in the memes, in the form of overlaid text (c.f. Fig. 5), suggests interesting characteristics. All categories except *surprise*, prominently exhibit affect-enriched lexicon, including nouns, adjectives and verbs, exemplified as: *disgust* (gross, yuck, eww), *anger* (punks, hate, mad), *fear* (afraid, screaming, therapist), *sadness* (depressed, anxiety, lonely) and *joy* (friend, happy, love). Whereas, as elucidated from Fig. 5 (e), *sur-*
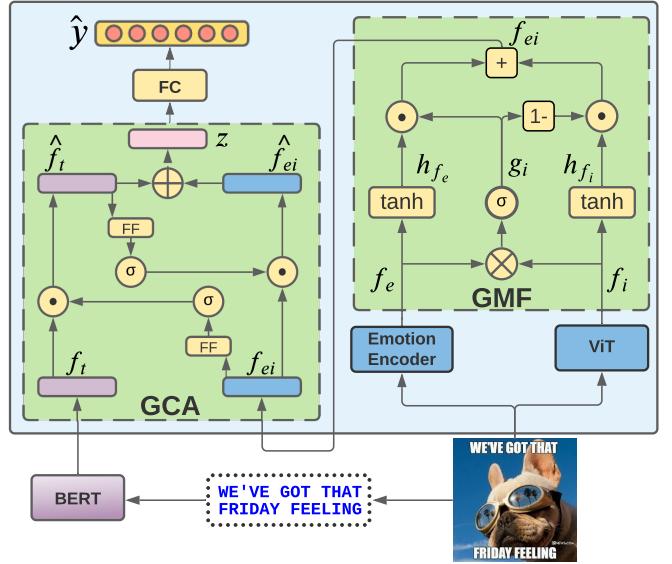


Fig. 7: ALFRED's model architecture. GCA: Gated Cross Attention module, GMF: Gated Multimodal Fusion module, $\otimes$: Low-rank Bilinear Pooling, $\oplus$: Concatenation.

*prise*-based emotion characterization relies significantly upon contextual cues, instead of lexical ones.

Additionally, Table 7 (c.f. Appendix A.3) shows top 10 frequent words in the meme text for each category after masking the 'category-keywords'. The corresponding TF-IDF scores are given in the parenthesis. Some category-specific relevant words can be observed at the top of the category columns, except for *anger*. We do not see any particular word that is usually used in the contexts conveying *anger*. This suggests the complex constructs people typically use while creating memes when conveying *anger*. Essentially, the anger might not be conveyed using explicit keywords but by using implicit and complex references instead. This indicates the complexity the affective content from social media can pose to multimodal systems. Fig. 3 shows the length distribution of meme text for each class; we observe no significant differences between the classes.

## 4 PROPOSED APPROACH

As shown in Fig. 7, our proposed approach utilizes the meme image and embedded text, extracted using Google GCV OCR[5] (GOCR) as primary inputs. For encoding image features, we use Vision Transformer (ViT) [51] and for meme text, we use BERT [44].

5. Google Cloud Vision OCR API ⤤

We model emotion-aware image representations by explicitly incorporating visually depicted emotions. We extract emotion features from the images using a ViT-based encoder network, pre-trained using the *extended* AffectNet dataset. We then pass the emotion and image features through a *gated* multimodal unit to obtain an emotion-aware image representation, as shown in Fig. 7. Further, we feed the emotion-aware image and textual representation through a gated cross-modal attention module that facilitates selective cross-modal attuning and blocking. Finally, we concatenate the two resultant updated representations and pass the joint representation of the meme through a feed-forward network toward the classification task. We describe each of these modules in detail below.

## 4.1 Unimodal Feature Extraction

We use pre-trained unimodal encoders to obtain representations of the image ($i$) and text ($t$) for a given meme ($M$), as described below.

**– Image Encoder:** We use ViT [51] initialised with ImageNet weights as the image encoder and obtain an $m \times 768$ dimensional output corresponding to $m$ image patches, i.e., $\mathbf{f_i} = \text{ViT}(i) \in \mathbb{R}^{m \times 768}$.

**– Text Encoder:** We use pre-trained BERT [44] as the text encoder. Specifically, we take the token-level representations corresponding to $n$ tokens from the last hidden layer: $\mathbf{f_t} = \text{BERT}(t) \in \mathbb{R}^{n \times 768}$.

## 4.2 Emotion Feature Extraction

There are well-established variants of approaches, specially tailored towards capturing semantic objects and salient features, like RCNNs [52], [53], YOLO [54], and CenterNet [55], [56], to name a few; yet, there is a dearth of solutions addressing complex emotion features from visuals, especially memes. Towards explicitly incorporating visually-depicted emotion features as part of our proposed methodology, we build an emotion feature extraction model by fine-tuning a ViT-based image patch encoder for the emotion expression classification task. We then freeze its weights toward extracting relevant emotion-enriched cues from a given meme. To this end, we leverage the AffectNet dataset, a large-scale dataset of typical human facial expressions; but we also extend it by adding non-human subjects.

After pretraining a ViT base model using over $50K$ emotion-enriched images from the *extended* AffectNet dataset for emotion classification, we freeze its weights for extracting emotion-enriched image features. These features are incorporated as part of `ALFRED`. This is expressed as $\mathbf{f_e} = \text{ViT}(i) \in \mathbb{R}^{m \times 768}$.

## 4.3 Gated Multimodal Fusion (GMF)

As part of effectively incorporating visual cues from memes for meme emotion recognition, it is crucial to optimally infuse emotion-enriched input signals while emphasizing other relevant visual cues. This becomes critical while fusing features from similar source modalities.

To induce selective processing of input features, we adapt a *gated multimodal unit* [57] by using low-rank bilinear pooling (LRBP) [58] while computing a sigmoid-based gating weight, instead of a simple concatenation based approach. The motivation for this change is the requirement to fuse the representations coming from the same input source (i.e., image), towards which a Hadamard product-based interaction is empirically observed to

be preferable over a concatenation-based approach. This module performs a fusion of the emotion features ($\mathbf{f_e}$) that are extracted using the emotion encoder, and the meme image features ($\mathbf{f_i}$) obtained using the image encoder to finally obtain emotion-aware image features ($\mathbf{f_{ei}}$). We do this as we have two different types of image encodings, $\mathbf{f_e}$ and $\mathbf{f_i}$. Such a fusion trades off on how much novel information is required from each encoding using a sigmoid-based gated fusion mechanism.

$$\mathbf{h_{f_i}} = \tanh(\mathbf{f_i W_i} + \mathbf{b_i}); \qquad \mathbf{h_{f_e}} = \tanh(\mathbf{f_e W_e} + \mathbf{b_e}) \quad (1)$$

$$\mathbf{g_i} = \sigma([\mathbf{h_{f_i}} \otimes \mathbf{h_{f_e}}]\mathbf{W_g}); \qquad \mathbf{f_{ei}} = \mathbf{g_i h_{f_e}} + (1 - \mathbf{g_i})\mathbf{h_{f_i}} \quad (2)$$

where $\mathbf{W_i}$, $\mathbf{W_e}$, $\mathbf{W_g} \in \mathbb{R}^{768 \times 768}$ are the weights for transforming image features ($\mathbf{f_i} \in \mathbb{R}^{m \times 768}$), emotion features ($\mathbf{f_e} \in \mathbb{R}^{m \times 768}$), and low-rank bilinear pooling based fusion ($\mathbf{g_i} \in \mathbb{R}^{m \times 768}$) of their latent features ($[\mathbf{h_{f_i}} \otimes \mathbf{h_{f_e}}] \in \mathbb{R}^{m \times 768}$), respectively. The bias terms, $\mathbf{b_i}$ and $\mathbf{b_e} \in \mathbb{R}^{768}$, correspond to the representation learning for image and its emotion-enriched signals, respectively. $\sigma$ denotes the *sigmoid* activation function. Also, numpy-like broadcasting is inherently applied wherever applicable via PyTorch API[6].

## 4.4 Gated Cross Attention (GCA)

Prominent conventional approaches that leverage co-attentional transformers-based layers have been observed to perform well in scenarios involving visual-linguistic grounding [59], [60]. However, they exhibit sub-optimal results while modeling memes [26]. This could be likely due to the cross-modal noise being captured and attended to while learning dissociated cues from modality-specific meme components. Towards regulating the inherent effect of cross-modal noise, we modify the cross-attention mechanism [61], by incorporating the adaptive co-attention strategy [62]. Instead of incorporating self-attention layers for cross-modal attention, we perform gating over one modality (visual) first, followed by weighting the other modality (textual). We then perform gated attention for the first modality (visual) using the weighted textual representation to obtain its feedback-based representation. We call this *Gated Cross Attention* mechanism. It facilitates the extraction of useful features from emotion-aware image ($\mathbf{f_{ei}} \in \mathbb{R}^{m \times 768}$) and textual ($\mathbf{f_t} \in \mathbb{R}^{n \times 768}$) features. Thus, we obtain new feature representations, $\hat{\mathbf{f_{ei}}} \in \mathbb{R}^{m \times 768}$ and $\hat{\mathbf{f_t}} \in \mathbb{R}^{m \times 768}$, as follows:

$$\mathbf{h_{f_{ei}}} = \sigma(\mathbf{f_{ei} W_{ei}} + \mathbf{b_{ei}}); \quad \alpha_{\mathbf{ei}} = \text{softmax}(\mathbf{h_{f_{ei}} W_{\alpha_{ei}}} + \mathbf{b_{\alpha_{ei}}}) \quad (3)$$

$$\hat{\mathbf{f_t}} = \alpha_{\mathbf{ei}} \times \mathbf{f_t} \in \mathbb{R}^{m \times 768} \quad (4)$$

$$\mathbf{h_{\hat{f_t}}} = \sigma(\hat{\mathbf{f_t}} \mathbf{W_t} + \mathbf{b_t}); \quad \alpha_{\mathbf{t}} = \text{softmax}(\mathbf{h_{\hat{f_t}} W_{\alpha_t}} + \mathbf{b_{\alpha_t}}) \quad (5)$$

$$\hat{\mathbf{f_{ei}}} = \alpha_{\mathbf{t}} \times \mathbf{f_{ei}} \in \mathbb{R}^{m \times 768} \quad (6)$$

where $\mathbf{W_{ei}}$, $\mathbf{W_t} \in \mathbb{R}^{768 \times 768}$, $\mathbf{W_{\alpha_{ei}}}$, and $\mathbf{W_{\alpha_t}} \in \mathbb{R}^{768 \times 1}$ are the weights for transforming emotion-aware image feature ($\mathbf{f_{ei}} \in \mathbb{R}^{m \times 768}$ repeated $n$ times to account for $n$ textual tokens), updated text feature ($\hat{\mathbf{f_t}} \in \mathbb{R}^{m \times 768}$ repeated $m$ times to account for $m$ image patches), intermediate representation of transformed emotion-aware image feature ($\mathbf{h_{f_{ei}}} \in \mathbb{R}^{m \times n \times 768}$), and intermediate representation of transformed text feature ($\mathbf{h_{\hat{f_t}}} \in \mathbb{R}^{m \times m \times 768}$), respectively. The bias terms, $\mathbf{b_{ei}}, \mathbf{b_t} \in \mathbb{R}^{768}$, and $\mathbf{b_{\alpha_{ei}}}, \mathbf{b_{\alpha_t}} \in \mathbb{R}^1$, correspond to the representation learning for $\mathbf{h_{f_{ei}}}, \mathbf{h_{\hat{f_t}}}, \alpha_{\mathbf{ei}}$, and $\alpha_{\mathbf{t}}$, respectively.

---

6. Broadcasting Semantics — PyTorch ↗

TABLE 3: Comparison of ALFRED and baselines on the MOOD dataset. The last row shows the improvement of ALFRED ($^*$) over the *early-fusion* based model, designated as the best baseline ($\dagger$). Class-wise accuracy for the six *Ekman* emotions in MOOD is also reported.

| Modality | Model | Acc. | Prec. | Rec. | F1 | FER | AGR | JOY | SDN | SPR | DGT |
|----------|-------|------|-------|------|-----|-----|-----|-----|-----|-----|-----|
| UM | BERT | 0.633 | 0.6537 | 0.6337 | 0.6387 | 0.4573 | 0.5587 | **0.7969** | 0.7484 | 0.4559 | 0.6869 |
| | ViT | 0.6713 | 0.6913 | 0.6713 | 0.6738 | **0.9065** | 0.5507 | 0.6884 | 0.6243 | 0.8766 | 0.7101 |
| MM | Early-fusion$^\dagger$ | 0.7836 | 0.8121 | 0.7836 | 0.7749 | 0.8991 | 0.7594 | 0.7405 | 0.657 | 0.8335 | 0.833 |
| | MMBT | 0.6337 | 0.6537 | 0.633 | 0.6352 | 0.581 | **0.7848** | 0.7818 | 0.6534 | 0.5252 | 0.7973 |
| | CLIP | 0.6378 | 0.8027 | 0.6378 | 0.6816 | 0.5351 | 0.5011 | 0.7797 | **0.764** | 0.5226 | 0.72 |
| | VisualBERT | 0.6725 | 0.7961 | 0.6725 | 0.7002 | 0.7075 | 0.5838 | 0.6969 | 0.7294 | 0.7787 | 0.5513 |
| | ALFRED$^\star$ | **0.8239** | **0.8314** | **0.8239** | **0.8243** | 0.8777 | 0.7835 | 0.7924 | 0.7625 | **0.8935** | **0.8392** |
| $\Delta_{(\star\text{-}\dagger)\times 100}(\%)$ | | ↑ 4.03% | ↑ 1.93% | ↑ 4.03% | ↑ 4.94% | ↓ 2.14% | ↑ 2.41% | ↑ 5.19% | ↑ 10.55% | ↑ 6.00% | ↑ 0.62% |

## 4.5 Prediction and Training Objective

Finally, we apply sum-pooling across the *first* dimension of the corresponding weight-aggregated features: $\hat{\mathbf{f}_{ei}} \in \mathbb{R}^{m \times 768}$ and $\hat{\mathbf{f}_t} \in \mathbb{R}^{m \times 768}$, followed by concatenating the sum-pooled features ($\hat{\mathbf{f}_{ei}} \in \mathbb{R}^{768}$ and $\hat{\mathbf{f}_t} \in \mathbb{R}^{768}$), to produce a joint meme representation ($\mathbf{f}_{\mathbf{z_1}} \in \mathbb{R}^{1536}$). This is given as input to a feed-forward network for the final classification.

$$\mathbf{f}_{\mathbf{z_1}} = [sum(\hat{\mathbf{f}_{ei}}), \mathbf{sum}(\hat{\mathbf{f}_t})] \in \mathbb{R}^{\mathbf{1536}} \quad (7)$$

$$\mathbf{h}_{\mathbf{f}_{\mathbf{z_1}}} = relu(\mathbf{f}_{\mathbf{z_1}} \mathbf{W}_{\mathbf{z_1}} + \mathbf{b}_{\mathbf{z_1}}) \in \mathbb{R}^{768} \quad (8)$$

$$\mathbf{h}_{\mathbf{f}_{\mathbf{z_2}}} = softmax(\mathbf{h}_{\mathbf{f}_{\mathbf{z_1}}} \mathbf{W}_{\mathbf{z_2}} + \mathbf{b}_{\mathbf{z_2}}) \in \mathbb{R}^{6} \quad (9)$$

$$\hat{\mathbf{y}} = argmax_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{y}; \theta, \mathbf{h}_{\mathbf{f}_{\mathbf{z_2}}}) \quad (10)$$

where $\mathbf{W}_{\mathbf{z_1}} \in \mathbb{R}^{1536 \times 768}$, $\mathbf{W}_{\mathbf{z_2}} \in \mathbb{R}^{768 \times 6}$, $\theta$ represents model parameters, $\hat{\mathbf{y}}$ is the predicted class index, and $\mathcal{Y}$ is the label-id set, with $|\mathcal{Y}| = 6$. The bias terms, $\mathbf{b}_{\mathbf{z_1}} \in \mathbb{R}^{768}$ and $\mathbf{b}_{\mathbf{z_2}} \in \mathbb{R}^{6}$, are corresponding to the last two layers. We use the cross-entropy loss for optimization. Moreover, we employ the online label smoothing [63] for regularization.

## 5 BASELINE MODELS

**Unimodal Baselines:** We restrict modality-specific encoders to the following choices, as the primary objective of this work is to investigate an optimal multimodal fusion strategy.

- **BERT [44]:** We use the BERT-base-uncased model as our text-only baseline.
- **ViT [51]:** We use ViT with ImageNet weights as our image-only baseline.

**Multimodal Baselines:** For multimodal systems, we explore the following competent approaches as comparative baselines. These systems endorse various multimodal interaction schemes, facilitating a robust assessment.

- **Early-fusion:** In this model, the features from ViT and BERT are concatenated and passed through a feed-forward network for classification.
- **MMBT [64]:** It is a supervised bi-modal transformer that projects image features from unimodally pre-trained image encoders to text tokens.
- **CLIP [65]:** CLIP is a contrastive learning-based approach that is designed to learn visual information through natural language supervision.
- **VisualBERT [60]:** VisualBERT is a transformer-based model for visuo-lingual modelling. It has been trained on the MS COCO dataset employing masked language modeling and the sentence-image prediction objective functions.

TABLE 4: Hyperparameters of different models.

| Model | BS | Ep | LR | Image Encoder | Text Encoder | # Param |
|-------|----|----|-----|---------------|--------------|---------|
| TextBERT | 64 | 20 | 0.0001 | - | bert-base-uncased | 110M |
| ViT | 32 | 20 | 0.0001 | vit-base | - | 86M |
| Early-fusion | 32 | 20 | 0.0001 | vit-base | bert-base-uncased | 196M |
| CLIP | 43 | 20 | 0.0001 | vit-base | bert-base-uncased | 151M |
| MMBT | 32 | 30 | 0.00001 | resnet152 | bert-base-uncased | 169M |
| V-BERT | 32 | 30 | 0.000001 | Faster RCNN | bert-base-uncased | 247M |
| ALFRED | 32 | 30 | 0.0001 | vit-base | bert-base-uncased | 282M |

## 6 EXPERIMENTS

Firstly, we compare ALFRED with both unimodal (image/text) and multimodal models. The comparisons are first made for meme emotion detection task using the MOOD dataset, followed by evaluation on three Memotion tasks [2] (c.f. Section 6.3). We further perform the ablation of our model. This is followed by examining the interpretability of ALFRED using GradCAM [66]. We then demonstrate the generalizability of ALFRED's performance on the HarMeme [11] and Dank Memes datasets [12]. Finally, we analyze the errors observed during performance evaluation. Our experimental setup's empirical examinations involve fine-tuning for the respective tasks and datasets. Since we primarily explore a multi-class classification setup, we are more interested in evaluating the system performances that factors-in the class-wise contributions equally. Therefore, we use macro-averaged formulations of accuracy, precision, recall, and F1-score as evaluation metrics. We also report class-wise F1

## 6.1 Implementation Details and Hyperparameter Values

We train all the models using Pytorch 1.10 on an Nvidia Tesla V100 GPU with 32 GB dedicated memory, CUDA-11.2 and cuDNN-8.1.1 installed. For the primary emotion classification, Memotion, HarMeme tasks, we use BERT as the text encoder and ViT as the image encoder. Specifically, we use the bert-base-uncased checkpoint for BERT and google/vit-base-patch16-224 checkpoint for ViT. However, for the Dank Memes task, we switch BERT with UmBERTo, which is a BERT-based model but pre-trained using Italian corpus. The linear layers in GCA and GMF modules are initialized using Xavier initialization, and the bias is set to zero. For the meme emotion classification task, we train all the models using the *online label smoothing* loss and Adam optimizer. For the Memotion, HarMeme and Dank Memes tasks, we use cross-entropy loss and Adam optimizer. We also present these details in Table 4.

## 6.2 Evaluation on the MOOD dataset

Among unimodal models, the image-only model is observed to perform better (c.f. Table 3) than the text-only model by 4% F1 score. Also, multimodal baselines are observed to perform either

TABLE 5: Task-wise and class-wise performance (Macro-F1) of different approaches on the `Memotion` tasks. The last row shows the improvement of ALFRED (*) over the *previous best* results (†) reported.

| Modality | Model | SENT | EMOT | | | | | EMOT-Q | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Sentiment | Humour | Sarcasm | Offensive | Motivation | **Average** | Humour | Sarcasm | Offense | Motivation | **Average** |
| UM | BERT | 0.3123 | 0.5235 | 0.4747 | 0.5012 | 0.5089 | 0.5021 | 0.2487 | 0.2309 | 0.2333 | 0.4641 | 0.2942 |
| | ViT | 0.3158 | 0.5114 | 0.4851 | 0.5119 | 0.519 | 0.5069 | 0.2434 | 0.249 | 0.2388 | 0.4972 | 0.307 |
| MM | Early-fusion | 0.3295 | 0.5122 | 0.5032 | 0.5059 | 0.4591 | 0.4951 | 0.2368 | 0.2426 | 0.235 | 0.4481 | 0.2906 |
| | MMBT | 0.3457 | **0.5393** | 0.5015 | 0.4970 | 0.4989 | 0.5092 | 0.2474 | 0.2364 | 0.2475 | 0.4838 | 0.3038 |
| | CLIP | 0.3261 | 0.4798 | 0.5133 | 0.5115 | 0.4967 | 0.5003 | 0.2652 | **0.2549** | 0.2448 | 0.4754 | 0.3101 |
| | Previous Best† | **0.3547** | 0.51587 | 0.5159 | **0.5225** | 0.51909 | 0.5183 | **0.27069** | 0.25028 | 0.25761 | 0.51126 | 0.3225 |
| | ALFRED* | 0.3486 | 0.5268 | **0.5272** | 0.5175 | **0.5375** | **0.5272** | 0.2646 | 0.2501 | **0.2598** | **0.5275** | **0.3253** |
| $\Delta_{(\star\text{-}\dagger)\times100}(\%)$ | | ↓ 0.61% | ↑ 1.09% | ↑ 1.13% | ↓ 0.5% | ↑ 1.84% | ↑ 0.89% | ↓ 0.6% | ↓ 0.01% | ↑ 0.22% | ↑ 1.62% | ↑ 0.28% |

at par or better than unimodal models. One of the top-performing baselines, as shown in Table 3, is the early-fusion model with BERT and ViT as its text and image encoders, respectively. This could be due to an optimal modeling requirement posed by the meme emotion detection task, which does not seem to favor complex co-attentive visual-linguistic grounding employed by models like MMBT, CLIP, and VisualBERT.

The early-fusion model (0.7749) yields a $7\%$ absolute improvement in F1-score over the sophisticated VisualBERT (0.7002). Overall, both perform better than the multimodal baselines like MMBT (0.6352) and CLIP (0.6816). In comparison, ALFRED registers $4.94\%$ F1 improvement over the early-fusion model. This improvement could be mainly attributed to the GMF-based explicit emotion modeling and GCA-based inter-modal fusion, facilitating preferential treatment for both input modalities conditioned upon emotion-enriched visual cues. Overall, ALFRED yields an improvement of 1.93%-4.94% across all four metrics.

ALFRED significantly increases accuracy for four classes – *anger* (↑2.41%), *joy* (↑5.19%), *sadness* (↑10.55%) and *surprise* (↑6.00%). In contrast, the accuracy for *disgust* improves slightly (0.62%), but not as much as for the classes above. This subtle enhancement observed could be likely due to the expressiveness of the emotion *disgust* via either text, image, or even both (See Fig. 2 (a)). Whereas the lower representation of the *disgust* class in the dataset explains the improvement that is minor compared to that of the categories mentioned above. Moreover, the performance for *fear* drops by $2\%$, wherein the discriminatory cues are predominantly image-based, the implication of which is also corroborated by the highest category-specific performance ($\approx 0.91$ F1-score), by ViT-based model. Besides posing challenges like object occlusion, complex pose, image quality, etc., the visual modeling is impacted by the category's under-representation in both `MOOD` and extended AffectNet datasets, leading to ALFRED's drop in accuracy for *fear* as against the enhancement observed for other categories.

### 6.3 Evaluation on the `Memotion` dataset

We then compare the performance of ALFRED, other baselines, and the state-of-the-art systems on the `Memotion` shared task [2]. The `Memotion` dataset contains approximately $8K$ memes. It was proposed for the three subtasks[7] – *sentiment analysis* (positive/negative), *emotion classification* (humour/sarcasm/offense/motivational), and *emotion class quantification* (slightly/mildly/very). The original average baseline F1 for

7. We use abbreviations SENT, EMOT and EMOT-Q for *sentiment analysis*, *emotion classification*, and *emotion class quantification*, respectively.

the three `Memotion` sub-tasks – SENT (0.2176), EMOT (0.5002), and EMOT-Q (0.3009) – indicate inherent non-triviality of the tasks. The previous best systems involve a word2vec [67], [68] based feed-forward neural network for SENT [69], a multimodal multi-tasking based setup for EMOT [70], and a feature-based ensembling approach for the EMOT-Q task [71]. The performance of the state-of-the-art systems for the three tasks (c.f. Table 5) are 0.3547, 0.5183, and 0.3225 F1, respectively. In comparison, ALFRED induces an increase of $0.89\%$ and $0.28\%$ F1-score for EMOT and EMOT-Q tasks, respectively. For SENT, however, ALFRED lags slightly behind the *previous best* score by 0.61% F1 score.

ALFRED's low score on the SENT task could be due to noise induced by the emotion-enriched feature that might complicate modeling a more straightforward task like SENT compared to simpler early-fusion-based state-of-the-art. ALFRED also performs relatively better with $1.09\% - 1.13\%$ F1-score increment in the *humour* and *sarcasm* categories for the EMOT task, as against that for EMOT-Q. Since the level of abstraction for the information being modeled for EMOT (emotion classification) is relatively higher as compared to that for EMOT-Q (emotion quantification), an explicit emotion modeling could help detect emotions for EMOT, and not necessarily for fine-grained emotion intensity quantification in EMOT-Q, especially for complex categories like *humor* and *sarcasm*. On average, ALFRED's performance on `Memotion` tasks are comparable – it reports better scores for the EMOT and EMOT-Q tasks; however, it yields inferior performance in the SENT task.

### 6.4 Ablation Studies

The incorporation of emotion features induces 5% improvement over ALFRED without emotion features. Further, replacing the GMF module with simple concatenation for incorporating emotion features causes a 2% performance drop. On the other hand, GCA is also observed to be pivotal as its replacement with dense

TABLE 6: Ablation study for ALFRED on `MOOD`; emotion features (EMO), GMF and GCA (via DCA: Dense Co-attention) modules.

| Approach | F1 | Acc |
|---|---|---|
| ALFRED | 82.43 | 82.39 |
| − EMO | 77.19 | 77.08 |
| − GMF | 80.60 | 80.38 |
| − GCA + DCA | 80.04 | 79.97 |

co-attention (DCA) [72] induces a drop of 2% performance. Effectively, the exclusion of GMF and GCA from ALFRED is empirically observed to induce a performance drop of $\approx 2\%$, as shown in Table 6.

We analyze the contribution of each module of `ALFRED` on `MOOD`: emotion encoder, gated multimodal fusion (GMF), and gated cross attention (GCA) in Table. 6. `ALFRED` without EMO is expected to perform worse due to the complexity of the model not being complemented by the required rich features. This corroborates the requirement of a solution that explicitly incorporates emotion-enriched feature modeling. The pre-training of the emotion encoder is discussed in detail in Section 4.2. On the other hand, the early fusion model, being efficient yet straightforward towards multimodal classification, yields impressive results, which we also consider the best baseline for comparison. Also, besides discussing the effect of `ALFRED` without EMO (c.f. Table 6), we specifically include evaluations without GMF and GCA to examine the optimal modeling of emotion features using these modules in `ALFRED`. For both exclusions, the performance is low. This assessment consolidates the boosting capacity of each module constituting `ALFRED`.

## 6.5 Discussion

**– Novelty Aspects of `ALFRED`:** We empirically establish the efficacy of using low-rank bilinear pooling-based non-linear gating fusion instead of simple concatenation for fusing emotion-aware image representations in GMF. This helps characterize intra-modal fusion against inter-modal fusion, for which concatenation has been the conventional fusion strategy. For the GCA module, we first perform the conditioning based upon the emotion-aware *image* representation to compute the *textual* attention and feed it back towards computing a final emotion-aware image representation. This is in contrast to the convention of performing either a text-based conditioning [62] or a parallel co-attention based strategy [59]. To our understanding, this is the first attempt to emphasize the visual cues toward overall modeling.

Essentially, through our proposed approach incorporating the adaptation of existing effective techniques like GMF and GCA, we present a strategy that has been observed to be empirically adequate for modeling intra-modal and inter-modal fusion. To the best of our understanding, incorporating emotion-oriented features explicitly via *visual* modality toward a task like a meme analysis has not been explored prior to this work.

**– Dataset Utility:** Meme datasets mostly encapsulate affective dimensions representing higher levels of abstraction ranging from categories like humor, sarcasm, offense, and motivation in the case of Memotion, to aspects like Harmfulness and Hatespeech within memes in `HarMeme` and `Dank Memes` datasets, respectively. Additionally, since memes in such datasets usually capture real-world events involving famous personalities and phenomena, they tend to be reasonably restricted in terms of the visual subjects they embody. For instance, a significant portion of such memes does not contribute visually towards affective adjudication, limiting the characterization due to cartoons, caricatures, and expressive personifications, for meme emotion detection. Most of them typically end up projecting textual cues as their characteristic feature.

In contrast, `MOOD` captures the affective dimension that objectively focuses on six *basic Ekman* emotions via multimodal cues for *generic* themes, capturing the expressivity differently from other datasets and soliciting an appropriate investigative framework. Such memetic configuration aptly represents the scope of this work – detecting basic emotions from *generic memes*, which tend to bear emotional expressivity via both image and text modalities.

## 6.6 Interpretability

We attempt to interpret the decisions made by `ALFRED` using GradCam [66]. It uses gradients flowing through a model to produce a rough attention map. This highlights the regions in the image that the model pays attention to while making a decision.

In Fig. 9(a), the text alone is not sufficient to detect the emotion of the meme, as the excerpt 'HAPPY NEW YEARS!...' could mislead the model's decision. The key aspect here lies with the *angry* ('grumpy cat meme') expression of the cat. In Fig. 9(b), we notice that the model pays attention to the cat's face to correctly predict that the emotion is *anger*.

## 6.7 Error Analysis

**– Visual Obscurity:** On analyzing the incorrect predictions from the test set, we find that most misclassifications involve complex memetic text or obscure visuals, including prominent visual occlusion. Another distinct trend for poor results is observed for the samples belonging to the *least represented* emotion categories *fear* and *disgust*, with 2.5 % decrease and a marginal enhancement of 0.62 % respectively, in the accuracy values compared with those from an *early-fusion* based model. An example that demonstrates the misclassification attempt of `ALFRED` due to both the aforementioned likely reasons is shown in Fig. 9(c), depicting the facial expressions of the man to be that of *fear*. Still, the model incorrectly predicts it as *surprise*. On interpreting the visual attention-map via GradCAM-based visualization, we observe that `ALFRED` does not pay attention to the subtle facial expressions indicating *fear*, as demonstrated by the misplaced visual attention, in Fig. 9(d).

**– Textual Obscurity:** Data sufficiency is also observed to play pivotal role towards class-wise performances observed in Table 3 w.r.t. *textual influence* within memes. A critical evidence corroborating this aspect is the lexical richness of memetic text for categories *disgust* and *fear*, as against lexical obscurity for *surprise*, as observed from Figs. 5(f), 5(a) and 5(e), respectively. The former two, being sparsely represented within `MOOD`, yield subpar performances for the corresponding categories, whereas the latter being densely represented, contributes a decent accuracy (c.f. Table 3), despite lexical obscurity. This suggests that multimodal (as against unimodal) contextual dependency is imperative towards emotion recognition from memes since it is not just the complex cross-modal interplay that encapsulates the intended message but the modality-specific intricacies that constitute complex memetic designs.

**– Analyzing Thematic Overlaps:** We further investigate the semantic complexity posed by the themes that various memes are based on. To this end, thematic structure and characteristics are derived via a clustering-based approach and are ascertained for semantic overlap [73] w.r.t. the six *Ekman* emotions for `MOOD`. Firstly, document embeddings are obtained via `all-MiniLM-L6-v2` based Transformer model [74], followed UMAP-based dimensionality reduction [75] and HDBSCAN-based clustering [76].

The hierarchical thematic sub-groupings obtained through HDBSCAN, analyzed for semantic similarity using similarity matrices (shown in Fig. 8), reveal distinct overlaps and proximity regarding six of Ekman's emotions. These patterns are illustrated through highlighted examples. Notably, *disgust* memes (topic id: 68, 75) represented by patterns #1 and #2 in Fig. 8 show variable overlaps with *surprise* (topic id: 37) and *joy* (topic id: 29, 30).
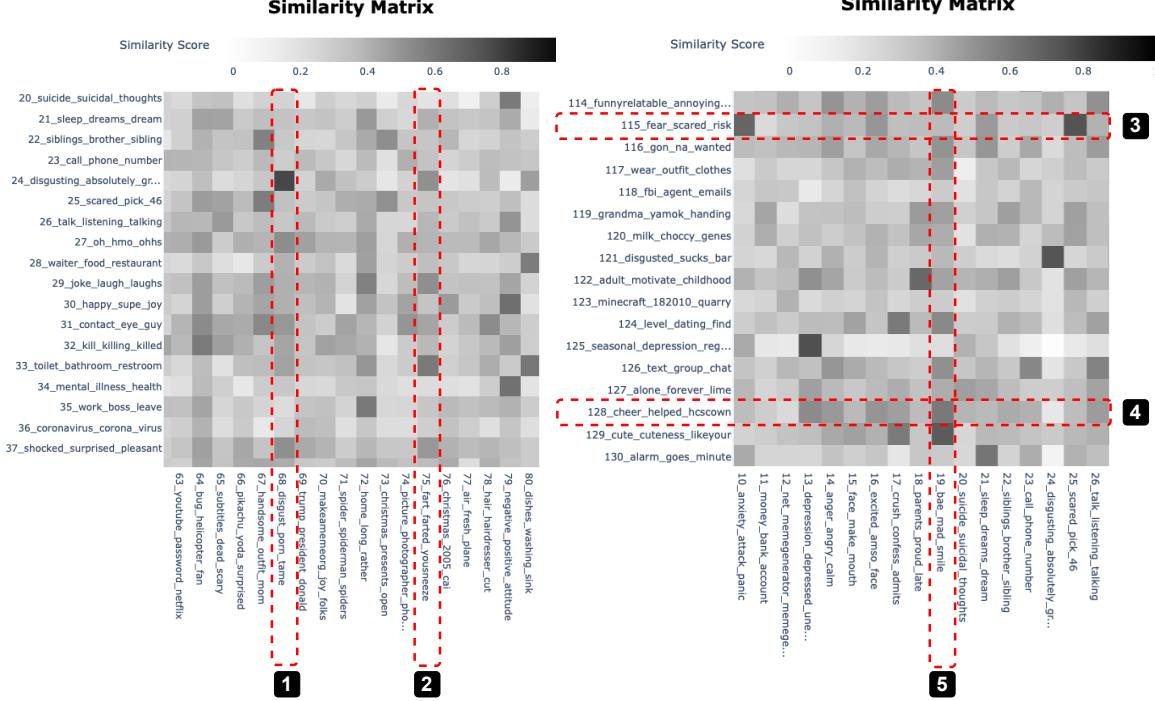
Fig. 8: Thematic overlap analysis via similarity matrix, with *five* example cases of inter-emotion overlap highlighted. Each x/y-axes label represents *topic id*, followed by a set of *three* corresponding *representative key-words* (`topicid_kw1_kw2_kw3`).



Fig. 9: Depiction of: Interpretability Analysis [subfigs. (a) and (b) – `ALFRED` correctly predicts an *anger* meme], `Error Analysis` [subfigs. (c) and (d) – `ALFRED` incorrectly predicts a *fear* meme as a *surprise* meme].

The second similarity matrix in Fig. 8 highlights patterns connecting *fear* (topic id: 115) with *anger* (topic id: 14) and *joy* (topic id: 16) (pattern #3), and *joy* (topic id: 128) with *sadness* (topic id: 13) and *anger* (topic id: 14) (pattern #4). Additionally, *anger* (topic id: 19) overlaps with *sadness* (topic id: 127) and others, along with *joy* (topic id: 128, 129) in pattern #5.

Overall, *joy* and *sadness* consistently emerge as common factors in emotion overlaps, aligning with Ekman's [10] and Plutchik's [77] theories. This suggests that the proximity of these emotions in memes stems from complex linguistic content. Determining the exact valence of this content remains challenging, underlining the need for detailed emotion analysis of memes.

## 6.8 Generalizability

Here, we establish the generalizability of `ALFRED` for `HarMeme` [11] and `Dank Memes` tasks [12]. The dataset for `HarMeme` constitutes $\approx 7K$ memes (in English) on Covid-19 and US Politics. This dataset captures annotations for harmfulness and the targeted entity types. The second dataset, `Dank Memes`, comprises $\approx 1K$ hateful memes (in Italian). The memes are about the 2019 Italian Government Crisis. There was an associated shared task involving three subtasks – a) meme detection, b) hate-speech identification, and c) event clustering. In this work, we focus on hate-speech identification to ensure evaluation consistency.

**– `HarMeme`:** The best performance on this dataset was reported by MOMENTA [26] which strongly outperformed sophisticated multimodal baselines such as V-BERT and ViLBERT. For two-class classification, `ALFRED` is observed to achieve an improvement of 3.08% and 1.8% F1 over MOMENTA, respectively, on the Harm-C and Harm-P datasets. For three-class classification, `ALFRED` achieves 6.43% and 23.86% F1 increment over MOMENTA on the Harm-C and Harm-P datasets, respectively (c.f. Fig. 10).

**– `Dank Memes`:** Dank Memes is an Italian hateful politics meme dataset. The top two submissions for the related shared-task were both early-fusion based: *Unitor* employs domain-specific pretraining before finetuning on Dank Memes; *UPB* uses VGCN-BERT for text modality [12]. Since the task deals with memes with embedded Italian content, we replace the BERT model with UmBERTo [78] within `ALFRED` while keeping other components same. `ALFRED` achieves an absolute increment of 2.02% and 4.37% in F1 and precision, respectively, over the best baseline, while the recall lags by 1.96% behind *Unitor*. These results not only highlight `ALFRED`'s generalizability, but also indicate its language-agnostic cross-lingual affinity (c.f. Fig. 10), especially w.r.t multimodal tasks like meme analysis.

## 7 CONCLUSION

In this work, we first introduced `MOOD`, a new dataset for detecting emotions in Internet memes. We then proposed `ALFRED`, which
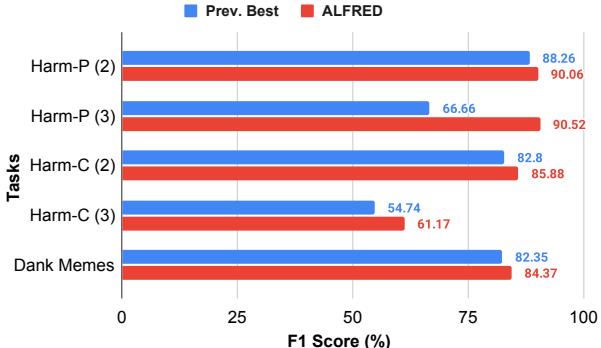
Fig. 10: Performance comparison for `ALFRED` and *previous best* on `HarMeme` (US Politics (P)/Covid (C); 2/3 class classification) and `Dank Memes` tasks, demonstrating `ALFRED`'s generalizability.

uses emotion-aware meme representations to detect emotions from memes. Extensive experiments indicated that `ALFRED` outperforms strong multimodal baseline with 4.94% F1 increment and yields robust performance on the `Memotion` task [2] dataset. Further, we investigated the interpretability of the model by establishing the correspondences between the correct emotion class being predicted and the expressive emotions being attended to within the meme image. We also highlighted the inherent limitations that explicit emotion modeling can develop. Finally, we established the generalizability of `ALFRED` by demonstrating its superiority over *previous best* baselines on the `HarMeme` and `Dank Memes` datasets. As part of the future extension to this work, we would like to explore a multi-task learning setup involving the detection of correlated fine and coarse-grained emotion features for memes. Moreover, tasks like explanation generation for various meme-emotions and network structure-based investigation of meme virality, w.r.t the emotions, are also promising avenues to explore.

## ACKNOWLEDGMENT

## REFERENCES

[1] A. Mogadala, M. Kalimuthu, and D. Klakow, "Trends in integration of vision and language research: A survey of tasks, datasets, and methods," *JAIR*, vol. 71, p. 1183–1317, Aug 2021.

[2] C. Sharma, D. Bhageria, W. Scott, S. PYKL, A. Das, T. Chakraborty, V. Pulabaigari, and B. Gambäck, "Semeval-2020 task 8: Memotion analysis–the visuo-lingual metaphor!" *arXiv preprint arXiv:2008.03781*, 2020.

[3] S. Sharma, S. Agarwal, T. Suresh, P. Nakov, M. S. Akhtar, and T. Chakraborty, "What do you meme? generating explanations for visual semantic role labelling in memes," in *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, B. Williams, Y. Chen, and J. Neville, Eds. AAAI Press, 2023, pp. 9763–9771. [Online]. Available: https://doi.org/10.1609/aaai.v37i8.26166

[4] S. Sharma, R. S, U. Arora, M. S. Akhtar, and T. Chakraborty, "MEMEX: detecting explanatory evidence for memes via knowledge-enriched contextualization," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 5272–5290. [Online]. Available: https://doi.org/10.18653/v1/2023.acl-long.289

[5] S. Sharma, M. K. Siddiqui, M. S. Akhtar, and T. Chakraborty, "Domain-aware self-supervised pre-training for label-efficient meme analysis," in *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing, AACL/IJCNLP 2022 - Volume 1: Long Papers, Online Only, November 20-23, 2022*, Y. He, H. Ji, Y. Liu, S. Li, C. Chang, S. Poria, C. Lin, W. L. Buntine, M. Liakata, H. Yan, Z. Yan, S. Ruder, X. Wan, M. Arana-Catania, Z. Wei, H. Huang, J. Wu, M. Day, P. Liu, and R. Xu, Eds. Association for Computational Linguistics, 2022, pp. 792–805. [Online]. Available: https://aclanthology.org/2022.aacl-main.60

[6] S. Pramanick, S. Sharma, D. Dimitrov, M. S. Akhtar, P. Nakov, and T. Chakraborty, "MOMENTA: A multimodal framework for detecting harmful memes and their targets," in *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, M. Moens, X. Huang, L. Specia, and S. W. Yih, Eds. Association for Computational Linguistics, 2021, pp. 4439–4455. [Online]. Available: https://doi.org/10.18653/v1/2021.findings-emnlp.379

[7] S. Zhao, G. Ding, Q. Huang, T.-S. Chua, B. W. Schuller, and K. Keutzer, "Affective image content analysis: A comprehensive survey," in *IJCAI-18*, 2018, pp. 5534–5541.

[8] P. Patwa, S. Ramamoorthy, N. Gunti, S. Mishra, S. S, A. Reganti, A. Das, T. Chakraborty, A. Sheth, A. Ekbal, and C. Ahuja, "Findings of memotion 2: Sentiment and emotion analysis of memes," *Proceedings of De-Factify: Workshop on Multimodal Fact Checking and Hate Speech Detection, collocated with AAAI'22*, 02 2022.

[9] H. Alves, C. Fernandes, and M. Raposo, "Social media marketing: A literature review and implications: Implications of social media marketing," *Psychology And Marketing*, vol. 33, pp. 1029–1038, 12 2016.

[10] P. Ekman and D. Cordaro, "What is meant by calling emotions basic," *Emotion Review*, vol. 3, no. 4, pp. 364–370, 2011.

[11] S. Pramanick, D. Dimitrov, R. Mukherjee, S. Sharma, M. S. Akhtar, P. Nakov, and T. Chakraborty, "Detecting harmful memes and their targets," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Online: ACL, Aug. 2021, pp. 2783–2796.

[12] M. Miliani, G. Giorgi, I. Rama, G. Anselmi, and G. Lebani, *DANKMEMES @ EVALITA 2020: The Memeing of Life: Memes, Multimodality and Politics*, 01 2020, pp. 275–283.

[13] M. Abdul-Mageed and L. Ungar, "EmoNet: Fine-grained emotion detection with gated recurrent neural networks," in *ACL*, Jul. 2017, pp. 718–728.

[14] A. Mollahosseini, B. Hasani, and M. H. Mahoor, "Affectnet: A database for facial expression, valence, and arousal computing in the wild," *IEEE TAC*, vol. 10, p. 18–31, 2019.

[15] S. Poria, I. Chaturvedi, E. Cambria, and A. Hussain, "Convolutional mkl based multimodal emotion recognition and sentiment analysis," in *ICDM*, 2016, pp. 439–448.

[16] C. T. Duong, R. Lebret, and K. Aberer, "Multimodal classification for analysing social media," *arXiv preprint arXiv:1708.02099*, 2017.

[17] A. Hu and S. Flaxman, "Multimodal sentiment analysis to explore the structure of emotions," *SIGKDD*, Jul 2018.

[18] J. Qiu, X. Li, and K. Hu, "Correlated attention networks for multimodal emotion recognition," in *BIBM*, 2018, pp. 2656–2660.

[19] J. Huang, J. Tao, B. Liu, Z. Lian, and M. Niu, "Multimodal transformer fusion for continuous emotion recognition," in *ICASSP*, 2020, pp. 3507–3511.

[20] S. Siriwardhana, T. Kaluarachchi, M. Billinghurst, and S. Nanayakkara, "Multimodal emotion recognition with transformer-based self supervised feature fusion," *IEEE Access*, vol. 8, pp. 176274–176285, 2020.

[21] T. Mittal, P. Guhan, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha, "Emoticon: Context-aware multimodal emotion recognition using frege's principle," 2020.

[22] S. Zannettou, T. Caulfield, J. Blackburn, E. De Cristofaro, M. Sirivianos, G. Stringhini, and G. Suarez-Tangil, "On the origins of memes by means of fringe web communities," in *IMC '18*, New York, NY, USA, 2018, p. 188–202.

[23] S. Suryawanshi, B. R. Chakravarthi, M. Arcan, and P. Buitelaar, "Multimodal meme dataset (MultiOFF) for identifying offensive content in image and text," in *Proceedings of the Second Workshop on Troll., Agg. and Cyb.*, May 2020, pp. 32–41.

[24] D. Kiela, H. Firooz, A. Mohan, V. Goswami, A. Singh, P. Ringshia, and D. Testuggine, "The hateful memes challenge: Detecting hate speech in multimodal memes," in *NeurIPS '20*, vol. 33, 2020.

[25] R. Gomez, J. Gibert, L. Gomez, and D. Karatzas, "Exploring hate speech detection in multimodal publications," in *WACV '20*, 2020, pp. 1459–1467.

[26] S. Pramanick, S. Sharma, D. Dimitrov, M. S. Akhtar, P. Nakov, and T. Chakraborty, "MOMENTA: A multimodal framework for detecting harmful memes and their targets," in *Findings of EMNLP 2021*, Punta Cana, Dominican Republic, Nov. 2021, pp. 4439–4455.

[27] J. Drakett, B. Rickett, K. Day, and K. Milnes, "Old jokes, new media – online sexism and constructions of gender in internet memes," *Fem. & Psy.*, vol. 28, no. 1, pp. 109–127, 2018.

[28] A. Williams, "Black memes matter: #livingwhileblack with Becky and Karen," *Social Media + Society*, vol. 6, no. 4, 2020.

[29] T. Askanius, "On frogs, monkeys, and execution memes: Exploring the humor-hate nexus at the intersection of neo-Nazi and alt-right movements in Sweden," *Tel. & New Media*, vol. 22, no. 2, pp. 147–165, 2021.

[30] S. Sharma, T. Suresh, A. Kulkarni, H. Mathur, P. Nakov, M. S. Akhtar, and T. Chakraborty, "Findings of the CONSTRAINT 2022 shared task on detecting the hero, the villain, and the victim in memes," in *CON-STRAINT*, 2022, pp. 1–11.

[31] S. Masud, S. Dutta, S. Makkar, C. Jain, V. Goyal, A. Das, and T. Chakraborty, "Hate is the new infodemic: A topic-aware modeling of hate speech diffusion on twitter," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*.  IEEE, 2021, pp. 504–515.

[32] T. Garg, S. Masud, T. Suresh, and T. Chakraborty, "Handling bias in toxic speech detection: A survey," *ACM Computing Surveys*, vol. 55, no. 13s, pp. 1–32, 2023.

[33] S. Masud, M. Bedi, M. A. Khan, M. S. Akhtar, and T. Chakraborty, "Proactively reducing the hate intensity of online posts via hate speech normalization," in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 3524–3534.

[34] A. Kulkarni, S. Masud, V. Goyal, and T. Chakraborty, "Revisiting hate speech benchmarks: From data curation to system deployment," in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023, pp. 4333–4345.

[35] M. S. Hee, S. Sharma, R. Cao, P. Nandi, P. Nakov, T. Chakraborty, and R. K. Lee, "Recent advances in hate speech moderation: Multimodality and the role of large models," *CoRR*, vol. abs/2401.16727, 2024. [Online]. Available: https://doi.org/10.48550/arXiv.2401.16727

[36] A. Das, J. S. Wahi, and S. Li, "Detecting hate speech in multi-modal memes," 2020.

[37] V. Sandulescu, "Detecting hateful memes using a multimodal deep ensemble," *arXiv:2012.13235*, 2020.

[38] X. Zhou, J. Wu, and R. Zafarani, "SAFE: Similarity-aware multi-modal fake news detection," in *PAKDD*.  Springer, 2020, pp. 354–367.

[39] P. Lippe, N. Holla, S. Chandra, S. Rajamanickam, G. Antoniou, E. Shutova, and H. Yannakoudakis, "A multimodal framework for the detection of hateful memes," *arXiv:2012.12871*, 2020.

[40] L. Shang, C. Youn, Y. Zha, Y. Zhang, and D. Wang, "KnowMeme: A knowledge-enriched graph neural network solution to offensive meme detection," in *eScience '21*, 2021, pp. 186–195.

[41] K. Karkkainen and J. Joo, "Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation," in *WACV '21*, 2021, pp. 1548–1558.

[42] S. Sharma, F. Alam, M. S. Akhtar, D. Dimitrov, G. Da San Martino, H. Firooz, A. Halevy, F. Silvestri, P. Nakov, and T. Chakraborty, "Detecting and understanding harmful memes: A survey," in *IJCAI-ECAI '22*, Vienna, Austria, 2022.

[43] J. H. French, "Image-based memes as sentiment predictors," *i-Society*, pp. 80–85, 2017.

[44] J. Devlin, M-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *ACL-HLT*, 2019, pp. 4171–4186.

[45] J. Pennington, R. Socher, and C. Manning, "GloVe: Global vectors for word representation," in *EMNLP '14*.  Doha, Qatar: ACL, Oct. 2014, pp. 1532–1543.

[46] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *ICML*, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97, 2019, pp. 6105–6114.

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[48] P. Singh, N. Bauwelinck, and E. Lefever, "Lt3 at semeval-2020 task 8: Multi-modal multi-task learning for memotion analysis," in *SEMEVAL*, 2020.

[49] G.-A. Vlad, G.-E. Zaharia, D.-C. Cercel, C.-G. Chiru, and S. Trausan-Matu, "Upb at semeval-2020 task 8: Joint textual and visual modeling in a multi-task learning architecture for memotion analysis," 2020.

[50] M. L. McHugh, "Interrater reliability: the kappa statistic," *Biochemia medica*, vol. 22, no. 3, p. 276—282, 2012.

[51] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[52] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," in *NeurIPS '16*.  Red Hook, NY, USA: Curran Associates Inc., 2016, p. 379–387.

[53] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV '17*, 2017, pp. 2980–2988.

[54] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," 2022.

[55] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *ArXiv*, vol. abs/1904.07850, 2019.

[56] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *ECCV '18*, September 2018.

[57] J. Arevalo, T. Solorio, M. Montes-y Gómez, and F. A. González, "Gated multimodal units for information fusion," *arXiv preprint arXiv:1702.01992*, 2017.

[58] J.-H. Kim, K.-W. On, W. Lim, J. Kim, J.-W. Ha, and B.-T. Zhang, "Hadamard product for low-rank bilinear pooling," *arXiv preprint arXiv:1610.04325*, 2016.

[59] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *NeurIPS '19*, vol. 32, 2019.

[60] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang, "Visualbert: A simple and performant baseline for vision and language," *arXiv preprint arXiv:1908.03557*, 2019.

[61] M. Abavisani, L. Wu, S. Hu, J. Tetreault, and A. Jaimes, "Multimodal categorization of crisis events in social media," in *CVPR*, 2020, pp. 14 679–14 689.

[62] Q. Zhang, J. Fu, X. Liu, and X. Huang, "Adaptive co-attention network for named entity recognition in tweets," in *AAAI'18/IAAI'18/EAAI'18*. AAAI Press, 2018.

[63] C.-B. Zhang, P.-T. Jiang, Q. Hou, Y. Wei, Q. Han, Z. Li, and M.-M. Cheng, "Delving deep into label smoothing," *IEEE Tran. on Image Proc.*, vol. 30, p. 5984–5996, 2021.

[64] D. Kiela, S. Bhooshan, H. Firooz, E. Perez, and D. Testuggine, "Supervised multimodal bitransformers for classifying images and text," *arXiv preprint arXiv:1909.02950*, 2019.

[65] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[66] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," *IJCV*, vol. 128, no. 2, p. 336–359, Oct 2019.

[67] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *NeurIPS*, vol. 26, 2013.

[68] T. Mikolov, K. Chen, G. S. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *ICLR*, 2013.

[69] V. Keswani, S. Singh, S. Agarwal, and A. Modi, "IITK at SemEval-2020 task 8: Unimodal and bimodal sentiment analysis of Internet memes," in *SemEval*, Barcelona (online), Dec. 2020, pp. 1135–1140.

[70] G. A. Vlad, G. E. Zaharia, D. C. Cercel, C. Chiru, and S. Trausan-Matu, "UPB at SemEval-2020 task 8: Joint textual and visual modeling in a multi-task learning architecture for memotion analysis," in *SemEval-2020*, Barcelona, Dec. 2020, pp. 1208–1214.

[71] Y. Guo, J. Huang, Y. Dong, and M. Xu, "Guoym at SemEval-2020 task 8: Ensemble-based classification of visuo-lingual metaphor in memes," in *SemEval-2020*, Barcelona, Dec. 2020, pp. 1120–1125.

[72] D.-K. Nguyen and T. Okatani, "Improved fusion of vis. and lang. representations by dense symmetric co-attention for vqa," in *CVPR*, June 2018.

[73] M. Grootendorst, "Bertopic: Neural topic modeling with a class-based tf-idf procedure," *arXiv preprint arXiv:2203.05794*, 2022.

[74] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using Siamese BERT-networks," in *EMNLP-IJCNLP '19*.  Hong Kong, China: ACL, Nov. 2019, pp. 3982–3992.

[75] L. McInnes, J. Healy, N. Saul, and L. Großberger, "Umap: Uniform manifold approximation and projection," *J. of Open Src. Soft.*, vol. 3, no. 29, p. 861, 2018.

[76] R. J. G. B. Campello, D. Moulavi, and J. Sander, "Density-based clustering based on hierarchical density estimates," in *PAKDD '13*, J. Pei, V. S. Tseng, L. Cao, H. Motoda, and G. Xu, Eds.  Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 160–172.

[77] R. Plutchik, *The Nature of Emotions: Clinical Implications*.  Boston, MA: Springer US, 1988, pp. 1–20.

[78] B. Magnini, E. Pianta, C. Girardi, M. Negri, L. Romano, M. Speranza, V. B. Lenzi, and R. Sprugnoli, "I - cab: the italian content annotation bank." in *LREC*, 2006, pp. 963–968.

[79] C. Jennifer, F. Tahmasbi, J. Blackburn, G. Stringhini, S. Zannettou, and E. D. Cristofaro, "Feels bad man: Dissecting automated hateful meme detection through the lens of facebook's challenge," 2022.

**Shivam Sharma** is a Ph.D. student in the Dept. of Electrical Engineering at Indian Institute of Technology Delhi (IIT Delhi), India. His research interests span around multimodal applications within NLP. He also works as a Lead Research at Wipro AI Research (Lab45), Wipro Ltd. He completed his MS (by Research) in Signal Processing (Acoustics) from the Indian Institute of Information Technology, Sri City (IIIT-S) in 2020.

**Ramaneswaran S** works as a Data Scientist at NVIDIA. He completed his Bachelor of Information Technology in the School of Information Technology and Engineering at Vellore Institute of Technology (VIT), Vellore, Tamil Nadu, India. His broad research interests include Natural Language Processing and Conversational AI.

**Md. Shad Akhtar** is an Assistant Professor at Indraprastha Institute of Information Technology, Delhi (IIIT-D). He completed his Ph.D. in Computer Science and Engineering from IIT Patna in 2019. He received his M.Tech from IIT (ISM), Dhanbad. He also has over two years of industry experience with HCL Tech. Ltd. His main research areas are Sentiment and Emotion Analysis in the Natural Language Processing domain. Currently, his area of interest focuses on Dialog Management and Multimodal Analysis.

**Tanmoy Chakraborty** is an Associate Professor of Electrical Engineering and an Associate Faculty member of the Yardi School of AI at IIT Delhi. He leads the Laboratory for Computational Social Systems (LCS2), a research group specializing in NLP and Computational Social Science. His current research primarily focuses on empowering frugal language models for improved reasoning, grounding, and prompting. He received numerous awards, including the Ramanujan Fellowship, PAKDD Early Career Award, ACL'23 Outstanding Paper Award, IJCAI'23 AI for Good Award, and several faculty awards/gifts from industries like Facebook, Google, LinkedIn, JP Morgan, and Adobe. He has authored a textbook on "Social Network Analysis". More details may be found at tanmoychak.com.

TABLE 7: Top 10 most frequent words in each emotion class. The TF-IDF score is in the parenthesis.

| Fear | Anger | Joy | Sadness | Surprise | Disgust |
|------|-------|-----|---------|----------|---------|
| mom (27.2419) | face (40.5642) | love (56.7116) | depression (45.3137) | realize (17.5789) | absolutely (38.0648) |
| scared (21.0282) | mad (33.5129) | day (38.6249) | life (40.7641) | like (16.5175) | face (11.9207) |
| pick (15.2956) | like (29.2126) | excited (38.2704) | like (38.8570) | oh (14.6915) | people (11.6430) |
| people (13.5623) | know (23.5728) | friend (36.5688) | anxiety (33.9166) | meme (13.0155) | make (9.8716) |
| hear (12.0607) | make (23.2726) | mom (35.4429) | day (31.8213) | people (12.3859) | food (7.7485) |
| spider (11.9407) | just (22.9047) | good (35.2838) | depressed (29.3055) | time (12.2300) | meme (7.3623) |
| afraid (11.3161) | people (22.6959) | friends (33.2909) | lonely (29.1180) | just (10.7769) | look (6.8671) |
| says (10.9221) | look (21.8237) | like (27.5951) | going (28.7829) | mom (10.2412) | like (6.3189) |
| home (9.0328) | time (20.6286) | make (26.5385) | friends (26.8488) | face (9.2913) | realize (6.1869) |
| time (8.6592) | say (20.5275) | best (25.1675) | feel (26.7687) | hell (8.8602) | just (6.0652) |

# APPENDIX A
## ADDITIONAL DETAILS OF MOOD

This section provides additional details on collecting and curating our proposed dataset MOOD.

### A.1 Filtering Criteria

For downloading meme images, we used the Mozilla Firefox extension tool, called Download All Images[8], with a few downloading specifications configured. These were file size (min): 4 KB, dimensions: 200X200, and format: JPGs and PNGs. We set these specifications after carefully observing the sample quality of memes available online and the requirements of the task at hand. Since despite pre-setting the required specifications, the download process ended up collecting images that were still unsuited towards manual annotations, the annotators were asked to further manually filter out images based on *filtering criteria* specified in Section 3.1 in the main text. These factors involved inadequate image resolution and text readability (perceptually ambiguity), absence of any of the six Ekman emotions, harmful memes containing personal information, and memes containing non-English textual content. Our primary heuristic for keeping a meme was the perceived intelligibility w.r.t. the textual and visual cues present in it. This ensured better interpretability of the model outputs as well. We did not consider any pre-defined (or otherwise) resolution threshold after collecting the raw meme images.

### A.2 Data Imbalance

Here, we want to highlight a popular effort towards investigating hateful memes via Hateful Memes Challenge [24]. This, although involved the curation of a balanced combination of hateful and non-hateful memes focusing on modality-specific nuances, did involve the inclusion of benign confounders towards evaluating the robustness of multimodal systems, but were created synthetically by adopting confounding strategies, essentially not reflecting the realistic data distribution. This effect has been empirically observed to exacerbate when evaluated for the content over other social media platforms like 4chan (/pol/) [22], [79]. Keeping in mind the adverse implication of the non-realistic dataset, we instead emphasized collecting and curating a dataset that not only captures the fine-grained aspects of the primary task we intended to address but also reflects the realistic distribution, offering the scope for imminent developments and hence novelties in the areas like un/self-supervised and few-shot learning. This has already demonstrated capabilities for characterizing harmful content over social media platforms[9].

8. Firefox Browser ADD-ONS — Download All Images ⌂
9. Meta AI — ML Applications ⌂

### A.3 Thematic Analysis

Towards performing thematic analysis for MOOD, we leverage a popular topic modelling technique, called BERTopic [73] that uses transformers and c-TF-IDF to create dense clusters. For the thematic analysis of visual objects, the overall pipeline first converts images into embeddings, followed by the performing dimensionality reduction, followed by HDBSCAN-based dense clustering. This is followed by captioning the images while weighting the cluster representative bag-of-words using c-TF-IDF and finding the best matching images based on most representative documents. Additionally, we also assess the tf-idf ranked set of words from each categorical distribution in MOOD, as shown in Table 7.
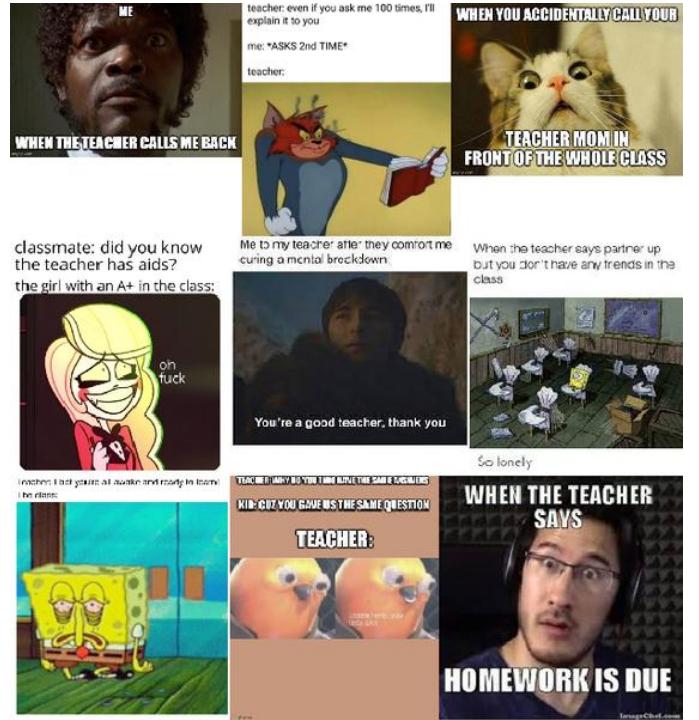


Fig. 11: A collection of meme examples featuring *human subjects, pop culture references, animated characters, and animals* from MOOD, with representative topics as *teacher, class, school, homework, test, kid, and exam.*

We look for the visual diversity captured within the MOOD dataset via manual and automated assessment. The manual review suggests a pre-dominant visual representation of human subjects, pop culture references, animated characters, and animals in the memes. In addition, the memes typically consist of various artistic modifications of these basic elements –*visual morphing* and *juxta-*

(a) ALFRED (w. Frozen Emotion Encoder)

(b) ALFRED (w. Fine-tuned Emotion Encoder)

(c) Category-wise Confusion ($\text{FN}_{\text{E}_i \to \text{N}} - \text{TP}_{\text{E}_i \to \text{E}_i}$)
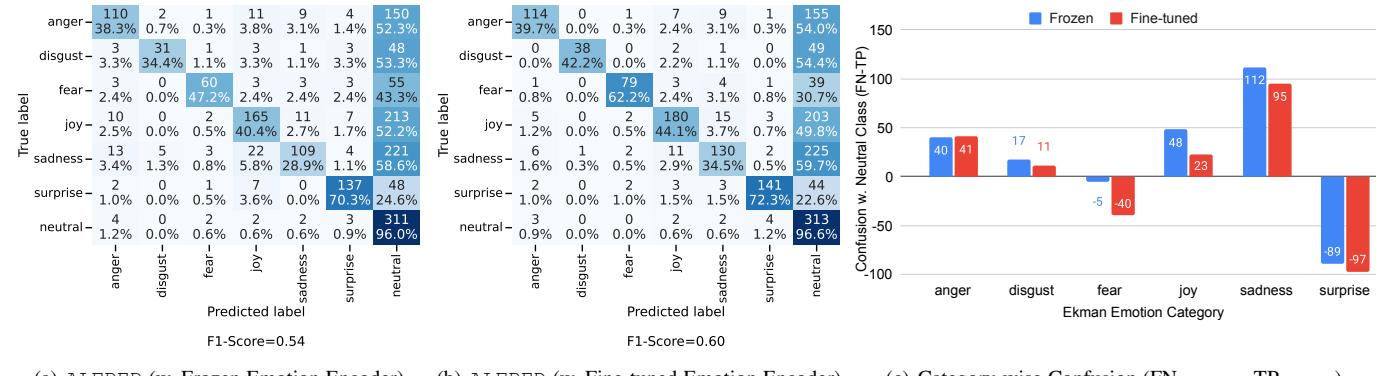
Fig. 12: Analyzing ALFRED's performance *with neutral* category. (a) and (b) Confusion Matrices and F1-score for ALFRED's two variants; (c) Quantifying Confusion: $\text{FN}_{\text{E}_i \to \text{N}} - \text{TP}_{\text{E}_i \to \text{E}_i}$: Difference between the *false-negatives* (FN) w.r.t *neutral* class (N) and *true-positives* (TP) for each Ekman emotion category ($\text{E}_i$).
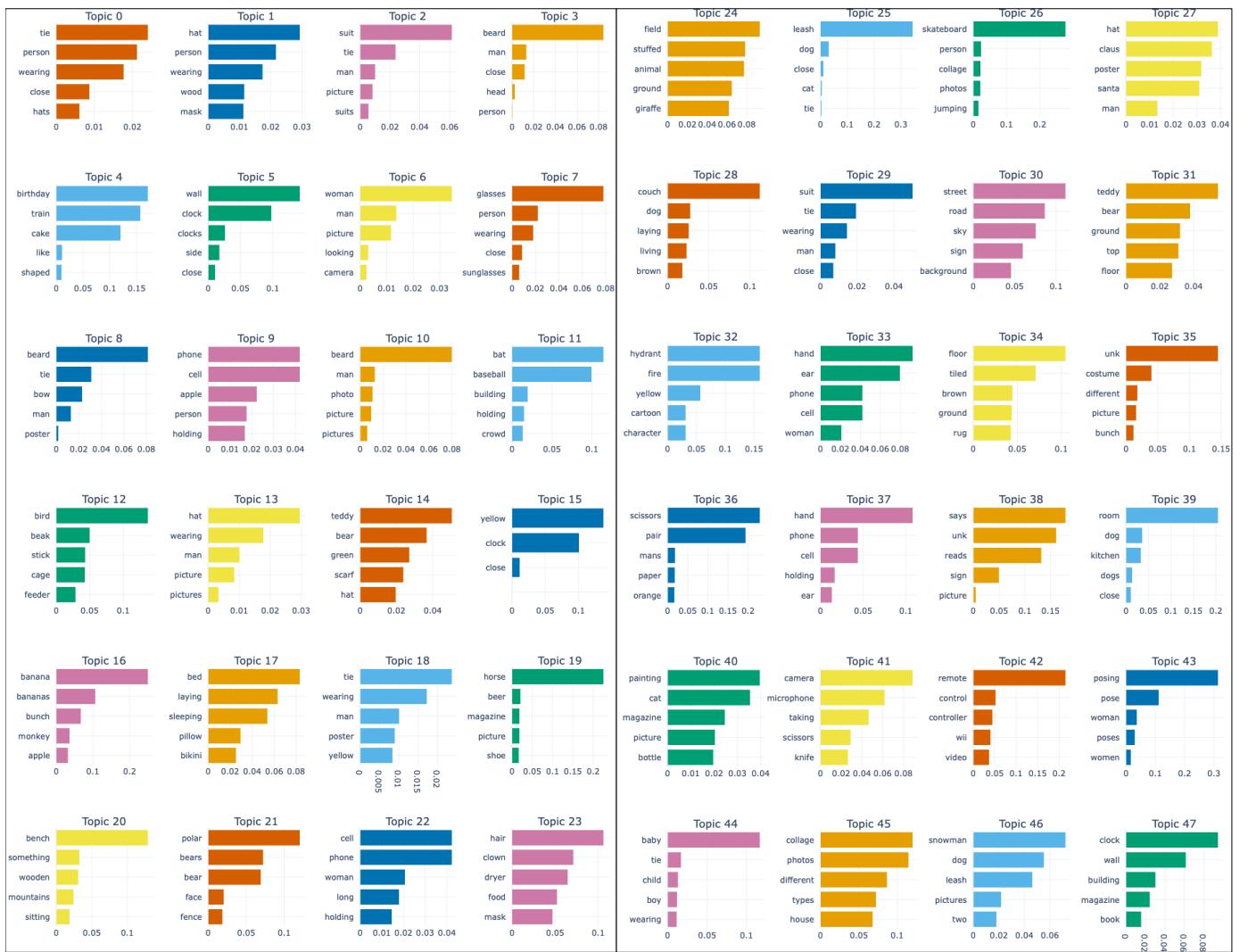


Fig. 13: Top$-48$ prominent topics representing themes of the *visually depicted content* in MOOD's memes.

*positioning*, along with diversified *textual overlays*. A representative set of such samples from MOOD is shown in Fig. 11 for *third largest topic cluster* associated with meme's visual embeddings. This topic is defined by words *teacher, class, school, homework,*

*test, kid, and exam*. We choose this set for exemplification of MOOD's visual diversity, as it has a relatively more diverse set of meme templates and visual subjects utilized as variants, in comparison to that within the larger topic clusters defined by
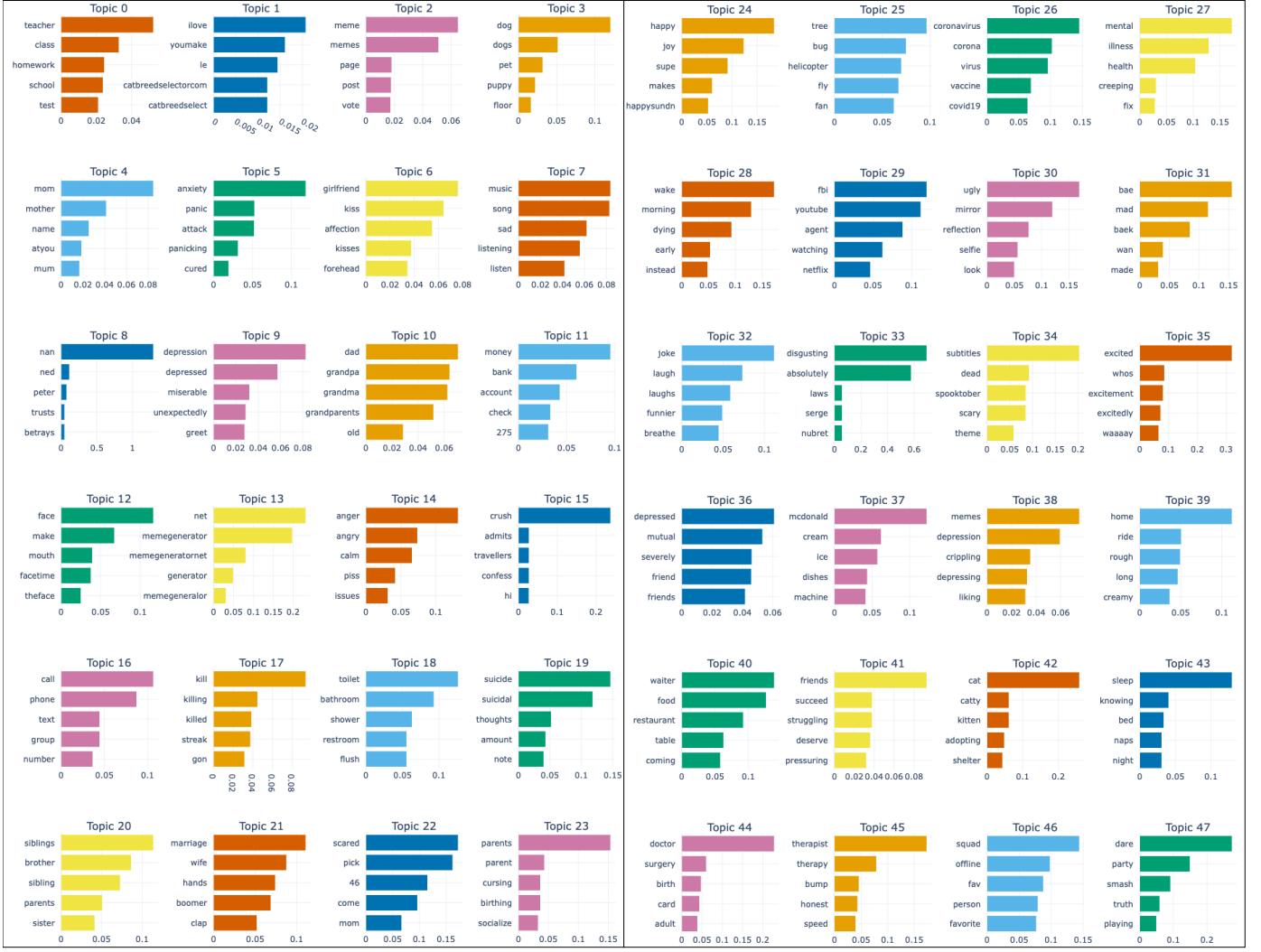
Fig. 14: Top−48 prominent topics representing themes of the *textually embedded content (OCR)* in MOOD's memes.

words - *disgusting, happy, she, crush, girlfriend, cute*, which are dominated by template-based meme designs. Further, we also analyze the bar charts of the top 48 topic clusters defined by the caption keywords corresponding to the image-embedding-based clusters. These are shown in Fig. 14.

## APPENDIX B
## ALFRED'S PERFORMANCE WITH 'NEUTRAL' MEMES

We examine ALFRED's efficacy/constraints towards the multi-class classification setup involving *fear, anger, joy, sadness, surprise,* and *disgust*, along with a seventh category, 'neutral'. We evaluate ALFRED, by training it using memes from *neutral* category as well. To this end, we utilize a total of 2197 *neutral* category memes (1549 memes in training, 324 memes in validations, and 324 memes in testing) from the publicly available dataset, Memotion [2], which (along with the other categories in the dataset) is systematically curated towards the designated categories. We do not use memes that we discard during our data collection process towards considering *neutral* class, as it mostly consists of *low-quality, noisy, or harmful* content, the generalizability towards which is accounted for in more realistic settings, as part of Section 6.8 (main text). We compare the performance of ALFRED when the emotion encoder weights are: (a) *frozen*, and (b) *fine-tuned*. As can be observed from the F1-scores in Figs. 12 (a) and (b), the overall performance of ALFRED gets reduced from 0.82 F1-score, when modeled for only *six* Ekman emotions, to the sub-par scores of 0.54 and 0.60 F1-scores for *frozen* and *fine-tuned* emotion-encoder-based scenarios, respectively. This highlights the limitations that just training on Ekman emotion-based samples, without considering the confounding effect of *neutral* class, can get induced within ALFRED's performance.

The confusion matrices for these experiments are also shown as part of Figs. 12(a) and (b), respectively. The overall relative performance pattern in terms of the difference b/w Ekman emotion category-specific *true-positives* (TP) and *neutral* class (FN) is distinctly reflected in Fig. 12(c). We observe that all the Ekman emotion categories get mixed up in different proportions with the *neutral* category, with the *most* confused class being *sadness* with an FN-rate of 58.6% and 59.7% for *frozen* and *fine-tuned* variants of ALFRED, respectively. At the same time, the *least* confused category is *surprise*, with an FN-rate of 24.6% and 22.6% for the corresponding variants. This observation, along with the slightly better accuracy produced by CLIP-based (reference) baseline (0.764) for class *sadness* (c.f. Fig. 3) hints at the

Fig. 15: Comparison b/w the quality of the OCR-extracted text via (a) Tesseract OCR, and (b) Google OCR.

utility of leveraging more *contextually enriched* representations towards discriminating it against the rest. Moreover, the distinct clarity of `ALFRED` towards discriminating a class like *surprise* is also corroborated by an imposing $6\%$ lead against our reference baseline (also having the second best category-specific) score. It is also worth noting that only minor confusions for *neutral* class, being predicted as any other emotion category, are observed.

Additionally, the general trend of distinct reduction as shown in Fig. 12(c), in *differences* between the *true-positive rate* (TPR) for Ekman emotions and *false-negative rate* (FNR) w.r.t the *neutral* class, when the emotion-encoder in `ALFRED` is *fine-tuned* (over the *frozen* variant), clearly prescribes the effect of *adapting* the emotion-encoder module, towards overall emotion classification. With the subtle exception of *anger* class (exhibiting the enhancement of TPR-FNR difference by one sample), all the other classes project reasonable reductions in the overall confusion between Ekman emotions and *neutral* category, quantified by the *absolute* differences of $6, 35, 25, 17,$ and $8$ samples for classes: *disgust, fear, joy, sadness,* and *surprise*, respectively.

The amount of confusion visible between Ekman emotions and the neutral category suggests further scope of improvement in terms of out-of-distribution generalizability for `ALFRED`.

# Appendix C
## Text Extraction via OCR

Text extraction via *optical character recognition* (OCR) is of critical importance when mining embedded text from memes. The quality of the OCR process utilized influences the overall modeling capacity of systems. Towards exploring an optimal OCR technique for our purpose, we compare the text extractions for *two* popular OCR-based text extraction APIs: Google Tesseract API[10] (TOCR) and Google GCV API (GOCR). We first qualitatively analyze the extraction quality for 30 random memes and find occasional mistakes by TOCR, and rare by GOCR. For TOCR, mistakes committed were mostly for difficult cases, like text-image embedded at the same location, poor quality graphics, small text, etc. Sometimes even for simple cases, we observe GOCR's text quality much better than TOCR's output. A couple of examples shown in Fig. 15 demonstrate the difference in the text-extraction quality for TOCR and GOCR. The first example shown in Fig. 15 (*left*) is the case consisting of a mix of simple

and complex regions like black text on white background and ambiguous visual-text overlap, respectively, that cannot be correctly mined by TOCR, while GOCR, is distinctly more accurate in its extraction. On the other hand, the second, relatively simpler meme in Fig 15 (*right*) poses more obscurity to TOCR, as compared to better visibility for GOCR.

We also examine `ALFRED`'s overall performance in terms of the macro F1-score for our primary task of emotion classification for six Ekman emotions, w.r.t the two choices of OCR techniques explored. In drastic contrast to the impressive F1-score of $0.82$ observed for the GOCR-based text extraction, we find an abysmal show of performance by TOCR, with an F1-score of $0.75$, which speaks volumes of its inferiority when compared with Google GCV-based OCR API. In addition to leveraging GOCR for our primary experiments, we also conclude the critical influence that the correct OCR extraction technique has over the downstream task at hand.

---

10. Google's Tesseract-OCR API ↗