# Highlights

## Toxic Memes: A Survey of Computational Perspectives on the Detection and Explanation of Meme Toxicities

Delfina S. Martinez Pandiani,Erik Tjong Kim Sang,Davide Ceolin

- We provide a survey of over 150 papers that computationally analyze toxic memes.

- We survey labels, task definitions, sources, and usage of 34 datasets of toxic memes.

- We categorize meme toxicities, providing a meta-model to characterize toxicity dimensions.

- We identify key challenges and trends in toxic meme detection and explanation.

# Toxic Memes: A Survey of Computational Perspectives on the Detection and Explanation of Meme Toxicities

Delfina S. Martinez Pandiani[a,*], Erik Tjong Kim Sang[b] and Davide Ceolin[a]

[a]*Centrum Wiskunde & Informatica, Science Park 123, Amsterdam, 1098 XG, The Netherlands*
[b]*Netherlands eScience center, Science Park 402, Amsterdam, 1098 XH, The Netherlands*

## ARTICLE INFO

## ABSTRACT

Internet memes, channels for humor, social commentary, and cultural expression, are increasingly used to spread toxic messages. Studies on the computational analyses of toxic memes have significantly grown over the past five years, and the only three surveys on computational toxic meme analysis cover only work published until 2022, leading to inconsistent terminology and unexplored trends. Our work fills this gap by comprehensively surveying content-based computational perspectives on toxic memes, and reviewing key developments until early 2024. Employing the PRISMA methodology, we systematically extend the previously considered papers, achieving a threefold result.

First, we include in the survey 119 new papers, analyzing 158 computational works focused on content-based toxic meme analysis. Also, we identify over 30 datasets used in toxic meme analysis and examine their labeling systems. Second, after observing the existence of unclear definitions of meme toxicity in computational works, we introduce a new taxonomy for categorizing meme toxicity types. We also note an expansion in computational tasks beyond the simple binary classification of memes as toxic or non-toxic, indicating a shift towards achieving a nuanced comprehension of toxicity. Third, we identify three content-based dimensions of meme toxicity under automatic study: target, intent, and conveyance tactics. We develop a framework illustrating the relationships between these dimensions and various meme toxicities.

The survey identifies and analyzes key challenges and recent trends, such as enhanced cross-modal reasoning, integrating expert and cultural knowledge, the demand for automatic toxicity explanations, and handling meme toxicity in low-resource languages. Also, it notes the rising utilization of Large Language Models (LLMs) and generative AI for detecting and generating toxic memes. Finally, the survey proposes pathways for advancing toxic meme detection and interpretation, addressing new opportunities and research gaps. <span style="color:red">Caution: This work includes toxic memes that may cause psychological distress. Viewer discretion is strongly advised. The example memes do not represent the views or opinions of the authors.</span>

## 1. Introduction

Memes, a term introduced by Richard Dawkins in 1976, serve as cultural replicators analogous to genes, rapidly transmitting ideas between human minds and shaping collective consciousness [1]. However, memes also harbor dangers that can profoundly impact individuals and society. They can be likened to insidious "flukes" that "hijack and infect the brain" with hazardous ideas, particularly affecting those vulnerable to their influence [2]. In contemporary discourse, the term "memes" has become synonymous with "internet memes" or "image memes", referring to text-image pairs disseminating swiftly across digital platforms [3]. These internet memes are powerful tools for communication and expression, rapidly sharing ideas, emotions, and cultural references across online communities. Though often humorous and lighthearted, internet memes significantly shape public and political discourse. For example, alt-right concepts on platforms like 4chan[1] and Encyclopedia Dramatica[2] create a subcultural language community linked to violent right-wing activism [4]. Conversely, progressive leftist meme makers use memes for counternarrative techniques, dialectical seeding, and fostering solidarity within leftist communities [5]. The widespread presence of memes on social media has sparked interest and concern about their societal impact, playing a crucial role in digital culture and reflecting the collective consciousness of online societies.

*Corresponding author

✉ delfina.martinez.pandiani@cwi.nl (D.S.M. Pandiani); davide.ceolin@cwi.nl (D. Ceolin)
ORCID(s): 0000-0003-2392-6300 (D.S.M. Pandiani); 0000-0002-8431-081X (E.T.K. Sang); 0000-0002-3357-9130 (D. Ceolin)

[1]https://www.4chan.org
[2]https://encyclopediadramatica.online

**Table 1**

Potential risks associated with toxic memes, identified in literature across various disciplines including computer science, human-computer interaction, internet pragmatics, multimedia, cybernetics, visual and media studies, and more.

| General Risk | Details | Citations |
|---|---|---|
| Violence in Public Discourse | Conveying toxic representations and messages, perpetuating harmful stereotypes via hate speech, harm, abuse, cyberbullying, offensiveness, and various other forms of toxicity. | [6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16] |
| Opinion Manipulation | Serving as potent tools for the dissemination of disinformation, propaganda, and trolling, leading to polarization and misunderstanding. | [18, 19, 20, 21, 22, 23, 24] |
| Psychological Impacts | Promoting groupthink and deindividuation, destructive thoughts and behaviors, desensitizing individuals to tragic news, fostering apathy, and exacerbating psychological distress. | [25, 26, 27, 28] |
| Material-World Effects | Exerting tangible impacts of disinformation in election outcomes and instances of physical violence, contributing to the normalization of extremist behaviors. | [29, 30, 31] |

The dangers commonly associated with memes are particularly pronounced with internet memes due to their accessibility and virality, which strongly contribute to the spread of toxic ideologies and narratives. The risks associated with toxic memes are extensive (see Table 1). Internet memes can convey toxic representations and messages, perpetuating harmful stereotypes, inciting divisive rhetoric, and contributing to a climate of violence in public discourse [6]. Internet memes can serve as conduits of hate speech [7, 8, 9], harm [10, 11, 12], abuse [13], cyberbullying [14], offensiveness [15, 16] and various other forms of toxicity [10]. These memes often cross into illegal territory, exhibiting characteristics of hate speech, incitement to violence, or posing systemic risks to public discourse [17]. Additionally, memes can serve as potent tools for opinion manipulation, including the dissemination of disinformation [18], propaganda [19, 20, 21], and trolling [22, 23]. Political internet memes, often humorous, can oversimplify complex issues, potentially leading to misinterpretation [24]. Moreover, exposure to and engagement with toxic memes can have profound psychological impacts, leading to phenomena like groupthink or deindividuation [25] and promoting destructive thoughts and behaviors in individuals [26]. The normalization of dark humor, self-deprecating jokes, and derogatory slang within meme culture can blur social norms and exacerbate psychological distress [27], as well as desensitizing individuals to tragic news, fostering apathy towards important issues [28]. By trivializing violence and desensitizing individuals, extremist behaviors are normalized with real-world manifestations, such as wearing meme-inspired clothing at extremist rallies or by perpetrators of violence, as seen in the Allen, Texas mass shooting [29]. Critically, the proliferation of toxic memes on online platforms extends their impact beyond the digital realm, influencing real-world events such as election results [30] and instances of physical violence.

Addressing the proliferation of toxic internet memes is crucial for online safety, requiring effective detection and moderation strategies considering their nuanced features. However, moderating toxic internet content is incredibly complex due to fuzzy decision boundaries influenced by cultural considerations [32], disagreements among annotators, unconscious biases, and the nuanced nature of harmful language [32]. Currently, moderation efforts often rely on human crowdworkers in regions like the Philippines, India, and Kenya, who face inadequate protection and significant mental health challenges due to their exposure to toxic content [33, 34, 35, 36, 37]. These perils to the workers and the sheer volume of online content exacerbate moderation challenges, prompting a growing demand for automated solutions to detect toxicity and explain their assessment. In essence, there is an urgent need for explainable toxic meme detection, including the identification of specific toxicity types and the provision of criteria that the toxicity is based on, such as slur words, hate symbols, or intricate rhetorical strategies.

In recent years, there has been a significant increase in computational analysis of internet memes, particularly in research on meme toxicity. This trend is evident from the rising number of Scopus-indexed computer science manuscripts addressing the harmfulness of memes (Figure 1). Despite this growth, the existing literature lacks comprehensive systematic reviews. Only three previous surveys on toxic meme analysis [10, 38, 39] exist, but they have limited scope and temporal coverage, ending in 2022. Recognizing the rapidly evolving nature of this field, our survey aims to incorporate the latest advances and insights. Our goal is to provide researchers, practitioners, and policymakers with a comprehensive understanding of toxic memes from a computational, content-based perspective, covering key developments up to early 2024. Following the PRISMA methodology, our survey systematically reviews 158 computational works selected for their focus on content-based analysis of toxic memes. Content-based analysis refers to analyzing meme content
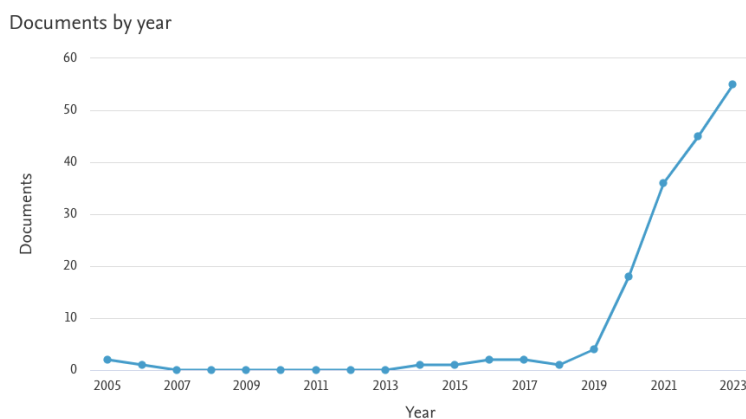
Documents by year



**Figure 1:** Graph depicting the exponential increase in publications within the field of computer science, as indexed by SCOPUS, focusing on research related to toxic memes. The data was gathered using a query targeting specific keywords associated with meme toxicities (see Section 4).

independently of their propagation. While topics like meme propagation and contextual dynamics within social networks, such as measuring visual similarity between image memes across polarized Web communities [40], and measuring the potential of harmful memes as precursors for spreading hate on platforms like Twitter [41], are increasingly studied, they are beyond the scope of this survey. The key contributions of this survey center around the following:

- **Analysis of Datasets and Tasks**: We provide an updated overview of available datasets containing toxic memes, including details such as labels, task definitions, sources, and their usage in research, highlighting the diversity in task definitions and labeling schemes.

- **Conceptual Overview and Standardized Taxonomy**: We comprehensively review concepts and definitions of computationally studied meme toxicities, addressing the lack of consensus in their definition. Additionally, we introduce a standardized taxonomy for categorizing meme toxicity types.

- **Dimensions of Meme Toxicities**: We identify specific dimensions of toxicity in memes (i.e., target, and conveyance tactics), outlining a framework illustrating the relationships between these dimensions and various meme toxicities.

- **Identification of Common Trends**: We identify prevalent challenges and recent trends, including enhanced multimodal reasoning and interpretability, integration of expert and cultural knowledge, and addressing meme toxicity in low-resource languages. We also note the increasing use of LLMs and generative AI in detecting and creating toxic memes.

- **Proposing Pathways for Advancement**: We suggest pathways to tackle challenges and advance the field of toxic meme detection and interpretation.

This survey is structured as follows. In Section 2, we provide an overview and definitions of (internet) memes and online toxicity. In Section 3, we discuss previous surveys on toxic memes. Section 4 outlines our approach to selecting the papers surveyed. In Section 5, we analyze and discuss the extent of coverage provided by this survey compared to three prior surveys. In Section 6, we provide a catalog of toxic meme datasets and associated label and task definitions. Section 7 introduces a novel taxonomy for categorizing the types of meme toxicity. Section 8 outlines a framework illustrating the relationships between dimensions of toxicity. Section 9 explores common challenges and recent trends in computational approaches to detect and interpret toxic memes. Lastly, in Section 10, we identify key future research directions and we conclude in Section 11.

## 2. Background

### 2.1. Defining (Internet) Memes
The concept of *memes* draws parallels to biological *genes* and evolution, describing ideas that self-replicate, evolve, and respond to selective Darwinian pressure, ultimately entering culture in ways similar to biological genes and

**Table 2**
Features of (internet) memes, along with their corresponding descriptions and relevant citations.

| Feature | Description |
|---|---|
| Multimodal | Combine visual and language information creatively [42, 46] |
| Succinct | Spread complex messages with a minimal information unit that connects virtual circumstances to real ones [46] |
| Fluid | Subject to variations and alterations [46] |
| Anomalous Juxtaposition/ Incongruity | Leverage lack of relevance in the arrangement of textual and/or visual constituents to produce unexpected outcomes [50] |
| Intertextual | Reference popular culture, symbols, artifacts, or events that hold meaning within the community of reference [46, 50] |
| Relatable/Tacit Background | Rely on viewer's familiarity with certain contextualized aspects of the world [53], shared knowledge, and implicit cultural references [54] |

potentially modifying human behavior [1]. However, the advent of the digital era has led to a resemantization of the term [42], with internet memes being intentional modifications of Dawkins' original concept, characterized by creative alterations rather than random mutation [43]. In the digital realm, memes are often defined as digital artifacts sharing common traits of content, form, or perspective, created by users and disseminated, imitated, or transformed via the Internet [44]. While a 'memetic construct' is the foundational structure comprising form, content, and perspective, a meme represents a specific instance or manifestation of a memetic construct, embodying a singular multimodal expression of a widespread cultural idea [45]. Thus, in this survey, we use the term 'meme' to specifically denote what others commonly refer to as 'internet' memes [46], 'visual memes' [47], or 'Image With Text' (IWT) memes [48].

Multimodality stands as a defining feature of memes, as memes rely on a combination of text and images to convey complex messages [45]. Typically, memes consist of an image paired with short text, allowing for easy sharing on social media [49]. They blend visual and verbal elements to convey humor, irony, or sarcasm, often referencing cultural symbols or events [42, 50]. Studies computationally operationalizing the term 'meme' have primarily utilized such visuo-linguistic association analysis to distinguish memes from non-meme images and relied on implicit judgement or crowdsourced annotations for differentiation [51]. Recent research [52] employs multi-channel convolutional neural networks to distinguish memes from not only photographs but also other image-with-text (IWT) formats, such as advertisements, movie posters, online news articles, or screenshots of posts, commonly circulated online. Table 2 summarizes key features of internet memes and associated citations.

## 2.2. Defining Online Toxicity

Online toxicity typically refers to negative behaviors on the internet that damage others' or even one's own self-image, hindering personal growth [55]. This definition is adopted by researchers studying text-image combinations [56] and those exploring 'toxic memes' [4]. Understanding online toxicity is challenging due to its multifaceted nature, including fine-grained categories and overlapping terms, requiring models capable of recognizing various aspects of toxic behavior [57]. Defining and detecting online toxicity, particularly from a computer science perspective, is further complicated by challenges in annotating data and developing machine learning models to identify toxicity types and relationships [58]. For instance, sharing misinformation on social media is associated with harmful language, highlighting the importance of integrating research on two types of toxicity (misinformation and harmful language) usually studied separately [59].

### 2.2.1. Textual Toxicity

Much of the research on online toxicity has focused on textual data, with extensive studies on detecting toxicity in texts [60] and tasks such as toxic comment classification, hate speech detection, and identification of offensive language. Different taxonomies have been proposed to categorize abusive language, distinguishing between abusive content directed at individuals or groups, and between explicit and implicit abusive content [61, 62, 63]. Even within 'hate' speech, there are various definitions and fine-grained labels, and thus a need for comprehensive datasets to train robust models for combating hate speech effectively [64]. Some works follow a three-level taxonomy considering the type

and target of offense [63], while other toxic comment detection systems follow a multi-label classification framework, with comments labeled as 'toxic,' 'severe toxic,' 'insult,' 'threat,' 'obscene,' and 'identity hate' [65]. Other datasets use labels like 'hate speech,' 'offensive but not hate speech,' and 'neither offensive nor hate speech,' highlighting biases and challenges in classifying offensive language in short-form content like tweets [66].

The differentiation of concepts associated with toxic speech presents complexities, often characterized by contentious definitions [67]. Efforts to delineate the hierarchy of hate speech concepts within computer science literature reveal fuzzy boundaries between toxicity-related concepts [68], and instances of interchangeability and even conflicting hierarchical relations. For instance, some perspectives consider toxicity as a subset of hate speech [69]. In response to these challenges, some initiatives aim to harmonize toxicity labels for textual content. For instance, in [70], researchers analyzed six publicly available datasets related to hate speech, aggression, and toxicity in text. Their objective was to standardize categories to ensure consistency and comparability across datasets. This involved merging related categories, clarifying ambiguous labels, and aligning similar concepts under common headings. The resulting taxonomy identifies various types of toxicity and recommends considering unique definitions and contexts before merging terms.

### 2.2.2. Multimodal Toxicity

Exploring multimodal toxicity from a computational perspective has not received as much attention as textual (unimodal) toxicity, yet understanding its dynamics is increasingly crucial [71]. Recent studies have begun to bridge this gap, offering varied perspectives on toxic multimodal content found online. This section provides an overview of these works, each presenting diverse taxonomies summarized in Table 3, with varying levels of granularity and emphasis on different facets of harmful content. For instance, Banko et al. [72] provide a taxonomy covering hate, harassment, self-inflicted harm, ideological harm, and exploitation, while Nakov et al. [73] offer a simpler list of harmful categories. Halevy et al. [74] focus on violations for malicious purposes, including misinformation and community standard violations such as hate speech and crimes. Pramanick et al. [49] distinguish between hateful, offensive, and generally harmful memes, while Sharma et al. [10] present a taxonomy for internet memes with categories like hateful, offensive, propaganda, harassment/cyberbullying, violence, self-inflicted harm, and exploitation. As seen in Table 3, there are commonalities and differences across these taxonomies. For instance, both [72] and [73] cover hate speech, harassment, and violence. However, [72] offers more detailed subcategories like doxing and identity attack, whereas [73] includes a broader range, such as dangerous organizations/people and glorifying crime.

To elucidate potential taxonomical relationships among the toxicities presented in Table 3, we constructed a Venn diagram, available in the appendix (refer to Appendix section A), to illustrate their intersections and semantic relationships. This process highlighted the complex nuances in defining toxic or harmful multimodal content. For example, we found that many behaviors are categorized by platforms as *misbehavior*, which includes actions that may not necessarily be toxic or harmful but are nonetheless restricted on many platforms, such as nudity. We identified diverse forms of harmful content, including ideological harm, hatefulness targeting protected groups, harassment, and more, emphasizing the multifaceted nature of online toxicity. Echoing the findings of [68] regarding textual toxicities, our analysis revealed the inherent challenges in delineating boundaries between different types of toxicities. This exploration also paves the way for identifying novel forms of toxic content that emerge uniquely in multimodal contexts and on previously unrecognized forms of harmful behavior that may not be contingent upon multimodality.

**Table 3**
Comparison of Categories and Subcategories in Different Taxonomies of Toxic Content

| Taxonomy | Source | Focus | Top-Level Distinctions | Second Level Distinctions | Granular Distinctions |
|---|---|---|---|---|---|
| Banko et al 2020 [72] | Research | Harmful content | Hate/ Harassment | Doxing | |
| | | | | Identity Attack | |
| | | | | Identity Misrepresentation | |
| | | | | Insult | |
| | | | | Sexual Aggression | |
| | | | | Threat of Violence | |
| | | | Self-Inflicted Harm | Eating Disorder Promotion | |
| | | | | Self-Harm | |
| | | | Ideological Harm | Misinformation | |
| | | | | Extremism, Terrorism & Organized Crime | White Supremacist Extremism |
| | | | Exploitation | Adult Sexual Services | |
| | | | | Child Sexual Abuse Materials | |
| | | | | Scams | |
| Nakov et al 2021 [73] / Arora et al 2023 [17] | Social Media | Policy clauses | Violence | | |
| | | | Dangerous Orgs/people | | |
| | | | Glorifying Crime | | |
| | | | Illegal Goods | | |
| | | | Self-Harm | | |
| | | | Child Sexual Abuse | | |
| | | | Sexual Abuse (Adults) | | |
| | | | Animal Abuse | | |
| | | | Human Trafficking | | |
| | | | Bullying and Harassment | | |
| | | | Revenge Porn | | |
| | | | Hate Speech | | |
| | | | Graphic Content | | |
| | | | Nudity and Pornography | | |
| | | | Sexual Solicitation | | |
| | | | Spam | | |
| | | | Impersonation | | |
| | | | Misinformation | | |
| | | | Medical Advice | | |
| Pramanick et al 2021 [49] | Research | Harmful memes | Hateful | | |
| | | | Offensive | | |
| | | | Other (Generally) Harmful | | |
| Halevy et al 2022 [74] | Social Media | Violations /Malicious Purposes | Misinformation | | |
| | | | Community Standards Violations | Hate Speech | |
| | | | | Crimes | Selling Illegal Drugs |
| | | | | | Coordinating Sex Trafficking |
| | | | | | Child Exploitation |
| Sharma et al 2022 [10] | Research | Harmful memes | Hateful | Doxxing | |
| | | | | Identity Attack | |
| | | | | Identity Misrepresentation | |
| | | | | Insult | |
| | | | | Racist | |
| | | | | Misogynistic/Sexist | |
| | | | | Sexual Aggression | |
| | | | | Extremism, Terrorism & Organized Crime | |
| | | | Offensive | | |
| | | | Propaganda | | |
| | | | Harassment/Cyberbullying | | |
| | | | Violence | | |
| | | | Self-inflicted Harm | Eating Disorder Promotion | |
| | | | | Self-harm | |
| | | | Exploitation | Adult Sexual Service | |
| | | | | Child Sexual Abuse Material | |
| | | | | Scams | |

# 3. Related Work

## 3.1. Surveys on Toxic Memes from Computational Perspective

Only three published works have surveyed the emerging field of computational toxic meme analysis [10, 38, 39]. Afridi et al. (2020) [38] were the first to survey efforts in automatic meme understanding, highlighting key challenges for future research. These challenges include defining hate in memes, distinguishing between humor and hate, and categorizing memes into subcategories targeting relevant issues. They stressed the need for detailed analysis within toxic memes, covering subcategories like hateful/non-hateful, rumor, fake news, and extremism. Additionally, they emphasized advancing techniques for cross-modal entailment in meme interpretation, crucial for tasks like automatic detection. Based on their literature review, they proposed a generic multimodal architecture for meme classification.

Building upon this foundation, Sharma et al. (2022) [10] conducted a thorough survey to categorize harmful meme types, proposing a new typology including hate, offensive, propaganda, harassment/cyberbullying, violence, and self-inflicted harm. This framework offers a structured approach for understanding and classifying harmful memes based on their characteristics and potential impact. The survey identified significant gaps in current research, highlighting challenges and future directions. Authors stress the need for detailed analysis in detection and interpretation, noting limited exploration of certain meme types due to dataset shortages, like those depicting self-harm and extremism. They also emphasize the global impact of memes, requiring research on cross-cultural implications. Furthermore, the study underscores the semiotic complexity of interpreting multimodal meme content, revealing a need for sophisticated analysis techniques to understand nuanced meanings conveyed through visual and linguistic elements.

Most recently, Hermida and dos Santos (2023) [39] surveyed methodologies for detecting hateful memes and introduced a taxonomy of machine learning architectures specifically for this purpose. This taxonomy operates on three levels: Level 1 distinguishes between non-attention mechanism-based and attention mechanism-based methods, with the latter generating multimodal representations for memes. Level 2 categorizes methodologies based on how they handle text and image components, with "restricted" methods directly using meme text and image information, while "extended" methods utilize indirectly extracted data like object tags and sentiment analysis. Both approaches are identified within attention mechanism-based methods, while non-attention mechanism-based methods fall under the restricted category. Level 3 considers the feature extraction process, distinguishing between auto-feature extraction techniques and hand-crafted techniques. Key insights from the survey include the importance of dataset quality and annotation guidelines in training models to detect toxic memes, highlighting limitations such as biases in datasets, small sample sizes, and the ongoing need for human moderation despite advancements in deep learning feature extractors.

## 3.2. Other Relevant Surveys On Hate Speech and Disinformation

Other surveys exploring areas such as multimodal disinformation and hate, while not specifically focusing on memes, provide valuable insights into broader issues surrounding multimodal toxicity. For instance, [75] survey computational approaches to multimodal disinformation and harm across different content types: text, speech, images, videos, and network data. Key findings underscore the importance of explainability in model interpretation, the necessity of considering personal preferences and cultural aspects beyond content and network signals, and the promise of knowledge-based approaches for factuality checking. [69] examine hate speech detection across multimodal and multilingual contexts, exploring various content types, including text, images, and videos, across platforms like private messages, stories, authorized account posts, comments, tweets, ads, user profiles, and sensitive content such as graphic violence and adult material. The authors advocate for proactive measures like blocking or reporting trolls on social media platforms, promoting data analysis before sharing posts, and stressing the importance of robust policy frameworks to counter abusive behavior by social media entities. A recent survey [76] primarily analyzing hate speech in text but also emphasizing the importance of multimodal features, identifies the prevalent reliance on surface-level features such as word frequency, punctuation usage, capitalization, specific keywords or phrases, and basic syntactic structures in hate speech detection approaches. The authors note that this reliance may limit the ability to fully capture the context and meaning of the text, highlight the lack of comparative studies, and stress the importance of open-source code and dataset links for evaluations. Although not addressing multimodality, two other surveys on textual toxicity are worth mentioning: [77] categorizes offensive language types into a hierarchical taxonomy, differentiating between explicit and implicit language, while [78] provides insights into the subjective nature of toxicity detection, biases in existing datasets, the influence of content source and topic on dataset characteristics, and challenges in collecting toxic comments.

## 3.3. Need for this Survey

While the surveys discussed in section 3.1 provide valuable insights, they leave a critical gap in the comprehensive analysis of what exactly constitutes "toxicity" or "harm" in memes from a computational perspective. While advocating for refinement of categorization schemes and clearer taxonomies and definitions, none delve deeply into the nuances of terminology interchangeability (e.g., toxic, harmful, malicious), their definitions, or the hierarchical relationships among different categories of harmful memes. Moreover, there remains a dearth of systematic examination regarding the alignment of dataset labels with these taxonomies. The surveys discussed have temporal limitations, covering research up to early 2022. However, in the subsequent two years, computational research on toxic memes has proliferated exponentially, and, impressively, with only three months into 2024, Scopus has indexed over 100 publications and preprints for the year. As such, several emerging trends and areas of research have not been thoroughly explored in existing surveys. These include the utilization of background knowledge and the emphasis on explainability in computational approaches [79, 80], the increasing use of LLMs for various tasks such as detecting hatefulness, misogyny, offensiveness, sarcasm, harmfulness, and specific harmful memes [79, 81], shifts towards more sophisticated evaluation methodologies [82, 83, 84, 85], novel approaches for generating toxic memes from benign prompts [86], and the emergence of new datasets, including GOAT-Bench and datasets in multiple languages beyond English [13, 87, 88].
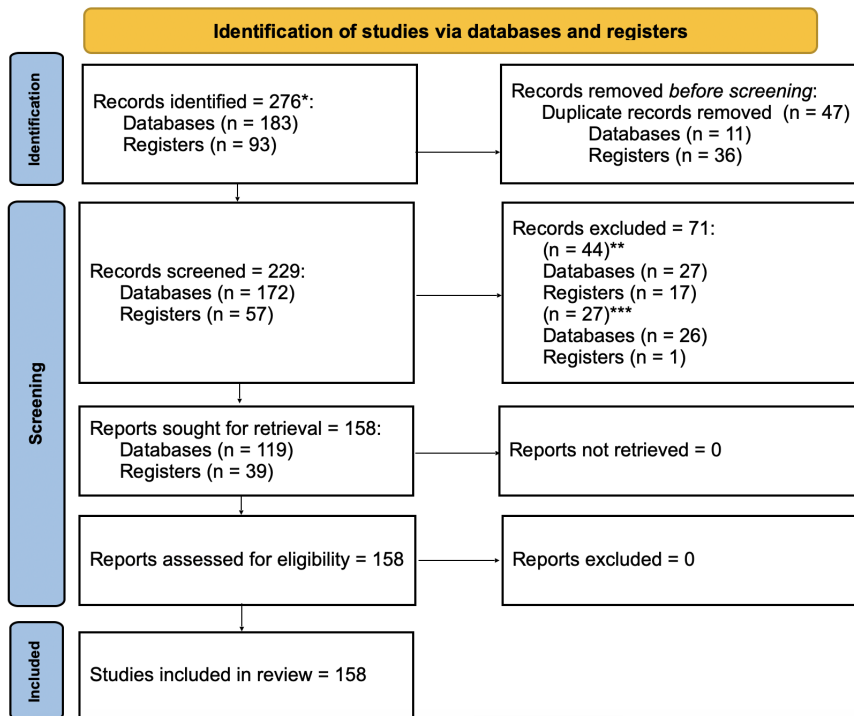
## 4. Methodology



**Figure 2:** PRISMA 2020 flow diagram for systematic reviews on *SCOPUS and Web of Science (WOS) databases. Registers refers to SCOPUS preprints. Records excluded due to: ** Topic Non-Relevance. *** Computational Non-Relevance.

*Selection of Database* We used Scopus[3] and Web of Science (WOS)[4] for our literature search because they are two of the largest and most reputable databases, indexing a wide range of high-quality journals across various disciplines and now including preprints.

---

[3]https://www.scopus.com/
[4]https://www.webofscience.com/

*Inclusion of Preprints* Acknowledging the novelty of this research field and the rapid pace of publication, we included preprints in our methodology. This decision aligns with the approach taken by the three previous surveys on this topic, emphasizing the importance of capturing the latest developments in this rapidly evolving field.

*Identification of Studies to Review* Inspired by PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses),[5] we used a structured approach to select the manuscripts. Figure 2 presents a diagram of our selection process following the PRISMA 2020 guidelines, which outlines the process of manuscript selection.

*Database Query* On March 6th, 2024, we queried the databases with a manually crafted key.[6] In SCOPUS, results were confined to "Computer Science", while in WOS, results were confined to "Computer Science Artificial Intelligence" or "Computer Science Information Systems," yielding 276 records (183 from databases and 93 from registers).

*Record Screening:* We removed 47 duplicate preprints, resulting in 229 manuscripts. We manually screened these based on abstracts and titles, excluding 71 irrelevant to our computational focus on toxic memes. We flagged 27 of those for discussion in the related work section, leaving 158 for further evaluation.

*Retrieval Assessment for Eligibility:* We retrieved all 158 records (119 peer-reviewed articles and 39 preprints). All 158 manuscripts were manually assessed and were considered eligible.

## 5. Coverage of this Survey

Our survey includes 158 papers from 2019 to 2024, as shown in Subfigure 3. The number of publications on toxic meme detection has steadily increased over time, peaking in 2023 with 52 papers. This rise indicates growing interest and research in the field. The apparent decline in 2024 is due to our survey covering only the first three months of the year. To evaluate the comprehensiveness of our systematic review on toxic memes, we conducted a systematic analysis to determine the extent to which the papers identified in our review were covered by three previous surveys conducted by Afridi et al. (2020) [38], Sharma et al. (2022) [10], and Hermida and dos Santos (2023) [39]. We manually cross-referenced each of the 158 papers in our survey with those covered in the previous works, noting whether each peer-reviewed paper or preprint was included. Detailed tracking information is provided in the project's Github page. This analysis quantified the number of our 158 papers covered by each prior survey, considering peer-reviewed papers, preprints, and the total number of papers, as well as the number of additional papers each survey reviewed compared to all previous surveys. Our results are shown in Table 4 and visually represented in Figure 3. We found that although Afridi et al. [38] call their work a survey on multimodal meme classification, they mostly discuss generic multimodal and visual understanding architectures due to a lack of specific research on memes at the time. Their surveyed papers focus mainly on hate speech detection from text and social media, with only 3 papers computationally addressing memes. We note that by the time [10] was published, the number of papers and preprints on computational meme analysis had increased. This survey extends the previous survey with 16 additional papers. Like Afridi et al. Sharma broadens their scope by including related research on textual and multimodal classification to supplement the limited work on memes. Finally, [39] survey 7 additional papers beyond those covered by [10]. Our survey reveals that most of the existing computational literature on toxic memes has not undergone a comprehensive review. Only 26 out of the 158 papers identified had been previously surveyed, meaning 84% (132 out of 158) of the papers we identified using the PRISMA methodology have not, to the best of our knowledge, been reviewed in the context of toxic meme detection and interpretation. These results are unsurprising considering the time gap and rapid rate of new publications; most papers were published after the completion of prior surveys. The most recent survey published is the work of Hermida and dos Santos [39] (2023), but the most comprehensive survey is from Sharma et al. [10], although it was published in 2022.

## 6. Toxic Meme Datasets

We thoroughly surveyed the selected works for this study, documenting the datasets used in their computational analyses. For each dataset, we recorded details such as its size (number of memes), language, sources (e.g., specific social media platforms, Google Images, etc.), the computational task(s) addressed, annotated features and classification

---

[5] https://www.prisma-statement.org/
[6] TITLE-ABS-KEY ((toxic OR harmful OR hateful OR unethical OR malicious OR malevolent OR offensive OR propaganda) AND meme)

**Table 4**
Number of papers out of the 158 identified manuscripts included in surveys about automatic detection of (toxic) memes.

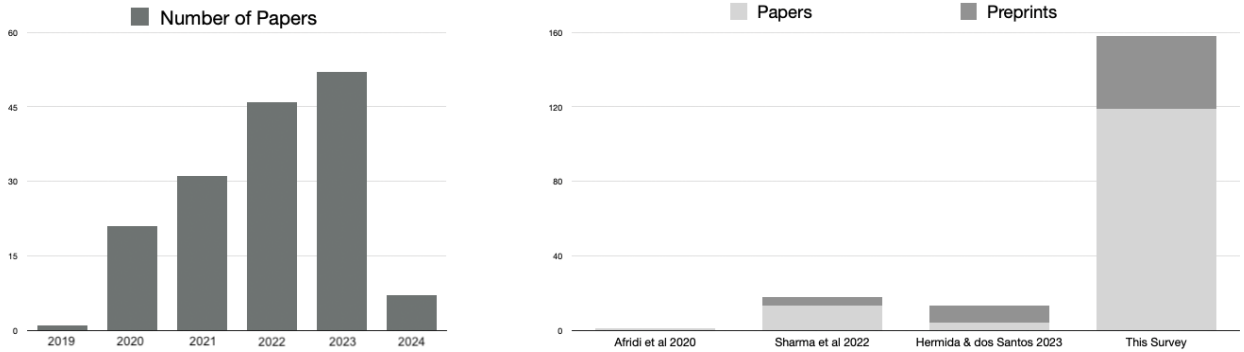| Survey | Papers | Preprints | Total | Novel Papers Considered |
|---|---|---|---|---|
| Afridi et al. (2021) [38] | 1 | 2 | 3 | 3 |
| Sharma et al. (2022) [10] | 13 | 6 | 19 | 16 |
| Hermida & dos Santos (2023) [39] | 3 | 10 | 13 | 7 |
| This survey (2024) | 119 | 39 | 158 | 132 |



**Figure 3:** Left: Distribution of surveyed papers by publication year. The figure illustrates a steady increase in the number of publications from year to year. Right: Comparison of coverage across previous surveys based on the papers surveyed here.

labels, task definitions (e.g., binary, single label multi-class, multi-label multi-class, etc.), and baseline macro F1 scores. This exhaustive examination yielded over 30 datasets (see Table 5). More details are available on GitHub.[7]

## 6.1. Dataset Characteristics

*Dataset Size* Our survey reveals a diverse spectrum of dataset magnitudes, spanning from a few hundred to tens of thousands of memes. Among the smaller datasets are Derogatory Facebook-Meme [90] with 650 memes and MultiOFF [112] with 743 memes. Conversely, larger datasets include Facebook Hateful Memes (HM) [53] with 9,540 memes, alongside MAMI [113], Memotion 1 [16], Memotion 2 [105], and MET-Meme [106], all hovering around the 10,000 mark. Notably, Innopolis Hateful Memes [102] boasts a substantial 23,000 memes. These findings underscore a considerable variance in dataset sizes within the domain. The distribution of dataset magnitudes demonstrates that most datasets lie within the range of 2500 to 10,000 memes (see Figure 4 (top)).

*Language* The majority of datasets primarily contain memes in English (see Figure 4 (middle)), with nearly 75% exclusively featuring English-language memes [15, 16, 19, 23, 49, 53, 87, 89, 90, 91, 92, 94, 95, 96, 97, 98, 99, 100, 102, 103, 104, 105, 107, 109, 110, 112]. However, we also observe a recent surge in datasets incorporating Asian languages such as Hindi [80, 93, 108], Bengali [13], and Chinese [106]. Notably, many of these datasets feature "code-mixed" memes, which blend these languages with English. Interestingly, beyond English and select Asian languages, no datasets have been utilized in the studies we surveyed that incorporate other widely spoken languages, such as Spanish.

*Data Origin* Examining the sources of meme collection revealed a diverse range of platforms, which can be categorized into four macro-categories: social media platforms (e.g., Facebook, WhatsApp), search engines (e.g., Google, Bing), image hosting platforms (e.g., Imgur, Pinterest), and dedicated meme creation and sharing resources (e.g., KnowYourMeme, 9gag) (see Figure 4 (bottom)). Social media platforms emerged as the most prominent sources, with Facebook and Reddit being the most used sources, followed by search engines, with Google being the most utilized search engine. Additionally, platforms such as Pinterest and Imgur were used for meme acquisition. Of particular interest was the utilization of meme-specific platforms like Memegenerator, KnowYourMeme, and 9gag, highlighting the importance of dedicated meme communities in the proliferation of memes. We provide detailed source information for each of the datasets in Table 9 in the Appendix. We also observed instances where new datasets were derived by

---

[7]https://github.com/delfimpandiani/toxic_memes

**Table 5**
Overview of the 34 datasets containing toxic memes identified in the literature. For each dataset, we specify the year of introduction, the corresponding manuscript, the languages included, the number of memes, the main focus/task, and the sources of the memes. Abbreviations used: SE - search engines; IP - image hosting platforms (e.g., Pinterest, Imgur); SM - social media platforms (e.g., Facebook, Reddit); MM - meme-specific resources (e.g., Know Your Meme, Memedroid).

| Dataset | Year | Language | Size | Main focus | Sources | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | SE | IP | SM | MM |
| AOMD Gab [15] | 2021 | English | 1965 | offensive meme detection | | | ✓ | |
| AOMD Reddit [15] | 2021 | English | 1094 | offensive meme detection | | | ✓ | |
| BanglaAbuseMeme [13] | 2021 | Bengali, English | 4043 | abusive meme detection | ✓ | | ✓ | |
| CrisisHateMM [89] | 2022 | English | 4700 | hateful meme detection | | | ✓ | |
| Derogatory Fb-Meme [90] | 2023 | Hindi, English | 650 | derogatory meme detection | | | ✓ | |
| DisinfoMeme [91] | 2022 | English | 1170 | disinformation meme detection | | | ✓ | |
| ELEMENT [92] | 2022 | English | 7912 | unethical meme detection | ✓ | | ✓ | ✓ |
| Emoffmeme [93] | 2023 | Hindi | 7500 | offensive meme detection | ✓ | | | |
| Ext-Harm-P [94] | 2023 | English | 4446 | harmful reference detection | ✓ | | ✓ | |
| Facebook Hateful Memes [53] | 2022 | English | 9540 | hateful meme detection | | | ✓ | |
| FAME dataset [95] | 2020 | English | 1000 | fake meme detection | ✓ | | | |
| Fine grained HM [96] | 2020 | English | 9540 | fine-grained hateful meme detection | | | ✓ | ✓ |
| GOAT-Bench [87] | 2021 | English | 6626 | abusive/toxic meme detection | ✓ | ✓ | ✓ | ✓ |
| Harm-C (aka HarMeme) [49] | 2024 | English | 3544 | harmful meme detection | ✓ | | ✓ | |
| Harm-P [97] | 2021 | English | 3552 | harmful meme detection | ✓ | | ✓ | |
| Hate Speech in Pixels [98] | 2021 | English | 5030 | hateful meme detection | ✓ | | ✓ | |
| HatReD [99] | 2019 | English | 3228 | hateful meme explanation | | | ✓ | ✓ |
| HVVMemes [100] | 2023 | English | 7000 | entity roles in harmful memes | ✓ | | ✓ | |
| Indian Political Memes [101] | 2022 | Hindi, English | 1218 | hateful meme detection | ✓ | | | |
| Innopolis Hateful Memes [102] | 2022 | English | 23000 | hateful meme detection | ✓ | | ✓ | ✓ |
| KAU-Memes [103] | 2022 | English | 2582 | offensive meme detection | | | ✓ | |
| Meme-Merge [104] | 2023 | English | 10000 | offensive meme detection | ✓ | | ✓ | |
| Memotion 1 [16] | 2023 | English | 10000 | offensive meme detection | ✓ | | | |
| Memotion 2 [105] | 2020 | English | 10000 | offensive meme detection | | ✓ | ✓ | |
| MET-Meme [106] | 2022 | English, Chinese | 10045 | offensive meme detection | ✓ | | ✓ | |
| Misogynistic-MEME [107] | 2022 | English | 800 | misogynous meme detection | | | ✓ | |
| MultiBully [108] | 2022 | Hindi, English | 5854 | cyberbullying meme detection | | | ✓ | |
| MultiBully-Ex [80] | 2022 | Hindi, English | 3222 | cyberbullying meme explanation | | | ✓ | |
| MAMI [109] | 2024 | English | 10000 | misogynous meme detection | | ✓ | ✓ | ✓ |
| MultiOFF [110] | 2022 | English | 743 | offensive meme detection | | | ✓ | |
| Pol_Off_Meme [111] | 2020 | Hindi, English | 7500 | offensive meme detection | ✓ | | | |
| SemEval-2021 Task 6 [19] | 2024 | English | 950 | propagandistic technique detection | | | ✓ | |
| TamilMemes [112] | 2021 | Tamil | 2969 | troll meme detection | | ✓ | ✓ | |
| TrollsWithOpinion [23] | 2020 | English | 8881 | troll meme detection | ✓ | | | |

augmenting previously introduced datasets: Ext-Harm-P [94] is derived from Harm-P [97], while Fine grained HM [96] and HatReD (Hateful meme with Reasons Dataset) [99] are derived from Facebook Hateful Memes (HM) [114]. Moreover, other datasets were formed by merging existing datasets: Meme-Merge [104] is a merge of MET-Meme [106], Memotion 1 [16], and Memotion 2 [105], while GOAT-Bench [87] was created by amalgamating data from several sources including Facebook HM [114], MAMI [109], MultiOFF [112], Harm-C [49], and Harm-P [97]. This analysis highlights the dynamic and varied meme distribution and consumption landscape across numerous online
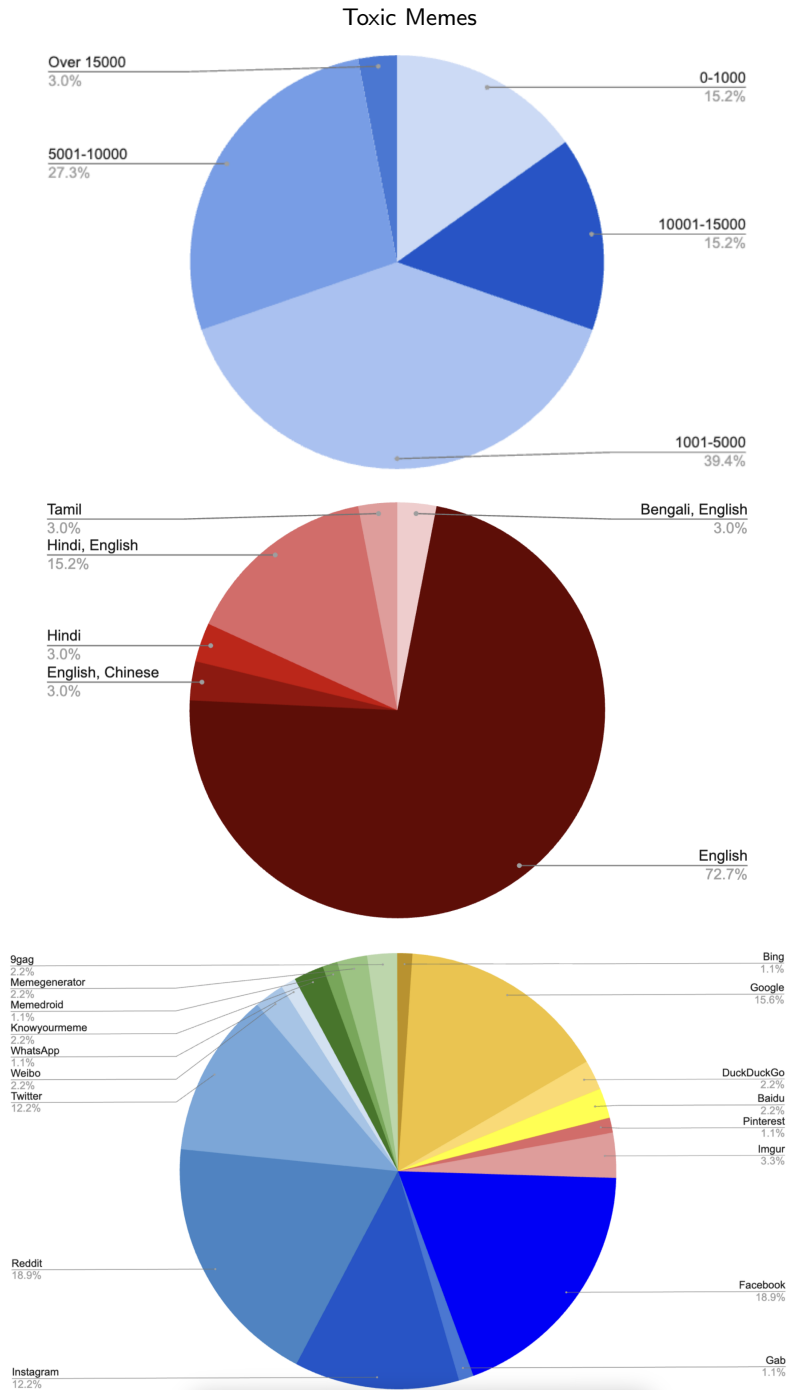
**Figure 4:** Top: Pie chart showing the distribution of datasets based on the number of memes they contain. Approximately 40% of datasets contain between 1000 to 5000 memes, followed by nearly 30% with 5000 to 10,000 memes. Less than 3% of the datasets have over 15,000 memes. Middle: Pie chart illustrating the distribution of dataset languages. English dominates with nearly 75% of the datasets exclusively in English, while the remaining datasets include Hindi, Bengali, Tamil, and code-mixed memes. Bottom: Pie chart displaying the sources of memes in the datasets. Social media platforms contribute the most (blue), followed by meme-specific sources (green), search engines (yellow), and image-hosting platforms (red). Over half of the datasets include memes from social media.

platforms. It also reveals a growing trend toward reusing existing datasets. Also, it emphasizes the evolution of dataset creation in meme research, showcasing how existing resources are being adapted to address new research questions.

## 6.2. Dataset Annotations

We carefully examined the annotation guidelines of each dataset to identify the types of labels assigned to toxic memes and the computational aspects addressed. Table 6 provides a comprehensive summary of the annotated features of meme images across multiple datasets. Each dataset is represented by a column, and each feature is represented by a row. A check mark indicates the presence of labels for a specific feature in each dataset, serving as a checklist to show which aspects of meme images are covered by each dataset. In this process, we found that different datasets employ various annotation methods for the same aspects of memes, including binary, single-label multi-class, or multi-label multi-class schemes, as described below.

**Abusiveness** is indicated in the BanglaAbuseMeme dataset [13] with each meme labeled as either abusive or not abusive in a binary (yes/no) format.

**Aggressiveness** is evaluated in the Misogynistic-MEME (MM) dataset [107] using a binary approach.

**Attack types** are annotated in various datasets. Fine-grained HM [96] adopts a multi-label multi-class approach, delineating different attack types within memes, including dehumanizing, inferiority, inciting violence, mocking, contempt, slurs, and exclusion. Similarly, the Multimedia Misogyny Dataset (MAMI) [109] provides attack type labels specific to misogyny in a multi-label multi-class format, encompassing general misogyny, shaming, stereotype, objectification, and violence.

**Bullying** categorization is found in both the MultiBully dataset [108] and its extension, MultiBully-Ex [80], with memes labeled as either bully or non-bully using binary labels.

**Disinformation** presence is annotated in the DisinfoMeme dataset [91] through a binary approach, with memes labeled as either containing disinformation or not (yes/no).

**Emotion** is a feature of memes that is labeled across datasets in various ways. Some datasets such as MultiBully [108], and MultiBully-Ex [80] utilize a single-label multi-class approach, tagging memes with emotions like joy, sadness, fear, surprise, anger, disgust, anticipation, trust, or ridicule. Similarly, MET-Meme [106] employs a single-label multi-class system, categorizing memes with emotions such as happiness, love, anger, sorrow, fear, hate, or surprise. Conversely, Emoffmeme [93] and Pol_Off_Meme [111] utilize a multi-label multi-class approach, encompassing annotations for fear, neglect, irritation, rage, disgust, nervousness, shame, disappointment, envy, suffering, sadness, joy, pride, and surprise. In contrast, Memotion 1 [16] and Memotion 2 [105] offer annotations for sarcastic, humorous, motivational, and offensive emotions within a multi-label multi-class framework. The term **Sentiment** is used to describe this same aspect of memes in Emoffmeme [93] which uses a multi-label multi-class paradigm for including annotations for fear, neglect, irritation, rage, disgust, nervousness, shame, disappointment, envy, suffering, sadness, joy, pride, and surprise. In contrast, most other datasets utilize sentiment labels in a Likert scale, single-label multi-class approach. For example, Memotion 1 [16] and Memotion 2 [105] categorize sentiments as very positive, positive, neutral, negative, and very negative. Similarly, BanglaAbuseMeme [13], MultiBully [108], and MultiBully-Ex [80] follow a single-label multi-class scheme with simpler labels: positive, neutral, and negative. Furthermore, the Derogatory Facebook-Meme dataset [90] employs a binary annotation to indicate negative sentiment (yes/no).

**Explanation** annotations are present in a couple of datasets: HatReD (Hateful meme with Reasons Dataset) [99] includes human-provided textual explanations. MultiBully-Ex [80] provides both textual explanations as textual rationales (highlighted words or phrases) and visual explanations via visual masks (image segmentation) for its categorization of memes as harmful.

**Fake/misattribution** is annotated in the FAME dataset [95], categorizing memes as either fake or real (binary).

**Harmfulness** is annotated in varied ways across the datasets. The GOAT-Bench dataset [87] uses a binary system to indicate harmfulness, labeling each instance as either harmful or not (yes/no). Several other datasets, including Harm-P [97], Harm-C (also known as HarMeme) [49], MultiBully [108], and MultiBully-Ex [80], categorize the degree of harmfulness using a single-label multi-class system. These datasets label content as very harmful, partially harmful, or harmless, providing a more nuanced understanding of harmfulness. Additionally, the Ext-Harm-P dataset [94] focuses on the harmfulness of the reference to a social entity. It employs a binary system, labeling references as either harmful or harmless, thus specifically targeting the impact on social entities.

**Hate speech** is annotated in a varied way across different datasets. Most datasets, including Hate Speech in Pixels [98], Facebook Hateful Memes (HM) [53], Fine grained HM, [96], CrisisHateMM [89], Innopolis Hateful Memes [102], GOAT-Bench [87], and Derogatory Facebook-Meme [90], use a binary labeling system to identify hate speech, categorizing content simply as either hateful or not (yes/no). However, the Indian Political Memes (IPM) dataset [101] employs a more nuanced approach with single-label multi-class annotations among non-offensive, hate-inducing, and satirical. Additionally, CrisisHateMM [89] includes an extra dimension by labeling the direction of hate. This dataset

differentiates between directed hate and undirected hate, using a binary system to specify whether the hate is aimed at a specific target or more general in nature.

**Humour** is assessed in Memotion 1 [16] and Memotion 2 [105] datasets on a continuum scale, spanning from not funny to hilarious, within a single-label multi-class framework.

**Intention** is addressed in a single-label multi-class format in the MET-Meme dataset [106], offering labels such as interactive, expressive, purely entertaining, offensive, or other.

**Irony** is annotated as a binary label (yes/no) in the Misogynistic-MEME dataset [107], while **Sarcasm** is similarly represented in binary format (yes/no) in BanglaAbuseMeme [13], MultiBully [108], MultiBully-Ex [80], and GOAT-Bench [87]. In contrast, sarcasm in Memotion 1 [16] and Memotion 2 [105] is classified along a continuum, with labels ranging from "not sarcastic" to "very twisted," within a single-label multi-class framework.

**Metaphors** in toxic memes are addressed in the MET-Meme dataset [106] annotates metaphor, which includes binary labels for **metaphorical expression** (metaphorical/literal), and also employs a single-label multi-class approach to classify metaphor types into text dominant, image dominant, or complementary.

**Misogyny** is labeled in a binary manner (yes/no) across multiple datasets, including the Misogynistic-MEME (MM) dataset [107], Multimedia Misogyny Dataset (MAMI) [109], and GOAT-Bench [87].

**Modality Class** in Memotion 2 [105] employs a single-label multi-class approach to label memes based on their modality, allowing them to be categorized as either image and text, image only, or text only.

**Motivation** is annotated in Memotion 2 [105] and Memotion 1 [16] datasets, distinguishing memes as either motivational or not motivational in a binary manner.

**Offensiveness** is assessed differently across datasets. Some datasets, such as MultiOFF [110], AOMD Gab [15], AOMD Reddit [15], Pol_Off_Meme [111], Emoffmeme [93], KAU-Memes [103], TrollsWithOpinion [23], and GOAT-Bench [87], use a binary labeling system (yes/no) to indicate offensiveness. However, in other datasets, such as Memotion 1 [16], Memotion 2 [105], and MET-Meme [106], offensiveness is assessed on a single-label multi-class scale, ranging from not offensive to hateful offensive. Meme-Merge [104] also adopts a degree-based labeling system, with offensiveness categorized from non-offensive to very offensive. Additionally, Pol_Off_Meme [111], and Emoffmeme [93] include binary labels for the explicitness of offensiveness, distinguishing between implicit and explicit offensiveness.

**Opinion manipulation** is addressed in TrollsWithOpinion [23], which employs a binary labeling approach, categorizing memes as either involving opinion manipulation or not. Additionally, it provides labels for different opinion manipulation types such as political, product, or other, utilizing a single-label multi-class manner to capture these distinctions.

**Political attributes** are identified in Pol_Off_Meme [111] through binary classification, distinguishing memes as either political or not political.

**Profanity** labels are included in the Derogatory Facebook-Meme dataset [90], while **vulgarity** labels are present in BanglaAbuseMeme [13]. Both are annotated in a binary framework of yes/no.

**Propagandistic techniques** are annotated in SemEval-2021 Task 6 [19] using a multi-label multi-class approach. This encompasses a wide range of techniques, including Loaded Language, Name Calling/Labeling, Smears, Doubt, Exaggeration/Minimization, Slogans, Appeal to Fear/Prejudice, and more. Each meme can be assigned multiple labels corresponding to the propagandistic techniques it employs.

**Target** identification in toxic memes is a common aspect annotated across datasets. While Derogatory Facebook-Meme [90] employs a binary label to indicate whether a meme is targeted or not, most datasets annotate the target in a single-label multi-class manner. Common categories include individual, community, organization, and society, as seen in Harm-P [97], Harm-C (also known as HarMeme) [49], and CrisisHateMM [89]. Some datasets provide more specific community targets: BanglaAbuseMeme [13] uses single-label multi-class labeling with options like gender, religion, national origin, individual, political, social sub-groups, and others. Additionally, Fine-grained HM [96] employs a multi-label multi-class approach, including labels for protected categories such as religion, race, sex, nationality, and disability. Relatedly, HVVMemes [100] utilizes a single-label multi-class framework to annotate each entity's role, with options including villain, victim, hero, or other.

**Troll** is annotated as a dimension in both TamilMemes [112] and TrollsWithOpinion [23] using a binary labeling.

**Unethical** is an aspect of memes annotated in the ELEMENT dataset [92] using a binary (yes/no) labeling system.

**Table 6**

Overview of the various aspects addressed in the identified datasets. Each row represents a dataset, and each column represents an aspect. Tick marks indicate that the dataset addresses the corresponding aspect.

| | abusiveness | aggressiveness | attack type | bullying | disinformation | emotion/sentiment | explanation | fake/misattribution | harmfulness | hate speech | humour | intention | irony/sarcasm | metaphor | misogyny | modality-class | motivation | offensiveness | opinion manipulation | political attributes | profanity/vulgarity | propagandistic technique | target | troll | unethical |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AOMD Gab [15] | | | | | | | | | | | | | | | | | | ✓ | | | | | | | |
| AOMD Reddit [15] | | | | | | | | | | | | | | | | | | ✓ | | | | | | | |
| BanglaAbuseMeme [13] | ✓ | | | | | ✓ | | | | | | | ✓ | | | | | | | | ✓ | | ✓ | | |
| CrisisHateMM [89] | | | | | | | | | ✓ | ✓ | | | | | | | | | | | | | ✓ | | |
| Derogatory Fb-Meme [90] | | | | | | ✓ | | | | ✓ | | | | | | | | | | | ✓ | | ✓ | | |
| DisinfoMeme [91] | | | | | ✓ | | | | | | | | | | | | | | | | | | | | |
| ELEMENT [92] | | | | | | | | | | | | | | | | | | | | | | | | | ✓ |
| Emoffmeme [93] | | | | | | ✓ | | | | | | | | | | | | ✓ | | | | | | | |
| Ext-Harm-P [94] | | | | | | | | | ✓ | | | | | | | | | | | | | | | | |
| Facebook Hateful Memes [53] | | | | | | | | | | ✓ | | | | | | | | | | | | | | | |
| FAME dataset [95] | | | | | | | | ✓ | | | | | | | | | | | | | | | | | |
| Fine grained HM [96] | | | ✓ | | | | | | | ✓ | | | | | | | | | | | | | ✓ | | |
| GOAT-Bench [87] | | | | | | | | | ✓ | ✓ | | | ✓ | | ✓ | | | ✓ | | | | | | | |
| Harm-C (aka HarMeme) [49] | | | | | | | | | ✓ | | | | | | | | | | | | | | ✓ | | |
| Harm-P [97] | | | | | | | | | ✓ | | | | | | | | | | | | | | ✓ | | |
| Hate Speech in Pixels [98] | | | | | | | | | | ✓ | | | | | | | | | | | | | | | |
| HatReD [99] | | | | | | | ✓ | | | ✓ | | | | | | | | | | | | | | | |
| HVVMemes [100] | | | | | | | | | | | | | | | | | | | | | | | ✓ | | |
| Indian Political Memes [101] | | | | | | | | | | ✓ | | | | | | | | | | | | | | | |
| Innopolis Hateful [102] | | | | | | | | | | ✓ | | | | | | | | | | | | | | | |
| KAU-Memes [103] | | | | | | | | | | | | | | | | | | ✓ | | | | | | | |
| Meme-Merge [104] | | | | | | | | | | | | | | | | | | ✓ | | | | | | | |
| Memotion 1 [16] | | | | | | ✓ | | | | | ✓ | | ✓ | | | | ✓ | ✓ | | | | | | | |
| Memotion 2 [105] | | | | | | ✓ | | | | | ✓ | | ✓ | | | ✓ | ✓ | ✓ | | | | | | | |
| MET-Meme [106] | | | | | | ✓ | | | | | | ✓ | | ✓ | | | | ✓ | | | | | | | |
| Misogynistic-MEME [107] | | ✓ | | | | | | | | | | | ✓ | | ✓ | | | | | | | | | | |
| MultiBully [108] | | | | ✓ | | ✓ | | | ✓ | | | | ✓ | | | | | | | | | | | | |
| MultiBully-Ex [80] | | | | ✓ | | ✓ | ✓ | | ✓ | | | | ✓ | | | | | | | | | | | | |
| MAMI [109] | | | ✓ | | | | | | | | | | | | ✓ | | | | | | | | | | |
| MultiOFF [110] | | | | | | | | | | | | | | | | | | ✓ | | | | | | | |
| Pol_Off_Meme [111] | | | | | | ✓ | | | | | | | | | | | | ✓ | | ✓ | | | | | |
| SemEval-2021 Task 6 [19] | | | | | | | | | | | | | | | | | | | | | | ✓ | | | |
| TamilMemes [112] | | | | | | | | | | | | | | | | | | | | | | | | ✓ | |
| TrollsWithOpinion [23] | | | | | | | | | | | | | | | | | | ✓ | ✓ | | | | | ✓ | |

**Table 7**
Meme toxicity-related terms as utilized in the surveyed literature focusing on computationally addressing toxic memes. It outlines each term alongside its corresponding definition provided by the papers, with the last column indicating the surveyed papers and preprints primarily focused on that specific term.

| Toxicity Type | Definition | Focus of |
|---|---|---|
| Abusive | Used interchangeably with "harmful" by [13, 115], who describe memes used by bad actors to threaten and abuse individuals or specific target communities. These memes contain words that target individuals or different protected communities, implicitly containing hateful, harmful, and antisemitic content. [116] uses it interchangeably with "offensiveness." | [13, 115, 116] |
| Cyberbullying | Defined as any communication that disparages an individual based on characteristics such as color, gender, race, sexual orientation, ethnicity, nationality, or other features [117]. This definition, used by [80], builds upon the work by [14], who consider cyberbullying as an antisocial activity where victims are targeted with malicious behavior, including posting cruel comments and messages without fear of being identified. It appears that both [14] and [115] consider cyberbullying as a subset of harmful behavior. | [14, 80, 118] |
| Derogatory | Term used by [90] with no explicit definition. Vaguely defined as posts that convey a derogatory notion about a recognized individual (or person) of the country; malicious content about a political, spiritual, and cultural entity; and/or posts that include hateful and negative sentiments of sentences. | [90] |
| Disinformation | [91] defines disinformation memes as those designed to actively spread inaccurate information. They differentiate instances that criticize misinformation from those actively spreading inaccurate information. | [18] |
| Fake/Misattribution | Defined by [95] as a type of disinformation meme, these contain messages, fabricated or otherwise, falsely attributed to specific individuals. Such memes could be deployed against political opponents during smear campaigns. | [95] |
| Harmful | [49] see harmful memes as those that "have the potential to cause harm to an individual, an organization, a community, or the society more generally. Here, harm includes mental abuse, defamation, psycho-physiological injury, proprietary damage, emotional disturbance, and compensated public image. [...] Harmful is a more general term than offensive and hateful: offensive and hateful memes are harmful, but not all harmful memes are offensive or hateful" (page 4). As defined in [94], harmful memes encompass various forms of harm expressed overtly or subtly towards target entities, socio-cultural or political ideologies, beliefs, principles, or doctrines associated with them. The harm may manifest as abuse, offense, disrespect, insult, demeaning, or disregard towards the target entity or its affiliations. It can also include more subtle attacks such as mockery or ridicule of a person or an idea. | [11, 12, 49, 79, 81, 94, 97, 100, 119, 120, 121, 122, 123, 124, 125, 126, 127] |
| Hateful | In line with [114], hateful memes are characterized as direct or indirect attacks on individuals based on protected characteristics such as ethnicity, race, nationality, immigration status, religion, caste, sex, gender identity, sexual orientation, disability, or disease. An attack is defined as containing violent or dehumanizing speech, statements of inferiority, or calls for exclusion or segregation. Additionally, mocking hate crimes is classified as hate speech. Exceptions to this definition include attacks on individuals/famous people not based on protected characteristics and attacks on groups perpetrating hate, such as terrorist groups. Detection of hate speech in memes often requires nuanced understanding of societal norms and context. More recently, fine-grained hateful types [128] have been studied, focusing on Protected Categories (PC) such as race, disability, religion, nationality, and sex, along with different Attack Types (AT) including contempt, mocking, statements of inferiority, slurs, exclusion, dehumanizing content, and incitement to violence. | [7, 8, 9, 42, 85, 86, 96, 98, 99, 101, 102, 114, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186] |
| Misogynous | Misogyny is a form of hate against women, and misogynous memes manifest different expressions of hate directed towards women, encompassing a broad spectrum [187]. These misogynous manifestations may be categorized into four main types. | [46, 113, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199] |
| Offensive | In some of the earliest meme analyses, "offensive" was conceived as a type of "emotion" alongside humor, sarcasm, and motivation, simply defined as something that aims to torment or disturb people [16]. The work introducing one of the most widely used datasets, multiOFF [110], follows [200]'s definition of offensive content as intending to upset or embarrass people by being rude or insulting. Offensive or abusive content on social media can be explicit or implicit [61], and could be classified as explicitly offensive or abusive if it is unambiguously identified as such. For example, it might contain racial, homophobic, or other offending slurs. In the case of implicit offensive or abusive content, the actual meaning is often obscured by the use of ambiguous terms, sarcasm, lack of profanity, hateful terms, or other means. | [15, 16, 88, 93, 103, 111, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221] |
| Propaganda | Defined by [19] as a form of communication to influence the opinions or actions of people towards a specific goal, achieved through well-defined rhetorical and psychological devices. | [19, 20, 21, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232] |
| Troll | Defined by [23] as a meme that is often provocative, distractive, digressive, or off-topic with the intent to demean or offend particular people, groups, or races, containing either I) Offensive text and non-offensive images; II) Offensive images with non-offensive text; or III) Sarcastically offensive text with non-offensive images, or sarcastic image with offensive text. | [22, 233, 234, 235] |
| Unethical | As defined by [92], unethical memes are those deemed "inconsistent with human values," failing to comply with ethical norms. They specify ethical memes as multi-modal units consisting of an image and embedded text aiming to promote fairness, justice, harmony, security, avoidance of bias, discrimination, privacy, and prevention of information leakage. While harmful and hateful memes are considered unethical, not all unethical memes are necessarily harmful or hateful. Detection of unethical content in memes is particularly challenging due to its often deeply implicit nature. Ethical memes focus not only on interpersonal principles but also on societal, human, and environmental relationships. | [92] |

## 7. Defining Meme Toxicities: A Taxonomy

Based on our investigation of the datasets being developed and utilized in the field, we discovered a diverse array of terms used to describe toxic memes, indicating that multiple types of toxicity are being explored computationally. To identify and classify various types of meme toxicity for computational analysis, we meticulously individually scrutinized each of the 158 research papers to extract explicit terms denoting the types of toxicity they addressed.

### 7.1. Meme Toxicities Identification and Definitions

We identified 12 meme toxicity terms: *abusive*, *cyberbullying*, *derogatory*, *disinformation*, *fake*, *harmful*, *hateful*, *misogynous*, *offensive*, *propaganda*, *troll*, and *unethical*. We then compiled the definitions associated with each term and harmonized them into Table 7. Our analysis of toxicity-related terms showed that some terms were used interchangeably in the literature. For example, *abusive* was used synonymously with *offensive* by [116] and with *harmful* by [13, 115]. Additionally, terms like *derogatory* lacked clear definitions and were vaguely described using other toxicity descriptors. Figure 5 provides an overview of research attention across various meme toxicity categories. Notably, there is a clear imbalance, with a significant focus on hateful memes. Nearly half of the research papers concentrate on this aspect, highlighting the prevalence of hate speech in online discourse. This trend indicates the significant impact of the hateful memes challenge proposed by Kiela et al. (2020) [53], as evidenced by the widespread use of its dataset and methodologies in subsequent research. Our analysis also reveals a lack of research on less prominent toxicity categories like troll, derogatory, and disinformation memes, suggesting gaps in understanding and addressing these forms of toxicity. No surveyed literature dealt with categories flagged in previous work [10], e.g., self-harm and exploitation.
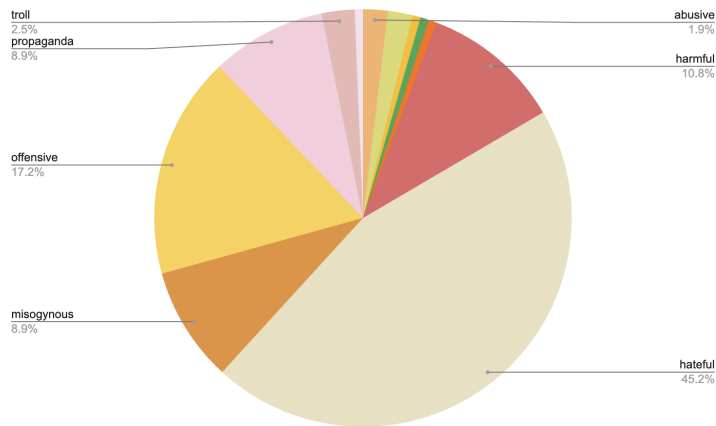


**Figure 5:** Proportion of papers focused on different meme toxicity categories, based on authors' explicit labels.

### 7.2. A Taxonomy of Meme Toxicities

In examining toxicity term definitions, we identified explicit and implicit taxonomical relationships among certain terms. For example, offensive, hateful, troll, and cyberbullying memes are typically defined as inherently harmful. Fake memes are seen as a subset of disinformation memes, while misogynous memes are considered a subtype of hateful memes targeting protected categories, such as sex. Recognizing the importance of clarifying these relationships, we aimed to establish a coherent taxonomy. We referred to the existing meme toxicity taxonomy by Sharma et al. (2022) [10], which included overlooked meme toxicities like self-harm and exploitation (e.g., adult sexual services, child sexual abuse, and scams). However, significant discrepancies emerged. Their taxonomy overlooked certain toxicities like disinformation, fake, and derogatory memes and lacked the concept of unethical memes, a superclass encompassing harmful memes and more. Additionally, the rationale for placing certain terms under specific macrocategories was unclear. For example, 'doxing'—defined as publicly publishing someone's private information as punishment or revenge—was categorized under hate rather than cyberbullying, even though it targets individuals rather than protected communities. Most critically, the taxonomy conflated two dimensions of meme toxicities: the type of toxicity (e.g.,
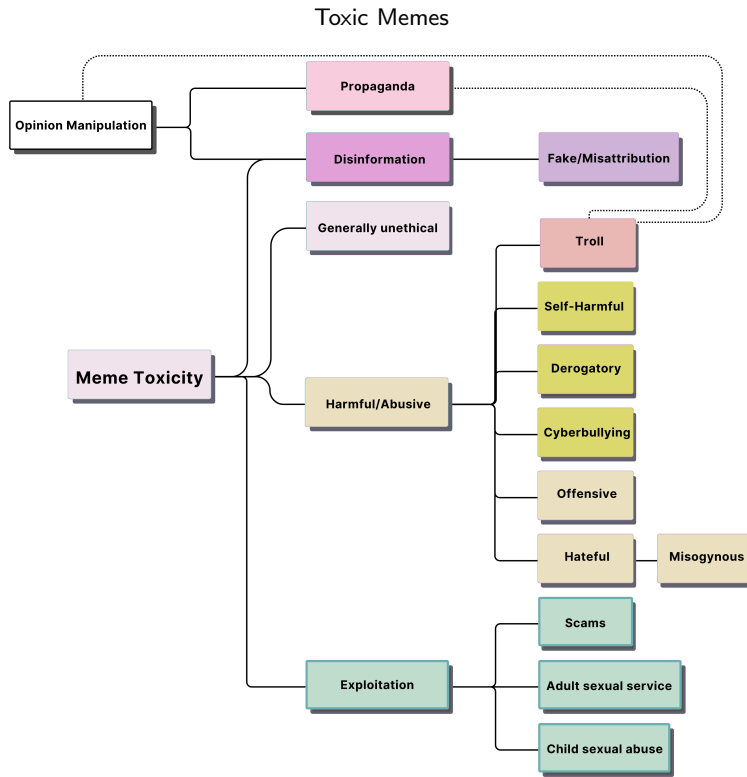
**Figure 6:** The taxonomy for meme toxicities that we propose is inspired by the taxonomy presented in [10], while addressing discrepancies, enhancing taxonomical clarity, and including the most recent types of toxicities being computationally studied.
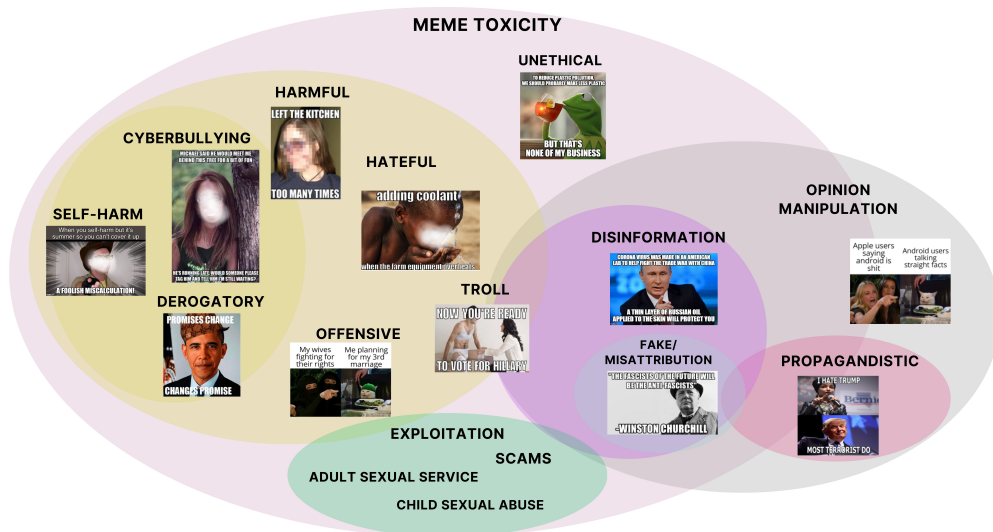


**Figure 7:** Venn diagram illustrating taxonomical relations and fuzzy boundaries between meme toxicities, along with examples. Memes in the exploitation category are not depicted due to the absence of corresponding datasets. The diagram is color-coded to match our taxonomy. Caution: Depicted memes contain toxic content, potentially inducing psychological distress. Viewer discretion is advised. These images do not reflect the views or endorsements of the authors.

hateful, disinformation) and the techniques or attacks used (e.g., identity attack, identity misrepresentation). We believe separating these dimensions is vital, as certain techniques can be used across different toxicities (see Section 8).

Building on the taxonomy by Sharma et al. (2022) [10], we refined it to address identified discrepancies and nuances, enhancing its adaptability for future meme toxicity studies. Our revised taxonomy, shown in Figure 6, is based on definitions from the literature, with each category as a subcategory of another. We found that certain toxicities, such as

disinformation and propaganda, stem more from opinion manipulation, while many computationally studied toxicities fall under harmful/abusive, characterized by intent to abuse, demean, offend, or exploit. Our review of the state of the art revealed complex relationships between categories and subcategories of meme toxicities, showing a level of nuance beyond what a taxonomy alone can convey. To visualize these relationships, we developed a Venn diagram (see Figure 7) that uses color coding to enhance clarity and illustrate the interconnectedness of various toxicity labels. This diagram provides insight into the complexity of meme toxicity relationships, highlighting the nuanced and sometimes fuzzy boundaries between different toxicities. While the taxonomical approach remains crucial for computational studies, the Venn diagram reflects the real-world complexity of these boundaries.

## 8. Dimensions of Meme Toxicities: Target, Intent, Tactic

While refining the taxonomy, we noticed that meme toxicity definitions often focused on three main aspects: the entity to whom the toxicity is directed, the intention behind a meme's creation or dissemination, and the rhetorical strategies employed to convey the toxic narratives. While terminology varied and their interrelation wasn't explicitly addressed, we identified these as three dimensions of meme toxicities: intent, target, and tactic(s) (see Figure 8). These dimensions could provide a structured framework for understanding and addressing meme toxicity.
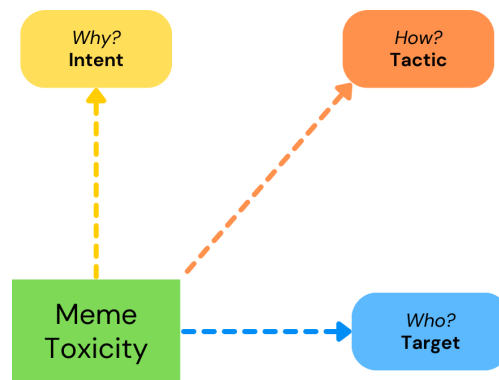


**Figure 8:** The toxicity of memes is a complex phenomenon with multiple dimensions, represented by dotted edges. In this survey, we have identified three content-based dimensions of meme toxicity: target, intent, and tactic. These dimensions address the fundamental questions: *who* is the toxicity directed towards, *why* the toxic meme is shared (i.e., the underlying goal), and *how* the toxicity is manifested and conveyed. Other dimensions also exist and should be studied, potentially answering questions like *where* the meme was posted and *when* it was shared.

### 8.1. Target

The most explicitly researched dimension of toxic memes is the *intended target* of the meme. In memes, targeting involves the harmful referencing of social entities [49, 94]. This harm can manifest in various forms such as mental abuse, psycho-physiological injury, proprietary damage, emotional disturbance, or damage to public image, based on the background of the entity (bias, social background, educational background, etc.) by the meme author. It is important to note that memes can reference social entities without causing harm; simply being referenced does not constitute being the target. Non-harmful referencing in memes includes benign mentions (or depictions) via humor, limericks, harmless puns, or any content that doesn't cause distress [94]. Research has focused on detecting the targets of, particularly, harmful memes [49, 89, 94, 97], based on the following target types:

*Individual* Toxic memes often target single individuals, typically celebrities like politicians, actors, artists, scientists, or environmentalists, as noted by Sharma et al. [94]. These individuals include figures such as Donald Trump, Joe Biden, Vladimir Putin, Hillary Clinton, Barack Obama, Chuck Norris, Greta Thunberg, and Michelle Obama. Other studies have noted specific subtypes of individuals as targets of certain meme categories: **famous or known individuals** are targeted in misattribution memes, with 20 well-known figures commonly featured in fake quote memes [95]. Additionally, **meme creators/posters** are implicitly recognized as targets of self-harmful memes, while **meme viewers/receivers** are implicitly targeted in scam memes. Individual **children** may be targeted in memes related to child sexual abuse, and **adults** are targeted in memes related to adult sexual service [10].

*Organization* Groups with specific purposes, such as businesses, government departments, companies, institutions, or associations are targeted, e.g., Facebook, the WTO, and the Democratic Party.

*Community* These targets are social units defined by shared personal, professional, social, cultural, or political attributes, including religious views, country of origin, or gender identity. These communities may manifest in geographical areas or virtual spaces facilitated by online communication platforms. Within this context, **Protected Category (PC) Communities** represent a focal point, as identified in Fine-grained version of Hateful Menes [96]. This dataset introduced fine-grained labels for protected categories, including **race**, **ethnicity**, **national origin**, **disability**, **religious affiliation**, **caste**, **sexual orientation**, **sex**, **gender identity**, and **serious disease**. However, the dataset provided primarily covers race, disability, religion, nationality, and sex [128]. A similar target community labeling scheme was done in the BanglaAbuseMeme dataset[13], in which communities are seen as related to gender, religion, national origin, individual, political, social sub-groups, or others.

*Society* The target is the entire societal fabric, e.g., when memes promote conspiracies, harming the general public.

## 8.2. Intent

Multimodal intent recognition plays a crucial role in understanding human language within real-world multimodal scenes. However, existing intent recognition methods face limitations in fully utilizing multimodal information, largely due to the constraints of benchmark datasets containing only text information [236]. In the case of memes, intent refers to the underlying purpose or motivation behind their creation or dissemination, and it has been sometimes annotated [106]. However, the labels offered to choose between are interactive, expressive, purely entertaining, offensive, or other. This implies that the only intent connected to toxicity is the *offensive* intent, defined as the intent for discrimination, satire, and abuse directed towards others' occupation, gender, appearance, or personality [106]. While it's widely assumed in the literature we surveyed that the intent of all toxic memes is unethical, involving the conveyance of information contrary to ethical societal values like fairness, justice, and privacy [92], specific types of toxic intents unique to toxic memes have not, to the best of our knowledge, been studied. However, as based on the collected definitions of meme toxicities, there is also an acknowledgment of specific intents, such as the intent to harm, misinform, or exploit, that partly characterize the type of toxicity that the meme carries. comprehensive research explicitly identifying and comparing these different intents of toxic memes to each other is lacking. Taking a step in this research direction, based on our survey of meme toxicities, we have identified three subtypes of intents in toxic memes:

*Harm/Abuse* The intent is to harm, offend, demean, abuse, or insult social entities [97].

*Disinform* The intent is to manipulate opinions by actively spreading false or inaccurate information [18, 95].

*Exploit* This intent is mentioned, but not explicitly defined, by Sharma et al (2022) [10], who identified memes attempting to exploit people, including children for sexual abuse, adults for sexual services, or viewers for scams.

## 8.3. Tactic(s)

We also noticed a growing emphasis on automatically identifying the rhetorical techniques or tactics used by memes to convey toxicity. This focus is justified as memes employ a wide range of rhetorical strategies, persuasion techniques, and tactics to convey messages within their narratives [19]. This aspect, explored through attack types [96, 113, 128], persuasion/propagandistic techniques [19], or entities' roles in harmful memes [121], represents a third dimension implicitly or explicitly mentioned in many papers, but these terms have not, to our knowledge, been collectively examined.

*Attack types* delineate how memes, particularly those with hateful connotations, target individuals, groups, or ideas, and serve to categorize how memes express their message [96, 128]. Attack types were defined within the context of the WOAH 5 shared task on fine-grained hateful memes detection [96]. Attack types have also been specialized to address specific forms of hateful messaging in the realm of misogynous memes [113]:

- **Dehumanization**: Explicitly or implicitly describing or presenting a group as subhuman [96]. It includes **misogynous objectification**, i.e., treating or regarding women as objects, devoid of agency/humanity [113].

---

- **Statements of Inferiority**: Claiming that a group is inferior, less worthy or less important than either society in general or another group [96].

- **Violence**: Explicitly or implicitly calling for harm to be inflicted on a group, including physical attacks [96]. It includes **misogynous violence**, i.e., the incitation or direct expression of acts of violence against women [113].

- **Mocking/Shaming**: Making jokes about, undermining, belittling, or disparaging a group [96]. It includes **misogynous shaming**, defined as the practice of criticising women who violate expectations of behaviour and appearance regarding issues related to gender typology (such as "slut shaming") or related to physical appearance (such as "body shaming"). This category focuses on content that seeks to insult and offend women because of some characteristics of their body or personality [109].

- **Expression of Contempt**: Expressing intensely negative feelings or emotions about a group [96].

- **Slurs**: Using prejudicial terms to refer to, describe or characterise a group [96].

- **Calls for Exclusion**: Advocating, planning or justifying the exclusion or segregation of a group from all of society or certain parts [96].

- **Misogynous Stereotyping**: Propagating generalized beliefs about women across various contexts, including societal roles, personalities, and behaviors [113]. Stereotypes are fixed, conventional ideas or set of characteristics assigned to a woman, a meme can use an image of a woman according to her role in the society (role stereotyping), or according to her personality traits and domestic behaviours (gender stereotyping) [109].

*Persuasion or propagandistic techniques* encompass various methods employed to influence the perception of a meme's audience. These methods include selective editing of images or text, framing narratives to elicit specific emotional responses, and employing symbols and motifs associated with particular ideologies or agendas. One significant study [19] identifies key propaganda techniques that serve as shortcuts in the argumentation process of memes specifically, leveraging audience emotions or logical fallacies to influence perception. Notably, the presence of these techniques does not inherently categorize a meme as propagandistic; rather, memes are only annotated based on the propaganda techniques they contain.

- **Loaded language**: Using emotionally charged words and phrases to influence the audience.

- **Name calling or labeling**: Assigning labels to entities that evoke strong reactions from the target audience.

- **Doubt**: Questioning the credibility of someone or something.

- **Exaggeration / Minimization**: Representing something in an excessive or diminished manner.

- **Appeal to fear / prejudices**: Instilling anxiety or panic to build support for an idea.

- **Slogans**: Brief and striking phrases acting as emotional appeals.

- **Whataboutism**: Discrediting an opponent's argument by charging them with hypocrisy.

- **Flag-waving**: Appealing to strong national or group sentiments to justify or promote an action or idea.

- **Straw man**: Substituting an opponent's proposition with a similar one, which is then refuted.

- **Causal oversimplification**: Assuming a single cause or reason when multiple causes exist for an issue.

- **Appeal to authority**: Stating that a claim is true simply because an authority on the issue said it was true.

- **Thought-terminating cliché**: Phrases discouraging critical thought and meaningful discussion on a given topic.

- **Black-and-white fallacy or dictatorship**: Presenting two alternative options as the only possibilities.

- **Reductio ad Hitlerum**: Disapproving an action or idea by suggesting it's popular with hated groups.

- **Repetition**: Repeating the same message to influence acceptance.

- **Obfuscation**, Intentional vagueness, Confusion: Using deliberately unclear words.

- **Presenting irrelevant data (Red Herring)**: Introducing irrelevant material to divert attention.

- **Bandwagon**: Persuading the target audience to join a course of action because "everyone else is doing it."

- **Smears**: Attempting to damage someone's reputation through negative propaganda.

- **Glittering generalities**: Using words or symbols that produce a positive image when attached to a person/issue.

- **Appeal to (strong) emotions**: Using emotionally charged images to influence the audience.

- **Transfer**: Evoking an emotional response by projecting qualities of a person, entity, or value onto another.

*Entity roles* are another feature investigated in the rhetorical dimension of meme toxicity, particularly within harmful memes. This task revolves around identifying which entities are glorified, vilified, or victimized within a meme. Framed from the perspective of the meme's author, the objective is to classify, for a given pair of a meme and an entity, whether the entity is depicted as a Hero, Villain, Victim, or falls into another category within that meme [121]:

- **Hero**: Entities portrayed in a positive light, often glorified for their actions or inferred from context.

- **Villain**: Entities depicted negatively, associated with adverse traits like wickedness, cruelty, or hypocrisy.

- **Victim**: Entities shown suffering from the negative consequences of someone else's actions or conveyed implicitly.

- **Other**: Entities that do not fit the categorization of hero, villain, or victim within the context of the meme.

Overall, these three aspects of memes—attack types, persuasion techniques, and entity roles—are all related to the rhetorical strategies used to convey meme toxicity, yet they have been studied separately thus far. However, more relationships may exist between specific attack types, persuasion techniques, and entity roles than previously recognized (see Figure 9). For instance, the attack type *slurs* may be equivalent or strongly related to the persuasion techniques of *name calling* and *smears*, as they all involve the use of prejudicial terms to influence perception. Similarly, the attack type of *contempt* may be equivalent or related to the persuasion technique of *loaded language*, as they both entail the use of strongly charged emotional language. Additionally, there may be specific attack types or persuasion techniques employed to portray certain entities as heroes or villains. Furthermore, as noted in Section 6.2, some toxic meme datasets include labels for additional features that may be part of the rhetorical strategy for conveying toxicity in memes: emotion [16, 80, 93, 105, 106, 108], sentiment [13, 16, 105, 108], humour [51, 105], irony/sarcsam [13, 16, 87, 105, 107, 108], metaphor [106], profanity/vulgarity [13, 90], and more. These connections remain largely unexplored, underscoring the need for a comprehensive examination of these aspects within the rhetorical dimension.

## 8.4. Dimensions' Intersections
*Intent x Target* While examining individual dimensions of meme toxicity is valuable, we propose that simultaneous analysis of multiple dimensions holds untapped potential, especially for tasks like automatic detection and explanation of toxicity. We particularly emphasize the value of identifying intent and target dimensions concurrently, as this can offer insights into the specific types of toxicity in memes (see Figure 10). This visual representation highlights the types of toxicity typically arising from specific combinations of targets and intents. For example, when the intent is to harm and the target is people with disabilities, this may result in ableism, a type of toxicity that is both hateful and harmful. It is worth noting that while we have included toxicities explicitly addressed or mentioned in the surveyed literature, this depiction may not be comprehensive. There may be manifestations of toxicities that combine all types of intent with all types of targets, suggesting that all parts of the two-dimensional plane can potentially be populated. This implies that while these types of memes may exist, they have not yet been categorized under a specific type of toxicity.

**Figure 9:** While relationships between attack types, persuasion/propagandistic techniques, and entity roles have not been thoroughly explored, insights can be gleaned from their definitions, suggesting potential connections, equivalences, and other relationships (illustrated by red dotted lines).



**Figure 10:** Intersections between intent and target dimensions hint at specific types of meme toxicity. The y-axis, highlighted in yellow, categorically presents the three identified intents. Meanwhile, the x-axis, depicted in blue, categorizes target types with varying levels of specificity (from less specific at the bottom to more specific at the top). The green plane illustrates the toxicity landscape observed through the lens of intent and target dimensions, with greener areas indicating more specific toxicity types determined by particular intersections of intent and target.

# 9. Recent Trends and Research Directions

## 9.1. Trend 1: Tackling Cross-Modal Entailments and Reasoning

Until recently, AI models for identifying toxic or hateful speech mainly relied on single-mode classifiers or text-limited datasets (e.g., [61, 66]). However, there's been advocacy in social sciences for a multi-modal critical discourse approach, stressing the importance of understanding meanings constructed through various sign systems, including language and visuals [42, 237, 238]. Memes, characterized by their co-creation of meaning through text and image, present significant hurdles for machine interpretation because of the complexity of multimodal understanding [44, 50, 239]. The Hateful Memes Challenge at NeurIPS 2020 [114] played a key role in revealing deficiencies in recent AI models lacking holistic, multi-modally informed reasoning. The challenge spurred significant efforts to deepen the multi-modal understanding of memes, leading to increased use of multi-modal deep learning models that integrate information from specialized neural networks analyzing specific modalities. This integration, often achieved through fusion techniques, enables a more comprehensive analysis and decision-making process [39]. Despite advancements, the performance of toxic meme detection algorithms remains suboptimal and many memes evade filters due to their complex, multi-modal nature [41].

In their recent study, Polli et al. (2023) [42] emphasize the complex interplay between textual and visual elements in toxic meme interpretation, revealing persistent challenges in automated detection beyond current computational models' capabilities. Drawing from sociosemiotics and critical multimodal studies, they demonstrate that meaning-making in hateful memes defies unimodal determination or basic multimodal fusion used in most computational approaches. Analyzing examples, they illustrate how seemingly harmless elements can combine to create toxic outcomes. For instance, in Figure 11, while both examples feature a non-hateful image, the second example becomes hateful and toxic due to text implicitly objectifying and dehumanizing the child, and thus perpetuating racist stereotypes. This example highlights the complex cross-modal reasoning needed to detect the toxicity of the meme: given the absence of explicit racial references in the text, AI-driven language-based classification systems struggle to identify it as hateful. These findings underscore the shortcomings of classification systems relying solely on text keywords or image features. Consequently, Polli et al. advocate for rejecting the simplistic view of single modalities, and instead advocate for computational models informed by semiotic and multimodal approaches that prioritize multiplicative meanings [240].

In this survey, we observe a rise in the exploration of sophisticated multimodal approaches to tackle the complex interplays between different modes in memes. Researchers are increasingly focusing on understanding the significance of each modality for specific meme types, experimenting with combining and exploring unimodal and multimodal learning and representations. For instance, in [153], hateful meme detection is approached via a multi-task learning method for hateful memes detection, comprising a primary multimodal task and two unimodal auxiliary tasks, while the authors in [165] delve into how textual and visual components contribute differently to hateful meme detection, shedding light on the nuanced interactions between modalities. Meanwhile, the authors in [12] propose a representation framework that facilitates inter-modal interaction and dynamically balances inter-modal and intra-modal relationships, providing a systematic way to disentangle memes into modality-invariant and modality-specific spaces.



**Figure 11:** Illustration of how altering text modifies the interpretation of an image. The left meme conveys a message about water crisis awareness. The right meme alters the text, resulting in a hateful interpretation of the child as "farm equipment", perpetuating dehumanization and racism. We blurred the face of the child for privacy concerns. Caution: The depicted memes contain toxic content. Viewer discretion is advised. These images do not reflect the views of the authors.

Attention frameworks are increasingly prevalent in cross-modal research for meme analysis. For instance, [212] employs an inter-modal attention framework to detect offensive memes by synergistically fusing visual and textual information, while the authors in [152] use a cross-attention network to explore connections between visual modalities

(using object detection and image caption models) and text features (using OCR) for hateful memes detection. Similarly, [231] proposes a Multimodal Visual-Textual Object Graph Attention Network (MViTO-GAT) for detecting propaganda techniques in memes. Their model learns semantic and positional relationships between visual and textual objects through attention-based sequential intra-modality and cross-modality graph reasoning, outperforming state-of-the-art baselines in performance. Many of the works we reviewed adopt CLIP embeddings, which are representations learned from diverse image-text pairs, providing meaningful embeddings for both images and text [241]. For example, [80] utilizes CLIP representations and modality-specific gating mechanisms to manage the interaction between textual and visual data. Similarly, the HateCLIPper architecture [9] explicitly models cross-modal interactions between image and text representations obtained using CLIP encoders and a feature interaction matrix to learn meaningful concepts. Also, [157] systematically analyzes semantic regularities in CLIP-generated embeddings, allowing for the study of how hateful memes evolve by fusing visual elements from multiple images or combining text with a hateful image.

Additionally, we notice that an increasing number of works explicitly focus on capturing more complex or nuanced connections between different modalities and discourse-intensive modeling of complex linguistic phenomena. For example, [127] propose the MemeFier approach, which employs a dual-stage modality fusion strategy to capture nuanced connections between text and image features. Leveraging Transformer encoders, MemeFier effectively learns inter-modality correlations at the token level, contributing to a deeper understanding of multimodal data. [162] introduces Topology-Aware Optimal Transport (TOT), a framework designed to handle complex cross-modal interactions by formulating optimal transportation plans and leveraging topology information for representation alignment. [15] develops the Analogy-aware Offensive Meme Detection (AOMD) framework, which focuses on accurately detecting offensive analogy memes. By capturing implicit analogy and aligning complex analogies across different modalities, AOMD achieves significant performance improvements over existing methods. Also, [228] present FigMemes, a dataset tailored for figurative language classification in politically-opinionated memes, providing benchmark results for both unimodal and multimodal models, while [242] utilize brain-like perceptual integration to reason about the subtle metaphors behind memes, exemplifying the trend toward more discourse-intensive modeling of cross-modal interplays.

## 9.2. Trend 2: Tackling Contextuality, Cultural and Background Knowledge

Memes are characterized by their intertextuality, often referencing elements from popular culture, symbols, events, or artifacts that hold significance within specific communities or affinity spaces [50]. This intertextuality renders memes highly contextual, necessitating background knowledge on various topics, including politics, current events, and cultural references, for accurate interpretation. This background knowledge, often referred to as "meme literacy" or "prior knowledge", is fundamental for comprehending the nuanced meanings embedded within memes [137]. In the realm of *toxic* memes, contextualization is even more critical: a seemingly innocuous image can take on a malicious connotation simply with a shift in context, which can completely change the interpretation and impact of the content.[8] Therefore, the direct incorporation of external knowledge into the classification process is recently emerging as a promising strategy to enhance harmful meme detection effectiveness and improve real-world applicability [11].

A prominent trend involves combining named entity recognition (NER) with background knowledge linking. The authors in [219] incorporate conceptual information from the external knowledge base Probase [243], a large-scale knowledge base that provides isA semantic relations, to enhance semantic representation and allows their model, MeBERT, to retrieve relevant concepts from meme text. This conceptual information is used as an extra modality to bridge the meme text to the image, using a concept-image attention module to align the concepts with corresponding image regions. This approach results in a concept-sensitive visual representation, where important image regions receive more attention, and significantly improves meme classification performance. However, challenges remain with text-only memes and long-tail entities with sparse information in the knowledge base. The authors in [150] enhance meme classification with their MemeGraphs method, which builds upon but goes beyond [219] by using automated NER and knowledge base (KB) augmentation not only on meme text but also on text extracted by first transforming memes into scene graphs. This process retrieves background knowledge for each identified entity from Wikidata, concatenates these augmentations to the meme text, and then feeds it into a Transformer for text classification. By integrating structured representations and interactions between internal entities and external knowledge, the MemeGraphs method consistently improves classification performance compared to models that rely solely on learned representations.

Other works utilize ConceptNet as the background knowledge base. [216] introduce KnowMeme, which constructs a graph representing meme content and related knowledge retrieved from ConceptNet, and then conducts graph

---

[8]For instance, the Nesquik bunny, originally featured in public information and awareness campaigns, was repurposed and recontextualized into a disparaging humor meme mocking the African water crisis [42].

classification to determine whether the meme is harmful. By capturing implicit meanings and cross-modal relations in memes, KnowMeme achieves significant performance enhancements compared to baseline methods. Similarly, [11] introduce KERMIT (Knowledge-Empowered Model), which builds a meme's knowledge-enriched network by merging internal meme entities with relevant external knowledge from ConceptNet. It starts by extracting the meme's text and caption using OCR and BLIP. Graph extraction involves identifying nodes (primarily nouns) from the text and caption through POS tagging, establishing relationships between them via dependency parsing, simplifying these relationships, and merging dependency trees to connect related nodes from the text and caption to form the final meme graph. For knowledge enrichment, KERMIT leverages ConceptNet iteratively. Subsequently, KERMIT employs a dynamic learning mechanism that leverages memory-augmented neural networks and attention mechanisms to discern the most informative segment of the knowledge-enriched information network to accurately classify harmful memes.

A key trend is the development of meme-specific external knowledge bases.The Image Meme Knowledge Graph (IMKG), recently published by [244], is a groundbreaking tool for studying Internet memes, providing a comprehensive and structured repository of meme-related semantics in text, vision, and metadata. It integrates data from diverse sources, enriches information through entity extraction and semantic links, and follows Semantic Web principles for effective knowledge representation. The authors used Wikidata [245] to extract background knowledge about meme seeds and entities, and KnowYourMeme (KYM) to collect and catalog meme-related data such as lore, interpretations, historical context, origins, and popularity. For textual enrichment, they used DBpedia Spotlight [246] to extract entities from KYM paragraphs (including sections about, origin, and spread) and from ImgFlip captions, linking these DBpedia entries to Wikidata entities. Additionally, they employed the Google Vision API for visual enrichment to detect objects and link them to Freebase [247]. While IMKG lacks image data and does not utilize the graph for downstream tasks, the authors explicitly suggest incorporating IMKG into methods to enhance the accuracy and explainability of neuro-symbolic methods for internet memes, such as hate speech detection and classification. The recent work by [248] introduces KYMKB, a knowledge base with 54,000 meme-related images and detailed information, notably focusing on meme templates. It distinguishes between templatic and non-templatic memes, a first in AI literature. KYMKB provides extensive data about templatic memes, including title, meaning, and origin, making it valuable for various tasks. The authors employ KYMKB for offensive meme detection, achieving a promising accuracy.

### 9.3. Trend 3: Interpretability, Explanations and Meme Literacy

We have observed a growing emphasis on interpretability in toxic meme detection, with many of the works we surveyed acknowledging it as a crucial aspect [11, 100, 132, 150, 201]. This trend seems to be part of a broader movement towards developing complementary measures alongside algorithms for automatic detection and removal of harmful content [249]. Specifically, there's a critical trend emerging where works focus on providing comprehensive rationales or explanations to aid users and content moderators in understanding the nuances of toxic memes [99]. This shift reflects a growing recognition of the importance of enhancing users' understanding of harmful content, particularly in the context of memes, to improve media literacy [137]. While straightforward indicators like nudges, labels, and red flags have been suggested for identifying hateful memes, their efficacy remains largely untested ([250]).

To enhance the interpretability of meme toxicity, recent trends include explicit identification of targeted entities, exploration of underlying themes and similarities, and utilization of pre-trained language models for enhanced analysis and explanation. The explicit identification of entities targeted by toxicity increasingly approached automatically, is exemplified in DisMultiHate ([151]), which disentangles target entities within multimodal memes to improve the classification and explainability of hateful content, as well as in DISARM ([94]), which employs named entity recognition and person identification for this purpose. Recent research efforts also aim to provide insights into underlying themes and similarities, such as [46], which proposes modular and explainable architectures for meme understanding using example- and prototype-based reasoning and the Hate-CLIPper approach by [9], which illuminates underlying themes and similarities across modalities. Furthermore, a recent trend in works addressing toxic memes involves enhancing interpretability by leveraging pre-trained language models, such as BERT and RoBERTa, to analyze linguistic patterns and contextual cues and generate insights into why certain content is classified as sexist [193] and to perform abductive reasoning to explore the interplay of multimodal information in memes ([81]).

Critically, recent datasets include ground truth explanation annotations for memes, as noted in Section 6.2. For instance, the authors in [137] propose quality-controlled crowdsourcing as an effective strategy for offering explanations and background knowledge for hateful memes through a Generate-Annotate-Revise workflow. The MultiBully-Ex dataset [80] provides multimodal explanations, combining visual cues like image segmentation with textual cues such as words relevant to or explaining the cyberbullying. Their experimental results demonstrate that training with multimodal

explanations improves performance in generating textual justifications and accurately identifying visual evidence. Additionally, in the HatReD dataset [99], hateful memes are annotated with textual explanations, and models are trained to decode the meaning of multimodal hateful memes and provide explicit explanations for the classifications.

## 9.4. Trend 4: Low-Resource Languages

As social media continues to expand globally, the spread of harmful content, including toxic memes, transcends linguistic and cultural boundaries. While considerable attention has been given to detecting such content in English, there is a growing realization of the urgency to address this issue in low-resource languages and contexts. Researchers have proposed integrating state-of-the-art deep learning models, such as BERT and Electra, for multilingual text analysis, alongside face recognition and optical character recognition models for comprehending meme images [90]. Additionally, Bengali BERT models have been deployed for automated Bengali abusive text classification, aiming to streamline the hate speech filtering process in resource-constrained languages, achieving notable accuracy and performance [116]. In low-resource settings, multi-modal prompt tuning has emerged as an effective approach for detecting propaganda techniques in memes. This method incorporates visual cues into language models through prompt-based multi-modal fine-tuning, showcasing efficacy in resource-limited scenarios [229]. Furthermore, transfer learning techniques have been explored to extend abusive meme detection to multiple languages. By leveraging model transfer techniques, researchers aim to bridge the language gap and establish baseline models for detecting abusive memes [115]. Dataset creation is also a significant effort, particularly in Asian languages such as Hindi, Bengali, Tamil, and Chinese. Efforts have been made to address the lack of benchmark datasets for specific languages by creating resources like the BanglaAbuseMeme dataset for Bengali abusive meme classification. These datasets facilitate the development and evaluation of models for detecting abusive content in low-resource languages ([13]).

## 9.5. Trend 5: Generative AI, Large Language Models (LLMs)

The emergence of Large Language Models (LLMs) and generative artificial intelligence (AI) has opened up new avenues for detecting and interpreting toxic memes. These advanced models offer capabilities for understanding implicit meanings and context in multimodal content, thereby addressing the challenges posed by harmful memes. In [81], LLMs are utilized for abductive reasoning to detect harmful memes. By distilling multimodal reasoning knowledge from LLMs and conducting lightweight fine-tuning, the model demonstrates effectiveness in identifying toxic content. Similarly, [79] explores the use of multimodal debate between LLMs for explainable harmful meme detection. By facilitating a debate between LLMs, the model gains insights into the context and implicit meanings of memes, leading to more interpretable results and a deeper understanding of meme content. Another notable application is in the correction of hate speech within multimodal memes, as proposed by [178]. Leveraging a large visual language model, the method effectively detects and corrects hate speech, contributing to the mitigation of online toxicity. Furthermore, [85] investigates the capability of GPT models to analyze the emotions conveyed through memes. By analyzing meme content using GPT models, the study provides insights into the emotional content and context of memes, showcasing the potential of LLMs for emotion analysis tasks. Additionally, [182] offers a comprehensive review of vision-language models and their performance on the Hateful Memes Challenge. Through an analysis of different models and techniques, valuable insights are provided for future research in the domain of detecting hateful memes. [139] propose PromptHate, a prompt-based model that leverages pre-trained language models for hateful meme classification. By constructing prompts and providing context examples, PromptHate exploits implicit knowledge in pre-trained models to achieve high classification accuracy, showcasing the potential of leveraging external knowledge for improved classification. Finally, [14] introduces a unified Multimodal Generative framework (MGex) for detecting cyberbullying in memes. This framework reframes the problem of meme detection as a multimodal text-to-text generation task, achieving competitive performance against baselines and state-of-the-art models. We also note that the growing number of guardrails in LLMs to prevent unsafe use will likely limit their ability to generate and analyze unsafe content in the future. Many examples in toxic datasets trigger these guardrails, designed to mitigate the processing of hateful or offensive content ([251]).

A related trend is the investigation of text-to-image models for the generation of unsafe images and hateful memes. In [158], the focus is on demystifying the generation of unsafe images and hateful memes from Text-to-Image models. The study assesses the proportion of unsafe images generated by advanced Text-to-Image models and evaluates the potential of generating hateful meme variants, highlighting risks associated with model misuse. Furthermore, [86] explores the proactive generation of unsafe images from Text-To-Image models using benign prompts. By studying the generation process, the research highlights potential risks associated with model misuse and emphasizes the importance of implementing safety measures to mitigate the generation of harmful content.

## 10. Discussion and Future Directions

This survey serves as a roadmap for researchers seeking to understand the landscape of computational perspectives on toxic internet memes. It focuses on multimodal toxic meme analysis and sheds light on the nuanced, complex taxonomical relationships within harmful online content (see Section 2). Our examination of existing surveys on computational toxic meme analysis emphasized the need for an up-to-date survey on toxic memes to address critical gaps in the literature while highlighting emerging trends and areas of research in Section 3. Employing the PRISMA 2020 guidelines (see Section 4), we surveyed 158 papers from 2019 to 2024, filling a crucial gap in the literature.

### 10.1. Contributions and Implications

This survey provides insights to support the advancement of computational models and tools for detecting, analyzing, and mitigating the proliferation of toxic memes in online environments. The theoretical implications of our study are described in Sections 7 and 8, where we harmonize the definitions of meme toxicities and provide a framework to identify their constituting elements. We identified and harmonized 12 meme toxicity terms, noting a significant focus on hateful memes and gaps in research on less prominent categories like troll, derogatory, and disinformation memes. We developed a taxonomy to clarify explicit and implicit taxonomical relationships among these terms, addressing discrepancies and enhancing adaptability for future studies. This offers a structured framework for comprehensively understanding and classifying various types of toxic content found in memes. This taxonomy serves as a valuable resource for researchers seeking to analyze, categorize, and compare different toxic memes. Furthermore, our taxonomy addresses discrepancies and nuances in existing frameworks, enabling researchers to guide their investigations effectively and explore the complex relationships between different toxicity categories with clarity and precision. This resulting framework has also significant practical implications since it identifies complex features that may be used by developers and practitioners to automatically or semi-automatically analyze memes to classify their toxicity. The datasets described and compared in Section 6 provide a map for possible training sets to be used to create such (semi-)automated models.

Another practical contribution is the synthesis of annotation guidelines across diverse toxic meme datasets detailed in Section 6.2. By characterizing the annotation methodologies and the range of computational task definitions utilized across various datasets for different meme attributes, we furnish researchers with a useful tool to inform their future investigations. Furthermore, this resource empowers researchers to pose targeted inquiries about meme attributes, including assessing the efficacy of different annotation schemes in capturing nuanced meme characteristics.

Meme toxicity definitions predominantly focus on three core elements: the targeted entity, the underlying intent behind meme creation or dissemination, and the rhetorical techniques used (see Section 8). Building upon these observations, we pinpoint three key dimensions of meme toxicity—target, intent, and rhetorical tactics—that are frequently examined in isolation within computational studies. Researchers can utilize these dimensions and the relationships among them to systematically categorize and analyze various types of toxic content found in memes.

In Section 9, we identified trends and research directions in automatic toxic meme detection and understanding. These trends encompass tackling complex cross-modal entailments by leveraging multi-modal deep learning models, attention frameworks, and multidisciplinary approaches. Also, there is a focus on addressing contextuality and cultural background knowledge through named entity recognition and connection with commonsense knowledge bases, and the development of meme-specific external knowledge bases. We observed an emphasis on interpretability and explainability. Furthermore, attention is directed towards low-resource languages and the potential of generative AI and Large Language Models (LLMs) for detecting and interpreting toxic memes.

### 10.2. Future Research Directions

The recent emphasis on providing explanations alongside toxic meme detection marks a shift towards more transparent and accountable AI systems. However, there is a need for further clarification regarding the types of interpretability required and the appropriate methods for achieving them. It is crucial to differentiate between methods ensuring transparency of computation (e.g., explaining what the machine computed) and those providing users with understandable explanations for the toxicity label assigned to a meme (e.g., explaining why a meme is considered propagandistic). We conclude that the field is shifting to fulfill the latter—offering reasonable explanations for assigned labels. Interpretable methods help build trust in automatic moderation systems, ensuring that toxicity labels are properly explained. Based on the contributions identified, we propose the following future research directions:

*Semiotics-Informed Cross-Modal Reasoning* Leveraging insights from semiotic research can significantly enhance the development of an advanced automatic understanding of cross-modality. Interdisciplinary collaboration with

semiotics experts can operationalize these insights, improving AI-driven meme analysis systems. Future research should focus on advancing these models to effectively capture the multiplicative nature of meaning-making in multimodal content. This includes exploring semiotically relevant features with feature engineering and utilizing advanced fusion techniques like modality-specific gating and attention mechanisms to balance inter-modal and intra-modal relationships. By developing cross-modal reasoning models, researchers can enable machines to perform critical semiotic-based analyses, providing valuable insights into toxic message construction and dissemination. Using these insights in explanations, automatic systems can help users understand how toxic messages are constructed and conveyed in multimodal content. The integration of semiotic insights can empower content moderators and users with clues about toxic message construction, enhancing AI interpretability and informed decision-making. Prioritizing semiotic-informed cross-modal reasoning models can create more effective and transparent solutions for addressing toxic content.

*Incorporating Toxicity-Specific Background Knowledge*  Building upon the current trend of integrating external knowledge bases into meme analysis, there are several promising research directions to explore. While resources like ConceptNet, WikiData, and Probase have enhanced meme interpretation, a critical gap remains: incorporating databases with cultural knowledge specifically related to toxicity. For instance, the Global Extremist Symbols Database[9] provided by the Global Project Against Hate and Extremism (GPAHE) offers a wealth of information on cultural icons associated with toxicity and hate, such as hate symbols, visual and numerical icons, flags, hand gestures, acronyms, and salutes used by extremist groups. Research efforts could focus on integrating such specialized knowledge bases into meme analysis frameworks and exploring hybrid approaches that combine external knowledge with pre-trained vision-language models, similar to KERMIT [11], to enhance classification performance by leveraging the strengths of both approaches. Additionally, given the dynamic nature of online content, developing adaptive learning mechanisms that continuously update knowledge representations to reflect evolving trends and cultural contexts would be beneficial. Overall, by broadening the scope of incorporated knowledge, future research can enhance the contextual understanding and analysis of memes within their socio-cultural contexts, leading to more effective meme classification.

*Embracing Linguistic Diversity in Toxic Meme Detection*  Expanding toxic meme detection beyond English to include other widely spoken languages is imperative, along with annotating datasets from diverse linguistic and cultural backgrounds. While leveraging models trained in English may provide a foundation, incorporating annotations from other languages is crucial to ensure the effectiveness and cultural relevance of detection methods. Additionally, exploring novel approaches that consider the unique linguistic and cultural characteristics of different languages can enhance the accuracy and applicability of toxic meme detection systems. Further work should investigate the transferability of existing models across languages and develop techniques to adapt them to new linguistic contexts, as embracing linguistic diversity and cultural nuances in toxic meme detection research can lead to more robust and inclusive frameworks.

*Leveraging LLMs*  Recent research highlights the potential of Large Language Models (LLMs) and generative AI in identifying and interpreting toxic memes. However, further investigation is needed to understand their strengths and weaknesses in meme detection, including their ability to differentiate between harmless and toxic memes and any inherent cultural or modal biases. Additionally, exploring LLMs' capacity for cross-modal analysis of meme content with multimodal data containing nuanced cross-modal interplays is crucial. Future research should also focus on assessing LLMs' cultural knowledge and understanding, exploring efficient methods like Retrieval-Augmented Generation (RAG) for incorporating background knowledge. Understanding how LLMs can adapt to evolving cultural references and incorporate diverse perspectives in meme interpretation is crucial. Equally important is further investigating the effectiveness and ethical implications of safeguard rails in LLMs for toxic meme detection.

*Multi-Dimensional Exploration of Meme Toxicity*  Recognizing that, for AI to detect multimodal toxicity, "it must learn to understand content the way that people do: holistically" [96, p. 201], we advocate for a multidimensional approach. Future research should focus on simultaneously analyzing multiple dimensions using an evidential reasoning approach grounded in decision theory, allowing for the accumulation and scrutiny of evidence to guide judgments or decisions. Additionally, there is a need to explore potential common relationships between meme toxicity types and the tactics used to convey them. Furthermore, further investigation into additional dimensions, such as the context of posting (user, forum, platform) and propagation features [252], is warranted, as these factors can provide a more comprehensive understanding of meme toxicity dynamics.

---

[9]https://globalextremism.org/global-extremist-symbols-database/

*Legal, Ethical and Collateral Impacts* Exploring legal and ethical considerations in automatic moderation is crucial. Automatic moderation systems must navigate legal and ethical landscapes to discern illegal content from harmful but legal content. Determining computational tasks based on content legality and harm potential raises complex ethical dilemmas that require careful consideration. Further research is also needed to address other issues, such as collateral damage in moderation algorithms. Evaluation of these algorithms should account for false positives and the potential unintended consequences they may have, particularly concerning the censorship of speech from marginalized groups or the unintentional suppression of resistance narratives. Other ethical challenges in meme analysis include addressing subjectivity, combating the use of AI to generate toxic content, and navigating the complexities of toxic positivity.

## 11. Conclusion

Toxic memes represent a spectrum of harmful or malevolent multimodal content disseminated across online platforms, often with the intention of promoting harm and hate, spreading disinformation, or promoting exploitation. Understanding the conceptualizations and distinguishing features of toxic memes is crucial for developing effective computational models capable of detecting and moderating such content. Our survey provides readers with a comprehensive understanding of toxic memes from a computational, content-based perspective, covering key developments up to early 2024. Our survey identified that the computational field has used a wide variety of terminology to refer to toxic memes, highlighting the increasing demand for automatic tools for identifying fine-grained toxicity types beyond a simple toxic/non-toxic classification. This includes specifying whether a meme is misogynistic, spreading disinformation, involved in cyberbullying, and so on. This complexity necessitates harmonizing term definitions and the establishment of a clearer taxonomy that delineates how these terms relate to each other. In response to this need, we provide a harmonized set of definitions and introduce a novel taxonomy in Section 7. We offer insights into various dimensions of meme toxicity, including intent, target, and conveyance tactics, along with a standardized taxonomy for categorizing meme toxicity types. Also, we catalog datasets containing toxic memes, analyze prevalent challenges, and identify emerging trends in computational approaches to toxic meme detection and interpretation. Through synthesizing existing knowledge and identifying research gaps, we aim to promote interdisciplinary collaboration and innovation to foster media literacy and potentially a safer, higher quality, and more inclusive online ecosystem. Our survey offers pathways for computational advancement in the field, including enhancing interpretability through sophisticated cross-modal reasoning, background knowledge integration, attention on low-resource languages, and refining the usage of LLMs.

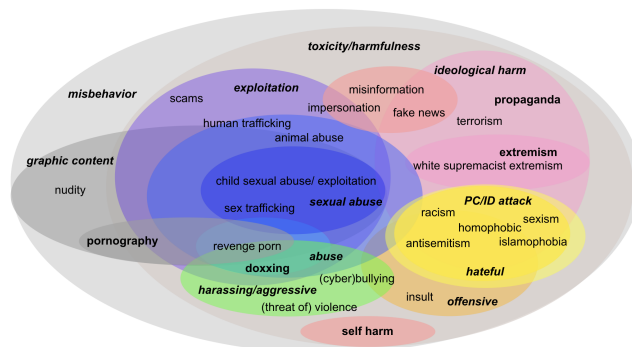## 12. Acknowledgements

## A. Appendix

**Figure 12:** Toxicity-related terms derived from our investigation of harmfulness and toxicity in multimodal data, illustrating the complex and overlapping nature of multimodal toxicities.

Toxic Memes

| Dataset | Search Engines | | | | Img Platforms | | Social Media Platforms | | | | | | | Meme Platforms | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B | G | D | B | Pin | Img | FB | Gab | IG | RD | TW | Wei | WA | KYM | MD | MG | 9gag |
| AOMD Gab | | | | | | | | ✓ | | | | | | | | | |
| AOMD Reddit | | | | | | | | | | ✓ | | | | | | | |
| BanglaAbuseMeme | ✓ | ✓ | | | | | ✓ | | ✓ | | | | | | | | |
| CrisisHateMM | | | | | | | ✓ | | | ✓ | ✓ | | | | | | |
| Derogatory Fb-Meme | | | | | | | ✓ | | | | | | | | | | |
| DisinfoMeme | | | | | | | ✓ | | | | | | | | | | |
| ELEMENT | | ✓ | | | | | | | | ✓ | | | | | | ✓ | |
| Emoffmeme | | ✓ | | | | | | | | | | | | | | | |
| Ext-Harm-P | | ✓ | | | | | ✓ | | ✓ | ✓ | | | | | | | |
| Facebook HM | | | | | | | ✓ | | | | | | | | | | |
| FAME dataset | | | ✓ | | | | | | | | | | | | | | |
| Fine grained HM | | | | | | | ✓ | | | ✓ | | | | | | ✓ | |
| GOAT-Bench | | ✓ | | | | ✓ | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | ✓ |
| Harm-C | | ✓ | | | | | ✓ | | ✓ | ✓ | | | | | | | |
| Harm-P | | ✓ | | | | | ✓ | | ✓ | ✓ | | | | | | | |
| Hate Speech in Pixels | | ✓ | | | | | | | | ✓ | | | | | | | |
| HatReD | | | | | | | ✓ | | | | | | | | | | |
| HVVMemes | | ✓ | | | | | ✓ | | ✓ | ✓ | | | | | | | |
| Indian Political Memes | | ✓ | | | | | | | | | | | | | | | |
| Innopolis Hateful Memes | | | ✓ | | | | | | | | ✓ | | | | ✓ | | |
| KAU-Memes | | | | | | | ✓ | | ✓ | ✓ | ✓ | | | | | | |
| Meme-Merge | | | | ✓ | | | | | | | ✓ | ✓ | | | | | |
| Memotion 1 | | ✓ | | | | | | | | | | | | | | | |
| Memotion 2 | | | | | | ✓ | ✓ | | ✓ | ✓ | | | | | | | |
| MET-Meme | | ✓ | | ✓ | | | | | | | ✓ | ✓ | | | | | |
| Misogynistic-MEME | | | | | | | ✓ | | ✓ | ✓ | ✓ | | | | | | |
| MultiBully | | | | | | | | | | ✓ | ✓ | | | | | | |
| MultiBully-Ex | | | | | | | | | | ✓ | ✓ | | | | | | |
| MAMI | | | | | | ✓ | | | | ✓ | ✓ | | | ✓ | | | ✓ |
| MultiOFF | | | | | | | ✓ | | ✓ | ✓ | ✓ | | | | | | |
| Pol_Off_Meme | | ✓ | | | | | | | | | | | | | | | |
| SemEval-2021 Task 6 | | | | | | | ✓ | | | | | | | | | | |
| TamilMemes | | | | | ✓ | | ✓ | | ✓ | | | | ✓ | | | | |
| TrollsWithOpinion | | ✓ | | | | | | | | | | | | | | | |

**Table 9**

Sources of meme acquisition, as documented in the papers introducing the datasets. Rows are datasets, while columns are sources categorized into search engines, image hosting platforms, social media platforms, or meme creation and sharing platforms. Abbreviations used: B: Bing, G: Google, D: DuckDuckGo, B: Baidu, Pin: Pinterest, Img Platforms: Image Hosting Platforms, Img: Imgur, FB: Facebook, IG: Instagram, Meme Platforms: Meme Creation and Sharing Platforms, RD: Reddit, TW: Twitter, Wei: Weibo, WA: WhatsApp, KYM: KnowYourMeme, MD: Memedroid, MG: MemeGenerator.

# References

[1] R. Dawkins, The Selfish Gene, new ed., Oxford University Press, 1989.

[2] D. Dennet, Dangerous memes - a ted talk, Youtube, 2007. URL: https://www.youtube.com/watch?v=KzGjEkp772s.

[3] C. Koutlis, M. Schinas, S. Papadopoulos, Memetector: Enforcing deep focus for meme detection, International Journal of Multimedia Information Retrieval 12 (2023) 11.

[4] S. Peeters, M. Tuters, T. Willaert, D. De Zeeuw, On the vernacular language games of an antagonistic online subculture, Frontiers in big Data 4 (2021) 718368.

[5] C. Arkenbout, Political meme toolkit: leftist dutch meme makers share their trade secrets, in: Critical meme reader II: memetic tacticality, Institute of Network Cultures, 2022, pp. 20–31.

[6] A. Wagener, Semiotic excess in memes: From postdigital creativity to social violence, Internet Pragmatics 6 (2023) 239–258.

[7] D. Kiela, H. Firooz, A. Mohan, V. Goswami, A. Singh, C. Fitzpatrick, P. Bull, G. Lipstein, T. Nelli, R. Zhu, R. Muennighoff, R. Velioglu, J. Rose, P. Lippe, N. Holla, S. Chandra, S. Rajamanickam, G. Antoniou, E. Shutova, H. Yannakoudakis, V. Sandulescu, U. Ozertem, P. Pantel, L. Specia, D. Parikh, The hateful memes challenge: Competition report, in: Proceedings of Machine Learning Res., volume 133, 2020.

[8] H. Kirk, Y. Jun, P. Rauba, G. Wachtel, R. Li, X. Bai, N. Broestl, M. Doff-Sotta, A. Shtedritski, Y. Asano, Memes in the wild: Assessing the generalizability of the hateful memes challenge dataset, in: Proceedings of WOAH 2021 - 5th Workshop on Online Abuse and Harms, volume 26, 2021.

[9] G. Kumar, K. Nandakumar, Hate-clipper: Multimodal hateful meme classification based on cross-modal interaction of clip features, in: NLP4PI 2022 - 2nd Workshop on NLP for Positive Impact, Proceedings of the Workshop, volume 171, 2022.

[10] S. Sharma, F. Alam, M. S. Akhtar, D. Dimitrov, G. Da San Martino, H. Firooz, A. Halevy, F. Silvestri, P. Nakov, T. Chakraborty, Detecting and understanding harmful memes: A survey, arXiv preprint arXiv:2205.04274 1 (2022) 1–9.

[11] B. Grasso, V. La Gatta, V. Moscato, G. Sperlì, Kermit: Knowledge-empowered model in harmful meme detection, Information Fusion 106 (2024).

[12] C. Yang, F. Zhu, J. Han, S. Hu, Invariant meets specific: A scalable harmful memes detection framework, MM 2023 - Proceedings of the 31st ACM International Conference on Multimedia 4788 (2023).

[13] M. Das, A. Mukherjee, Banglaabusememe: A dataset for bengali abusive meme classification, arXiv preprint arXiv:2310.11748 (2023) 15498–15512.

[14] R. Jain, K. Maity, P. Jha, S. Saha, Generative models vs discriminative models: Which performs better in detecting cyberbullying in memes?, in: Proceedings of the International Joint Conference on Neural Networks, 2023. doi:10.1109/ijcnn54540.2023.10191363.

[15] L. Shang, Y. Zhang, Y. Zha, Y. Chen, C. Youn, D. Wang, Aomd: An analogy-aware approach to offensive meme detection on social media, Information Processing & Management (2021).

[16] C. Sharma, D. Bhageria, W. Scott, S. Pykl, A. Das, T. Chakraborty, V. Pulabaigari, B. Gamback, Semeval-2020 task 8: Memotion analysis–the visuo-lingual metaphor!, arXiv preprint arXiv:2008.03781 (2020) 759–773.

[17] A. Arora, P. Nakov, M. Hardalov, S. M. Sarwar, V. Nayak, Y. Dinkov, D. Zlatkova, K. Dent, A. Bhatawdekar, G. Bouchard, I. Augenstein, Detecting harmful content on online platforms: What platforms need vs. where research efforts go, 2023. arXiv:2103.00153.

[18] A. Williams, M. Dupuis, I don't always spread disinformation on the web, but when i do i like to use memes: An examination of memes in the spread of disinformation, in: Proceedings of the 11th International Multi-Conferences on Complexity, Informatics and Cybernetics: IMCIC, 2020, pp. 165–172.

[19] D. Dimitrov, B. B. Ali, S. Shaar, F. Alam, F. Silvestri, H. Firooz, G. D. S. Martino, Detecting propaganda techniques in memes, arXiv preprint arXiv:2109.08013 (2021) 6603–6617.

[20] M. Abdullah, D. Abujaber, A. Al-Qarqaz, R. Abbott, M. Hadzikadic, Combating propaganda texts using transfer learning, IAES International Journal of Artificial Intelligence 12 (2023) 956–965.

[21] D. Rodríguez, P. Nakov, V. Dankers, E. Shutova, Paper bullets: Modeling propaganda with the help of metaphor, European Chapter of the Association for Computational Linguistics, Findings of EACL 2023 (2023) 472–489.

[22] M. G. Shridara, D. Hládek, M. Pleva, R. Haluška, Identification of trolling in memes using convolutional neural networks, in: 2023 33rd International Conference Radioelektronika (RADIOELEKTRONIKA), IEEE, 2023, pp. 1–6.

[23] S. Suryawanshi, B. R. Chakravarthi, M. Arcan, P. Buitelaar, Trollswithopinion: A taxonomy and dataset for predicting domain-specific opinion manipulation in troll memes, Multimedia Tools and Applications 82 (2023) 9137–9171.

[24] D. Bebić, M. Volarevic, Do not mess with a meme: the use of viral content in communicating politics, Communication & Society 31 (2018) 43–56.

[25] G. Mazambani, M. A. Carlson, S. Reysen, C. F. Hempelmann, Impact of status and meme content on the spread of memes in virtual communities, Human Technology: An Interdisciplinary Journal on Humans in ICT Environments 11 (2015) 148–164.

[26] Netanya, The Dangers of Memes — netanyataitague, https://medium.com/@netanyataitague/the-dangers-of-memes-b1bb67e10083, 2019. [Accessed 17-04-2024].

[27] M. Vang, Is Meme Culture Problematic - The Current — thecurrentmsu.com, https://thecurrentmsu.com/2021/08/07/is-meme-culture-problematic/, 2021. [Accessed 17-04-2024].

[28] K. Rojas, The Toxicity of Online Meme Culture: When Is It Too Far? — studybreaks.com, https://studybreaks.com/thoughts/meme-culture-2/, 2022. [Accessed 17-04-2024].

[29] Z. D. Roberts, How the 'Free Helicopter Rides' Meme Went Viral — progressive.org, https://progressive.org/magazine/how-the-free-helicopter-rides-meme-went-viral-roberts-20230907/, 2023. [Accessed 17-04-2024].

[30] F. J. A. Serna, Los memes como simbolos del discurso de odio: La influencia del humor gráfico en la libertad de expresión y la política, VISUAL REVIEW. International Visual Culture Review/Revista Internacional de Cultura Visual 16 (2024) 241–253.

[31] L. Needham, How the toxic went mainstream — pursuit.unimelb.edu.au, https://pursuit.unimelb.edu.au/articles/how-the-toxic-went-mainstream, 2019. [Accessed 17-04-2024].

[32] K. M. Duchscherer, J. F. Dovidio, When memes are mean: Appraisals of and objections to stereotypic memes, Translational Issues in Psychological Science 2 (2016) 335–345.

[33] P. M. Bennet, Who Moderates the Social Media Giants? A Call to End Outsourcing, Technical Report, NYU STERN, Center for Business and Human Rights, 2020.

[34] N. Nondo, Facing disturbing content daily, online moderators in africa want better protections and a fair wage, CBC Radio (2023).

[35] B. Perrigo, Inside facebook's african sweatshop, Time (2022).

[36] N. Rowe, "it's destroyed me completely": Kenyan moderators decry toll of training of ai models, The Guardian (2023).

[37] N. Mbagathi, In africa, taking on viral hate, Open Society Foundations (2023).

[38] T. H. Afridi, A. Alam, M. N. Khan, J. Khan, Y.-K. Lee, A multimodal memes classification: A survey and open research issues, in: Innovations in Smart Cities Applications Volume 4: The Proceedings of the 5th International Conference on Smart City Applications, Springer, 2021, pp. 1451–1466.

[39] P. C. d. Q. Hermida, E. M. D. Santos, Detecting hate speech in memes: a review, Artificial Intelligence Review 56 (2023) 1–19.

[40] S. Zannettou, T. Caulfield, J. Blackburn, E. De Cristofaro, M. Sirivianos, G. Stringhini, G. Suarez-Tangil, On the origins of memes by means of fringe web communities, in: Proceedings of the internet measurement conference 2018, 2018, pp. 188–202.

[41] T. Chakraborty, S. Masud, Nipping in the bud: detection, diffusion and mitigation of hate speech on social media, ACM SIGWEB Newsletter 2022 (2022) 1–9.

[42] C. Polli, M. G. Sindoni, Multimodal computation or interpretation? automatic vs. critical understanding of text-image relations in racist memes, Ssrn (2023).

[43] O. Solon, Richard dawkins on the internet's hijacking of the word 'meme', Wired UK 20 (2013).

[44] L. Shifman, Memes in Digital Culture, MIT Press Essential Knowledge Series, MIT Press, 2013.

[45] M. Dynel, The life of covid-19 mask memes: a diachronic study of the pandemic memescape, Comunicar 30 (2022) 73–85.

[46] A. Thakur, F. Ilievski, H.-A. Sandlin, Z. Sourati, L. Luceri, R. Tommasini, A. Mermoud, Explainable classification of internet memes, in: CEUR Workshop Proceedings, volume 3432, 2023, pp. 395–409.

[47] L. Xie, A. Natsev, J. R. Kender, M. Hill, J. R. Smith, Visual memes in social media: tracking real-world news in youtube videos, in: Proceedings of the 19th ACM international conference on Multimedia, 2011, pp. 53–62.

[48] Y. Du, M. A. Masood, K. Joseph, Understanding visual memes: An empirical analysis of text superimposed on memes shared on twitter, in: Proceedings of the International AAAI Conference on Web and Social Media, volume 14, Aaai, 2020, pp. 153–164.

[49] S. Pramanick, D. Dimitrov, R. Mukherjee, S. Sharma, M. S. Akhtar, P. Nakov, T. Chakraborty, Detecting harmful memes and their targets, arXiv preprint arXiv:2110.00413 (2021).

[50] M. Knobel, C. Lankshear, Online memes, affinities, and cultural production, A new literacies sampler 29 (2007) 199–227.

[51] C. Sharma, V. Pulabaigari, A. Das, Meme vs. non-meme classification using visuo-linguistic association., in: WEBIST, 2020, pp. 353–360.

[52] V. Sherratt, K. Pimbblet, N. Dethlefs, Multi-channel convolutional neural network for precise meme classification, in: Proceedings of the 2023 ACM International Conference on Multimedia Retrieval, 2023, pp. 190–198.

[53] D. Kiela, H. Firooz, A. Mohan, V. Goswami, A. Singh, P. Ringshia, D. Testuggine, The hateful memes challenge: Detecting hate speech in multimodal memes, in: Advances in neural information processing systems, volume 33, 2020, pp. 2611–2624.

[54] B. Kostadinovska-Stojchevska, E. Shalevska, Internet memes and their socio-linguistic features, European journal of literature, language and linguistics studies 2 (2018).

[55] N. Lapidot-Lefler, A. Barak, Effects of anonymity, invisibility, and lack of eye-contact on toxic online disinhibition, Computers in human behavior 28 (2012) 434–443.

[56] D. Gordeev, V. Potapov, Toxicity in texts and images on the internet, in: International Conference on Speech and Computer, Springer, 2020, pp. 156–165.

[57] A. Sheth, V. L. Shalin, U. Kursuncu, Defining and detecting toxicity on social media: context and knowledge are key, Neurocomputing 490 (2022) 312–318.

[58] N. Carlisle, Toxicity, memes and raids - Power of Zero — powerof0.org, https://powerof0.org/toxicity-memes-and-raids/, 2022. [Accessed 17-04-2024].

[59] M. Mosleh, R. Cole, D. G. Rand, Misinformation and harmful language are interconnected, rather than distinct, challenges, PNAS nexus (2024) pgae111.

[60] S. Ghosh, S. Lepcha, S. Sakshi, R. R. Shah, S. Umesh, Detoxy: A large-scale multimodal dataset for toxicity classification in spoken utterances, arXiv preprint arXiv:2110.07592 (2021).

[61] Z. Waseem, T. Davidson, D. Warmsley, I. Weber, Understanding abuse: A typology of abusive language detection subtasks, arXiv preprint arXiv:1705.09899 (2017).

[62] M. Wiegand, M. Siegel, J. Ruppenhofer, Overview of the germeval 2018 shared task on the identification of offensive language, in: Proceedings of GermEval 2018, 14th Conference on Natural Language Processing (KONVENS 2018), Austrian Academy of Sciences, 2019, pp. 1 – 10. URL: https://nbn-resolving.de/urn:nbn:de:bsz:mh39-84935.

[63] M. Zampieri, S. Malmasi, P. Nakov, S. Rosenthal, N. Farra, R. Kumar, Semeval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval), arXiv preprint arXiv:1903.08983 (2019).

[64] P. Piot, P. Martín-Rodilla, J. Parapar, Metahate: A dataset for unifying efforts on hate speech detection, arXiv preprint arXiv:2401.06526 (2024).

[65] C. Adams, J. Sorensen, J. Elliott, L. Dixon, M. McDonald, n., W. Cukierski, Toxic comment classification challenge, 2017. URL: https://kaggle.com/competitions/jigsaw-toxic-comment-classification-challenge.

[66] T. Davidson, D. Warmsley, M. Macy, I. Weber, Automated hate speech detection and the problem of offensive language, in: Proceedings of the international AAAI conference on web and social media, volume 11, 2017, pp. 512–515.

[67] F. E. Ayo, O. Folorunso, F. T. Ibharalu, I. A. Osinuga, Machine learning techniques for hate speech classification of twitter data: State-of-the-art, future challenges and research directions, Computer Science Review 38 (2020) 100311.

[68] F. Alkomah, X. Ma, A literature review of textual hate speech detection methods and datasets, Information 13 (2022) 273.

[69] A. Chhabra, D. K. Vishwakarma, A literature survey on multimodal and multilingual automatic hate speech identification, Multimedia Systems 29 (2023) 1203–1230.

[70] P. Fortuna, J. Soler, L. Wanner, Toxic, hateful, offensive or abusive? what are we really classifying? an empirical analysis of hate speech datasets, in: Proceedings of the Twelfth Language Resources and Evaluation Conference, 2020, pp. 6786–6794.

[71] M. Yankoski, W. Scheirer, T. Weninger, Meme warfare: Ai countermeasures to disinformation should focus on popular, not perfect, fakes, Bulletin of the atomic scientists 77 (2021) 119–123.

[72] M. Banko, B. MacKeen, L. Ray, A unified taxonomy of harmful content, in: Proceedings of the fourth workshop on online abuse and harms, 2020, pp. 125–137.

[73] P. Nakov, V. Nayak, K. Dent, A. Bhatawdekar, S. M. Sarwar, M. Hardalov, Y. Dinkov, D. Zlatkova, G. Bouchard, I. Augenstein, Detecting abusive language on online platforms: A critical analysis, arXiv preprint arXiv:2103.00153 (2021).

[74] A. Halevy, C. Canton-Ferrer, H. Ma, U. Ozertem, P. Pantel, M. Saeidi, F. Silvestri, V. Stoyanov, Preserving integrity in online social networks, Communications of the ACM 65 (2022) 92–98.

[75] F. Alam, S. Cresci, T. Chakraborty, F. Silvestri, D. Dimitrov, G. D. S. Martino, S. Shaar, H. Firooz, P. Nakov, A survey on multimodal disinformation detection, arXiv preprint arXiv:2103.12541 (2021).

[76] Anjum, R. Katarya, Hate speech, toxicity detection in online social media: a recent survey of state of the art and opportunities, International Journal of Information Security 23 (2024) 577–608.

[77] B. Lewandowska-Tomaszczyk, A. Bączkowska, C. Liebeskind, G. Valunaite O., S. Žitnik, An integrated explicit and implicit offensive language taxonomy, Lodz Papers in Pragmatics 19 (2023) 7–48.

[78] T. Garg, S. Masud, T. Suresh, T. Chakraborty, Handling bias in toxic speech detection: A survey, ACM Computing Surveys 55 (2023) 1–32.

[79] H. Lin, Z. Luo, W. Gao, J. Ma, B. Wang, R. Yang, Towards explainable harmful meme detection through multimodal debate between large language models, arXiv preprint arXiv:2401.13298 (2024).

[80] P. Jha, K. Maity, R. Jain, A. Verma, S. Saha, P. Bhattacharyya, Meme-ingful analysis: Enhanced understanding of cyberbullying in memes through multimodal explanations, arXiv preprint arXiv:2401.09899 (2024).

[81] H. Lin, Z. Luo, J. Ma, L. Chen, Beneath the surface: Unveiling harmful memes with multimodal reasoning distilled from large language models, arXiv preprint arXiv:2312.05434 9114 (2023).

[82] J. Chen, D. Zhu, X. Shen, X. Li, Z. Liu, P. Zhang, R. Krishnamoorthi, V. Chandra, Y. Xiong, M. Elhoseiny, Minigpt-v2: large language model as a unified interface for vision-language multi-task learning, arXiv preprint arXiv:2310.09478 (2023).

[83] W. Dai, J. Li, D. Li, A. Tiong, J. Zhao, W. Wang, B. Li, P. N. Fung, S. Hoi, Instructblip: Towards general-purpose vision-language models with instruction tuning, in: Poster at Advances in Neural Information Processing Systems, 2023.

[84] Y. Hu, O. Stretcu, C.-T. Lu, K. Viswanathan, K. Hata, E. Luo, R. Krishna, A. Fuxman, Visual program distillation: Distilling tools and programmatic reasoning into vision-language models, arXiv preprint arXiv:2312.03052 (2023).

[85] J. Wang, J. Luo, G. Yang, A. Hong, F. Luo, Is gpt powerful enough to analyze the emotions of memes?, arXiv preprint arXiv:2311.00223 (2023).

[86] Y. Wu, N. Yu, M. Backes, Y. Shen, Y. Zhang, On the proactive generation of unsafe images from text-to-image models using benign prompts, arXiv:2310.16613 (2023).

[87] H. Lin, Z. Luo, B. Wang, R. Yang, J. Ma, Goat-bench: Safety insights to large multimodal models through meme-based social abuse, arXiv preprint arXiv:2401.01523 (2024).

[88] E. Hossain, O. Sharif, M. Hoque, M. Akber Dewan, N. Siddique, M. Hossain, Identification of multilingual offense and troll from social media memes using weighted ensemble of multimodal features, in: Journal of King Saud University - Computer and Information Sciences, volume 34, 2022, pp. 6605–6623. doi:10.1016/j.jksuci.2022.06.010.

[89] A. Bhandari, S. B. Shah, S. Thapa, U. Naseem, M. Nasim, Crisishatemm: Multimodal analysis of directed and undirected hate speech in text-embedded images from russia-ukraine conflict, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 1993–2002.

[90] R. Bhowmick, I. Ganguli, J. Paul, J. Sil, A multimodal deep framework for derogatory social media post identification of a recognized person, ACM Transactions on Asian and Low-Resource Language Information Processing 21 (2022) 3447651.

[91] J. Qu, L. H. Li, J. Zhao, S. Dev, K.-W. Chang, Disinfomeme: A multimodal dataset for detecting meme intentionally spreading out disinformation, arXiv preprint arXiv:2205.12617 (2022).

[92] N. Zhang, X. Feng, T. Gu, L. Chang, Mvlp: Multi-perspective vision-language pre-training model for ethically aligned meme detection, Authorea Preprints (2023).

[93] G. Kumari, D. Bandyopadhyay, A. Ekbal, Emoffmeme: identifying offensive memes by leveraging underlying emotions, in: Multimedia Tools and Applications, volume 82, 2023, pp. 45061–45096. doi:10.1007/s11042-023-14807-1.

[94] S. Sharma, M. S. Akhtar, P. Nakov, T. Chakraborty, Disarm: Detecting the victims targeted by harmful memes, arXiv preprint arXiv:2205.05738 (2022) 1572–1588.

[95] B. Jabiyev, J. Onaolapo, G. Stringhini, E. Kirda, e-game of fame: Automatic detection of fake memes., in: TTO, 2021, pp. 1–11.

[96] L. Mathias, S. Nie, A. M. Davani, D. Kiela, V. Prabhakaran, B. Vidgen, Z. Waseem, Findings of the woah 5 shared task on fine grained hateful memes detection, in: Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021), 2021, pp. 201–206.

[97] S. Pramanick, S. Sharma, D. Dimitrov, P. Nakov, T. Chakraborty, Momenta: A multimodal framework for detecting harmful memes and their targets, arXiv (2021).

[98] B. O. Sabat, C. C. Ferrer, X. Giro-I-Nieto, Hate speech in pixels: Detection of offensive memes towards automatic moderation, arXiv (2019).

[99] M. Hee, W.-H. Chong, R.-W. Lee, Decoding the underlying meaning of multimodal hateful memes, in: IJCAI International Joint Conference on Artificial Intelligence, volume 2023-August, 2023.

[100] S. Sharma, T. Suresh, A. Kulkarni, H. Mathur, P. Nakov, M. S. Akhtar, T. Chakraborty, Findings of the constraint 2022 shared task on detecting the hero, the villain, and the victim in memes, in: Proceedings of the Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situations, Constraint, 2022, pp. 1–11.

[101] K. Rajput, R. Kapoor, K. K. Rai, P. Kaur, Hate me not: Detecting hate inducing memes in code switched languages, arXiv (2022).

[102] J. Badour, J. Brown, Hateful memes classification using machine learning, in: 2021 IEEE Symposium Series on Computational Intelligence, SSCI 2021 - Proceedings, volume 2, 2021. doi:10.1109/ssci50451.2021.9659896.

[103] J. Bacha, F. Ullah, J. Khan, A. Sardar, S. Lee, A deep learning-based framework for offensive text detection in unstructured data for heterogeneous social media, in: IEEE Access, volume 11, 2023, pp. 124484–124498. doi:10.1109/access.2023.3330081.

[104] A. Alzu'bi, L. Bani Younis, A. Abuarqoub, M. Hammoudeh, Multimodal deep learning with discriminant descriptors for offensive memes detection, ACM Journal of Data and Information Quality 15 (2023) 1–16.

[105] S. Ramamoorthy, N. Gunti, S. Mishra, S. Suryavardan, A. Reganti, P. Patwa, A. DaS, T. Chakraborty, A. Sheth, A. Ekbal, C. A, Memotion 2: Dataset on sentiment and emotion analysis of memes, in: Proceedings of De-Factify: Workshop on Multimodal Fact Checking and Hate Speech Detection, CEUR, 2022.

[106] B. Xu, T. Li, J. Zheng, M. Naseriparsa, Z. Zhao, H. Lin, F. Xia, Met-meme: A multimodal meme dataset rich in metaphors, in: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2022, pp. 2887–2899.

[107] F. Gasparini, G. Rizzi, A. Saibene, E. Fersini, Benchmark dataset of memes with text transcriptions for automatic detection of multi-modal misogynistic content, Data in brief 44 (2022) 108526.

[108] K. Maity, P. Jha, S. Saha, P. Bhattacharyya, A multitask framework for sentiment, emotion and sarcasm aware cyberbullying detection from multi-modal code-mixed memes, in: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2022, pp. 1739–1749.

[109] E. Fersini, F. Gasparini, G. Rizzi, A. Saibene, B. Chulvi, P. Rosso, A. Lees, J. Sorensen, Semeval-2022 task 5: Multimedia automatic misogyny identification, in: Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022), 2022, pp. 533–549.

[110] S. Suryawanshi, B. R. Chakravarthi, M. Arcan, P. Buitelaar, Multimodal meme dataset (multioff) for identifying offensive content in image and text, in: Proceedings of the second workshop on trolling, aggression and cyberbullying, 2020, pp. 32–41.

[111] G. Kumari, A. Sinha, A. Ekbal, A. Chatterjee, B. Vinutha, Enhancing the fairness of offensive memes detection models by mitigating unintended political bias, Journal of Intelligent Information Systems (2024) 1.

[112] S. Suryawanshi, B. R. Chakravarthi, P. Verma, M. Arcan, J. P. McCrae, P. Buitelaar, A dataset for troll classification of tamilmemes, in: Proceedings of the WILDRE5–5th workshop on indian language data: resources and evaluation, 2020, pp. 7–13.

[113] E. Fersini, G. Rizzi, A. Saibene, F. Gasparini, Misogynous meme recognition: A preliminary study, in: Lecture Notes in Computer Science, volume 13196 of *Lecture Notes in Computer Science*, 2022, pp. 279–293. doi:10.1007/978-3-031-08421-8\_19.

[114] D. Kiela, H. Firooz, A. Mohan, V. Goswami, A. Singh, P. Ringshia, D. Testuggine, The hateful memes challenge: Detecting hate speech in multimodal memes, in: Advances in Neural Information Processing Systems, volume 2020-December, 2020.

[115] M. Das, A. Mukherjee, Transfer learning for multilingual abusive meme detection, in: ACM International Conference Proceeding Series, 2023, pp. 245–250. doi:10.1145/3578503.3583607.

[116] S. R. Titli, S. Paul, Automated bengali abusive text classification: Using deep learning techniques, in: 2023 International Conference on Advances in Electronics, Communication, Computing and Intelligent Information Systems (ICAECIS), Ieee, 2023, pp. 1–6.

[117] P. K. Smith, J. Mahdavi, M. Carvalho, S. Fisher, S. Russell, N. Tippett, Cyberbullying: Its nature and impact in secondary school pupils, Journal of child psychology and psychiatry 49 (2008) 376–385.

[118] M. N. Kumar, D. M. Ahmed, J. Prashanth, V. Vinaykumar, J. A. Babu, T. K. Kumar, An efficient deep learning approach to deal with cyberbullying, in: 2023 2nd International Conference on Computational Modelling, Simulation and Optimization (ICCMSO), IEEE, 2023, pp. 253–258.

[119] J. Ji, W. Ren, U. Naseem, Identifying creative harmful memes via prompt based approach, in: ACM Web Conference 2023 - Proceedings of the World Wide Web Conference, WWW 2023, volume 3868, 2023.

[120] S. Fharook, S. Ahmed, G. Rithika, S. Budde, S. Saumya, S. Biradar, Are you a hero or a villain? a semantic role labelling approach for detecting harmful memes, in: CONSTRAINT 2022 - 2nd Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation, Proceedings of the Workshop, Constraint, 2022, pp. 19–23.

[121] T. Chakraborty, M. S. Akhtar, K. Shu, H. R. Bernard, M. Liakata, P. Nakov (Eds.), CONSTRAINT 2022 - 2nd Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation, Proceedings of the Workshop, 2022.

[122] R. Nandi, F. Alam, P. Nakov, Detecting the role of an entity in harmful memes: Techniques and their limitations, CONSTRAINT 2022 - 2nd Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation, Proceedings of the Workshop (2022) 43–54.

[123] S. Sharma, A. Kulkarni, T. Suresh, H. Mathur, P. Nakov, M. Akhtar, T. Chakraborty, Characterizing the entities in harmful memes: Who is the hero, the villain, the victim?, in: EACL 2023 - 17th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings of the Conference, 2023, pp. 2141–2155.

[124] P. Singh, A. Maladry, E. Lefever, Combining language models and linguistic information to label entities in memes, in: CONSTRAINT 2022 - 2nd Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation, Proceedings of the Workshop, Constraint, 2022, pp. 35–42.

[125] Z. Zhou, H. Zhao, J. Dong, J. Gao, X. Liu, Dd-tig at constraintacl2022: Multimodal understanding and reasoning for role labeling of entities in hateful memes, in: Proceedings of CONSTRAINT 2022 - 2nd Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation, Constraint, 2022, pp. 12–18.

[126] S. Sharma, M. K. Siddiqui, M. S. Akhtar, T. Chakraborty, Domain-aware self-supervised pre-training for label-efficient meme analysis, arXiv (2022).

[127] C. Koutlis, M. Schinas, S. Papadopoulos, Memefier: Dual-stage modality fusion for image meme, in: ICMR 2023 - Proceedings of the 2023 ACM International Conference on Multimedia Retrieval, 2023, pp. 586–591. doi:10.1145/3591106.3592254.

[128] H. B. Zia, I. Castro, G. Tyson, Racist or sexist meme? classifying memes beyond hateful, in: Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021), volume 215, 2021, pp. 215–219.

[129] J. Armenta-Segura, C.-J. Núñez Prado, G. Sidorov, A. Gelbukh, R. Román-Godínez, Ometeotlmultimodal hate speech event detection 2023: Hate speech and text-image correlation detection in real life memes using pre-trained bert models over text, in: CASE 2023 - Proceedings of the 6th Workshop on Challenges and Applications of Automated Extraction of Socio-Political Events from Text at RANLP, volume 53, 2023. doi:10.26615/978-954-452-089-2\_007.

[130] A. Aggarwal, V. Sharma, A. Trivedi, M. Yadav, C. Agrawal, D. Singh, V. Mishra, H. Gritli, Two-way feature extraction using sequential and multimodal approach for hateful meme classification, Complexity 2021 (2021).

[131] G. Arya, M. Hasan, A. Bagwari, N. Safie, S. Islam, F. Ahmed, A. De, M. Khan, T. Ghazal, Multimodal hate speech detection in memes using contrastive language-image pre-training, IEEE Access 12 (2024) 22359–22375.

[132] P. Aggarwal, P. Chawla, P. Das, P. Saha, B. Mathew, T. Zesch, A. Mukherjee, Hateproof: Are hateful meme detection systems really robust?, in: ACM Web Conference 2023 - Proceedings of the World Wide Web Conference, WWW 2023, volume 3734, 2023. doi:10.1145/3543507.3583356.

[133] P. Aggarwal, M. Liman, D. Gold, T. Zesch, Vl-bert+: Detecting protected groups in hateful multimodal memes, in: WOAH 2021 - 5th Workshop on Online Abuse and Harms, Proceedings of the Workshop, volume 207, 2021.

[134] M. Ahmed, N. Bhadani, I. Chakraborty, Hateful meme prediction model using multimodal deep learning, in: 2021 International Conference on Computing, Communication and Green Engineering, CCGE 2021, volume 2, 2021. doi:10.1109/ccge50943.2021.9776440.

[135] A. Bhat, V. Varshney, V. Bajlotra, V. Gupta, Detection of hatefulness in memes using unimodal and multimodal techniques, in: Proceedings - 2022 6th International Conference on Intelligent Computing and Control Systems, ICICCS 2022, volume 65, 2022.

doi:10.1109/iciccs53718.2022.9788376.

[136] A. Bhat, V. Vashisht, V. Sahni, S. Meena, Hate speech detection using multimodal meme analysis, in: Proceedings of the 2nd International Conference on Applied Artificial Intelligence and Computing, ICAAIC 2023, volume 1137, 2023. doi:10.1109/icaaic56838.2023.10140393.

[137] N. Bi, Y.-C. Huang, C.-C. Han, J.-J. Hsu, You know what i meme: Enhancing people's understanding and awareness of hateful memes using crowdsourced explanations, in: Proceedings of the ACM on Human-Computer Interaction, volume 7, 2023. doi:10.1145/3579593.

[138] E. Blaier, I. Malkiel, L. Wolf, Caption enriched samples for improving hateful memes detection, in: EMNLP 2021 - 2021 Conference on Empirical Methods in Natural Language Processing, Proceedings, volume 9350, 2021.

[139] R. Cao, R.-W. Lee, W.-H. Chong, J. Jiang, Prompting for multimodal hateful meme classification, in: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, volume 321, 2022.

[140] R. Cao, M. Hee, A. Kuek, W.-H. Chong, R.-W. Lee, J. Jiang, Pro-cap: Leveraging a frozen vision-language model for hateful meme detection, in: MM 2023 - Proceedings of the 31st ACM International Conference on Multimedia, volume 5244, Association for Computing Machinery, 2023. doi:10.1145/3581783.3612496.

[141] A. Chhabra, D. Vishwakarma, Multimodal hate speech detection via multi-scale visual kernels and knowledge distillation architecture, Engineering Applications of Artificial Intelligence 126 (2023).

[142] M. Constantin, D.-S. Parvu, C. Stanciu, D. Ionascu, B. Ionescu, Hateful meme detection with multimodal deep neural networks, in: ISSCS 2021 - International Symposium on Signals, Circuits and Systems, volume 9497374, 2021. doi:10.1109/isscs52333.2021.9497374.

[143] T. Deshpande, N. Mani, An interpretable approach to hateful meme detection, in: ICMI 2021 - Proceedings of the 2021 International Conference on Multimodal Interaction, volume 723, 2021. doi:10.1145/3462244.3479949.

[144] H. Fang, F. Zhu, J. Han, S. Hu, Multimodal hateful memes detection via image caption supervision, in: Proceedings - 2022 IEEE SmartWorld, Ubiquitous Intelligence and Computing, Autonomous and Trusted Vehicles, Scalable Computing and Communications, Digital Twin, Privacy Computing, Metaverse, SmartWorld/UIC/ATC/ScalCom/DigitalTwin/PriComp/Metaverse 2022, volume 1530, 2022. doi:10.1109/SmartWorld-UIC-ATC-ScalCom-DigitalTwin-PriComp-Metaverse56740.2022.00221.

[145] B. Gaikwad, B. Kurma, M. Patwardhan, S. Karande, N. Pedanekar, Can a pretrained language model make sense with pretrained neural extractors? an application to multimodal classification, in: CEUR Workshop Proceedings, volume 3168, 2022. doi:10.1109/iceic57457.2023.10049865.

[146] A. Goswami, A. Rawat, S. Tongaria, S. Jhingran, Detection of hate speech in multi-modal social post, in: Artificial Intelligence, Blockchain, Computing and Security: Volume 1, volume 1, 2023. doi:10.1201/9781003393580-50.

[147] M. Hee, R.-W. Lee, W.-H. Chong, On explaining multimodal hateful meme detection models, in: WWW 2022 - Proceedings of the ACM Web Conference 2022, volume 3651, 2022. doi:10.1145/3485447.3512260.

[148] A. Kiran, M. Shetty, S. Shukla, V. Kerenalli, B. Das, Getting around the semantics challenge in hateful memes, in: Lecture Notes on Data Engineering and Communications Technologies, volume 142, 2023. doi:10.1007/978-981-19-3391-2\_26.

[149] V. Kougia, J. Pavlopoulos, Multimodal or text? retrieval or bert? benchmarking classifiers for the shared task on hateful memes, in: WOAH 2021 - 5th Workshop on Online Abuse and Harms, Proceedings of the Workshop, volume 220, 2021.

[150] V. Kougia, S. Fetzel, T. Kirchmair, E. Çano, S. Baharlou, S. Sharifzadeh, B. Roth, Memegraphs: Linking memes to knowledge graphs, in: Lecture Notes in Computer Science, volume 14187 of *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Springer, 2023, pp. 534–551. doi:10.1007/978-3-031-41676-7\_31.

[151] R. K.-W. Lee, R. Cao, Z. Fan, J. Jiang, W.-H. Chong, Disentangling hate in online memes, in: Proceedings of the 29th ACM International Conference on Multimedia, volume 5138, Association for Computing Machinery, 2021, pp. 5138–5147. doi:10.1145/3474085.3475625.

[152] X. Liang, Y.-C. Huang, W. Liu, H. Zhu, Z. Liang, L. Chen, Trican: Multi-modal hateful memes detection with triplet-relation information cross-attention network, in: Proceedings of the International Joint Conference on Neural Networks, volume 2022-July, 2022. doi:10.1109/ijcnn55064.2022.9892164.

[153] Z. Ma, S. Yao, L. Wu, S. Gao, Y. Zhang, Hateful memes detection based on multi-task learning, Mathematics 10 (2022).

[154] G. MacRayo, W. Casino, J. Dalangin, J. Gabriel Gahoy, A. Christian Reyes, C. Vitto, M. Abisado, S. Lor Huyo-A, G. Avelino Sampedro, Please be nice: A deep learning based approach to content moderation of internet memes, in: 2023 International Conference on Electronics, Information, and Communication, ICEIC 2023, volume 0, Ieee, 2023, pp. 1–5. doi:10.1109/iceic57457.2023.10049865.

[155] L. Mookdarsanit, P. Mookdarsanit, Combating the hate speech in thai textual memes, Indonesian Journal of Electrical Engineering and Computer Science 21 (2021) 1493–1502.

[156] A. Nayak, A. Agrawal, Detection of hate speech in social media memes: A comparative analysis, in: Proceedings of the 2022 3rd International Conference on Intelligent Computing, Instrumentation and Control Technologies: Computational Intelligence for Smart Systems, ICICICT 2022, volume 1179, 2022. doi:10.1109/icicict54557.2022.9917633.

[157] Y. Qu, X. He, S. Pierson, M. Backes, Y. Zhang, S. Zannettou, On the evolution of (hateful) memes by means of multimodal contrastive learning, in: Proceedings - IEEE Symposium on Security and Privacy, volume 2023-May, 2023. doi:10.1109/sp46215.2023.10179315.

[158] Y. Qu, X. Shen, X. He, M. Backes, S. Zannettou, Y. Zhang, Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models, in: Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security, Ccs '23, Association for Computing Machinery, 2023, p. 3403–3417. doi:10.1145/3576915.3616679.

[159] A. Sethi, U. Kuchhal, Anjum, R. Katarya, Study of various techniques for the classification of hateful memes, in: 2021 6th International Conference on Recent Trends on Electronics, Information, Communication and Technology, RTEICT 2021, volume 675, 2021. doi:10.1109/rteict52294.2021.9573926.

[160] L. Wanbo, L. Suying, Research on multi-modal hateful meme detection, in: ACM International Conference Proceeding Series, volume 3470385, 2021. doi:10.1145/3469213.3470385.

[161] P. Wu, W. Mebane, Marmot a deep learning framework for constructing multimodal representations for vision-and-language tasks, Computational Communication Research 4 (2022).

[162] L. Zhang, L. Jin, X. Sun, G. Xu, Z. Zhang, X. Li, N. Liu, Q. Liu, S. Yan, Tot: Topology-aware optimal transport for multimodal hate detection, in: Proceedings of the 37th AAAI Conference on Artificial Intelligence, AAAI 2023, volume 37, 2023.

[163] Y. Zhou, Z. Chen, H. Yang, Multimodal learning for hateful memes detection, 2021 IEEE International Conference on Multimedia and Expo Workshops, ICMEW 2021 28 (2021).

[164] J. Zhu, R.-W. Lee, W. Chong, Multimodal zero-shot hateful meme detection, in: ACM International Conference Proceeding Series, volume 382, 2022. doi:10.1145/3501247.3531557.

[165] P. Aggarwal, J. Mehrabanian, W. Huang, O. Alacam, T. Zesch, Text or image? what is more important in cross-domain generalization capabilities of hate meme detection models?, arXiv (2024).

[166] Y. Chen, F. Pan, Multimodal detection of hateful messages using visual-linguistic pre-trained deep learning models, Research Square (2022).

[167] A. Das, J. S. Wahi, S. Li, Detecting hate speech in multi-modal memes, arXiv (2020).

[168] I. Evtimov, R. Howes, B. Dolhansky, H. Firooz, C. C. Ferrer, Adversarial evaluation of multimodal models under realistic gray box assumptions, arXiv (2020).

[169] A. Gao, B. Wang, J. Yin, Y. Tian, Hateful memes challenge: An enhanced multimodal framework, arXiv (2021).

[170] C. Jennifer, F. Tahmasbi, J. Blackburn, S. Zannettou, E. De Cristofaro, Feels bad man: Dissecting automated hateful meme detection through the lens of facebook's challenge, arXiv (2022).

[171] W. Jin, L. Wilhelm, The hateful memes challenge next move, arXiv (2022).

[172] Y. Li, Z. Zhang, H. Huang, Enhance multimodal model performance with data augmentation: facebook hateful meme challenge solution, arXiv (2021).

[173] P. Lippe, N. Holla, S. Chandra, S. Rajamanickam, G. Antoniou, E. Shutova, H. Yannakoudakis, A multimodal framework for the detection of hateful memes, arXiv preprint arXiv:2012.12871 (2020).

[174] J. Mei, J. Chen, W. Lin, B. Byrne, M. Tomalin, Improving hateful memes detection via learning hatefulness-aware embedding space through retrieval-guided contrastive learning, arXiv (2023).

[175] Y. Miyanishi, M. Le Nguyen, Causal intersectionality and dual form of gradient descent for multimodal analysis: a case study on hateful memes, arXiv (2023).

[176] N. Muennighoff, Vilio: State-of-the-art visio-linguistic models applied to hateful memes, arXiv preprint arXiv:2012.07788 (2020).

[177] V. Sandulescu, Detecting hateful memes using a multimodal deep ensemble, arXiv preprint arXiv:2012.13235 (2020).

[178] M.-H. Van, X. Wu, Detecting and correcting hate speech in multimodal memes with large visual language model, arXiv (2023).

[179] R. Velioglu, J. Rose, Detecting hate speech in memes using multimodal deep learning approaches: Prize-winning solution to hateful memes challenge, arXiv preprint arXiv:2012.12975 (2020).

[180] J. Yuan, Y. Yu, G. Mittal, S. Sajeev, M. Chen, Rethinking multimodal content moderation from an asymmetric angle with mixed-modality, arXiv (2023).

[181] W. Zhang, G. Liu, Z. Li, F. Zhu, Hateful memes detection via complementary visual and linguistic networks, arXiv (2020).

[182] B. Zhao, A. Zhang, B. Watson, G. Kearney, I. Dale, A review of vision-language models and their performance on the hateful memes challenge, arXiv (2023).

[183] X. Zhong, Classification of multimodal hate speech -the winning solution of hateful memes challenge, arXiv (2020).

[184] Y. Zhou, Z. Chen, Multimodal learning for hateful memes detection, arXiv (2020).

[185] R. Zhu, Enhance multimodal transformer with external label and in-domain pretrain: Hateful meme challenge winning solution, arXiv preprint arXiv:2012.08290 (2020).

[186] G. Burbi, A. Baldrati, L. Agnolucci, M. Bertini, A. Del Bimbo, Mapping memes to words for multimodal hateful meme classification, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 2832–2836. doi:10.1109/iccvw60793.2023.00303.

[187] E. Fersini, F. Gasparini, S. Corchs, Detecting sexist meme on the web: A study on textual and visual cues, in: 2019 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW), Ieee, 2019, pp. 226–231.

[188] G. Attanasio, D. Nozza, F. Bianchi, Milanlp at semeval-2022 task 5: Using perceiver io for detecting misogynous memes with text and image modalities, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 654–662.

[189] M. Behzadi, A. Derakhshan, I. Harris, Mitra behzadi at semeval-2022 task 5: Multimedia automatic misogyny identification method based on clip, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 724–727.

[190] L. Chen, H. Chou, Rit boston at semeval-2022 task 5: Multimedia misogyny detection by using coherent visual and language features from clip model and data-centric ai principle, in: Proceedings of SemEval 2022 - 16th Intl. Workshop on Semantic Evaluation, 2022, pp. 636–641.

[191] Y. Gu, I. Castro, G. Tyson, Mmvae at semeval-2022 task 5: A multi-modal multi-task vae on misogynous meme detection, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 700–710.

[192] M. Kalkenings, T. Mandl, University of hildesheim at semeval-2022 task 5: Combining deep text and image models for multimedia misogyny detection, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 718–723.

[193] D. Obeidat, H. Nammas, M. Abdullah, Just_one at semeval-2023 task 10: Explainable detection of online sexism (edos), in: Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), 2023, pp. 526–531.

[194] A. Paraschiv, M. Dascalu, D.-C. Cercel, Upb at semeval-2022 task 5: Enhancing uniter with image sentiment and graph convolutional networks for multimedia automatic misogyny identification, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 618–625.

[195] J. Ravagli, L. Vaiani, Jrlv at semeval-2022 task 5: The importance of visual elements for misogyny identification in memes, in: SemEval 2022 - 16th International Workshop on Semantic Evaluation, Proceedings of the Workshop, 2022, pp. 610–617.

[196] G. Rizzi, F. Gasparini, A. Saibene, P. Rosso, E. Fersini, Recognizing misogynous memes: Biased models and tricky archetypes, Information Processing and Management 60 (2023) 103474.

[197] S. Singh, A. Haridasan, R. Mooney, "female astronaut: Because sandwiches won't make themselves up there!": Towards multi-modal misogyny detection in memes, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2023, pp. 150–159.

[198] P. Tung, N. Viet, N. Anh, P. Hung, Semimemes: A semi-supervised learning approach for multimodal memes analysis, in: Lecture Notes in Computer Science, volume 14162 of *Lecture Notes in Computer Science*, 2023, pp. 565–577. doi:10.1007/978-3-031-41456-5\_43.

[199] N. K. Singh, P. Das, A. Manderna, S. Chand, Devi deep learning framework for misogyny identification in multimodal data, Research Square (2023).

[200] J. Drakett, B. Rickett, K. Day, K. Milnes, Old jokes, new media–online sexism and constructions of gender in internet memes, Feminism & psychology 28 (2018) 109–127.

[201] A. Alzu'bi, L. Bani Younis, A. Abuarqoub, M. Hammoudeh, Multimodal deep learning with discriminant descriptors for offensive memes detection, in: Journal of Data and Information Quality, volume 15, 2023, p. 3597308. doi:10.1145/3597308.

[202] A. Aman, G. Krishna, T. Anand, A. Lal, Identification of offensive content in memes, in: Lecture Notes in Networks and Systems, volume 290, 2021, pp. 438–445. doi:10.1007/978-981-16-4486-3\_49.

[203] A. Baruah, K. Das, F. Barbhuiya, K. Dey, Iiitg-adbu at semeval-2020 task 8: A multimodal approach to detect offensive, sarcastic and humorous memes, in: Proceedings of SemEval 2020 - 14th International Workshops on Semantic Evaluation, 2020, pp. 885–890.

[204] I. Bejan, Memosys at semeval-2020 task 8: Multimodal emotion analysis in memes, in: Proceedings of SemEval 2020 – 14th International Workshops on Semantic Evaluation, 2020, pp. 1172–1178.

[205] S. Boinepelli, M. Shrivastava, V. Varma, Sisiiith at semeval-2020 task 8: An overview of simple text classification methods for meme analysis, in: Proceedings of SemEval 2020 – 14th International Workshops on Semantic Evaluation, 2020, pp. 1190–1194.

[206] A.-M. Bucur, A. Cosma, I.-B. Iordache, Blue at memotion 2.0 2022: You have my image, my text and my transformer, in: CEUR Workshop Proceedings, volume 3168, 2022.

[207] V. Sharma, V. Kushwaha, S. Jaiswal, G. Nandi, Meme detection for sentiment analysis and human robot interactions using multiple modes, in: 9th IEEE Uttar Pradesh Section International Conference on Electrical, Electronics and Computer Engineering, UPCON 2022, 2022. doi:10.1109/upcon56432.2022.9986453.

[208] G. de la Peña Sarracén, P. Rosso, A. Giachanou, Prhlt-upv at semeval-2020 task 8: Study of multimodal techniques for memes analysis, in: 14th International Workshops on Semantic Evaluation, SemEval 2020 - co-located 28th International Conference on Computational Linguistics, COLING 2020, Proceedings of the Workshop, 2020, pp. 908–915.

[209] R. Giri, S. Gupta, U. Gupta, An approach to detect offence in memes using natural language processing(nlp) and deep learning, in: 2021 International Conference on Computer Communication and Informatics, ICCCI 2021, 2021, p. 9402406. doi:10.1109/iccci50826.2021.9402406.

[210] A. Gupta, H. Kataria, S. Mishra, T. Badal, V. Mishra, Bennettnlp at semeval-2020 task 8: Multimodal sentiment classification using hybrid hierarchical classifier, in: Proceedings of SemEval 2020 - 14th International Workshops on Semantic Evaluation, 2020, pp. 1085–1093.

[211] S. Hakimov, G. Cheema, R. Ewerth, Tib-va at semeval-2022 task 5: A multimodal architecture for the detection and classification of misogynous memes, in: Proceedings of SemEval 2022 - 16th International Workshop on Semantic Evaluation, 2022, pp. 756–760.

[212] E. Hossain, M. Hoque, M. Hossain, An inter-modal attention framework for multimodal offense detection, in: Lecture Notes in Networks and Systems, volume 569 Lnns, 2023, pp. 853–862. doi:10.1007/978-3-031-19958-5\_81.

[213] K. Myilvahanan, B. Shashank, T. Raj, C. Attanti, S. Sahay, A study on deep learning based classification and identification of offensive memes, in: Proceedings of the 3rd International Conference on Trends in Electronics and Informatics, ICOEI 2019, 2023, pp. 214–218. doi:10.1109/icoei.2019.8862647.

[214] T. Nguyen, N. Pham, N. Nguyen, H. Nguyen, L. Nguyen, Y.-G. Kim, Hcilab at memotion 2.0 2022: Analysis of sentiment, emotion and intensity of emotion classes from meme images using single and multi modalities, in: CEUR Workshop Proceedings, volume 3168, 2022.

[215] K. Phan, G.-S. Lee, H.-J. Yang, S.-H. Kim, Little flower at memotion 2.0 2022: Ensemble of multi-modal model using attention mechanism in memotion analysis, in: CEUR Workshop Proceedings, volume 3168, 2022.

[216] L. Shang, Y. Zhang, Y. Zha, Knowmeme: A knowledge-enriched graph neural network solution to offensive meme detection, in: Proceedings - IEEE 17th International Conference on eScience, eScience 2021, 2021, pp. 186–195. doi:10.1109/eScience51609.2021.00029.

[217] U. Walinska, J. Potoniec, Urszula walińska at semeval-2020 task 8: Fusion of text and image features using lstm and vgg16 for memotion analysis, in: Proceedings of SemEval 2020 – 14th International Workshops on Semantic Evaluation, 2020, pp. 1215–1220.

[218] W. Yu, D. Kolossa, wentaorub at memotion 3: Ensemble learning for multi-modal meme classification, in: CEUR Workshop Proceedings, volume 3555, 2022.

[219] Q. Zhong, Q. Wang, J. Liu, Combining knowledge and multi-modal fusion for meme classification, in: Lecture Notes in Computer Science, volume 13141, 2022, pp. 599–611. doi:10.1007/978-3-030-98358-1\_47.

[220] S. Pramanick, M. S. Akhtar, T. Chakraborty, Exercise? i thought you said 'extra fries' ☺: Leveraging sentence demarcations and multi-hop attention for meme affect analysis, arXiv (2021).

[221] D. Gaurav, S. Shandilya, S. Tiwari, A. Goyal, A machine learning method for recognizing invasive content in memes, in: Communications in Computer and Information Science, 2020, pp. 195–213. doi:10.1007/978-3-030-65384-2\_15.

[222] S. Gundapu, R. Mamidi, Detection of propaganda techniques in visuo-lingual metaphor in memes, arXiv (2022).

[223] D. Abujaber, A. Qarqaz, M. Abdullah, Lecun at semeval-2021 task 6: Detecting persuasion techniques in text using ensembled pretrained transformers and data augmentation, in: Proceedings of SemEval 2021 - 15th Intl. Workshop on Semantic Evaluation, 2021, pp. 1068–1074.

[224] F. Alam, H. Mubarak, W. Zaghouani, G. Da San Martino, P. Nakov, Overview of the wanlp 2022 shared task on propaganda detection in arabic, WANLP 2022 - 7th Arabic Natural Language Processing - Proceedings of the Workshop (2022) 108–118.

[225] J. Cui, L. Li, X. Zhang, J. Yuan, Multimodal propaganda detection via anti-persuasion prompt enhanced contrastive learning, in: ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings, 2023, p. 0. doi:10.1109/icassp49357.2023.10096771.

[226] T. Hossain, J. Naim, F. Tasneem, R. Tasnia, A. Chy, Csecu-dsg at semeval-2021 task 6: Orchestrating multimodal neural architectures for identifying persuasion techniques in texts and images, in: Proceedings of SemEval 2021 - 15th International Workshop on Semantic Evaluation, 2021, pp. 1088–1095.

[227] P. Li, X. Li, X. Sun, 1213li at semeval-2021 task 6: Detection of propaganda with multi-modal attention and pre-trained models, in: SemEval 2021 - 15th International Workshop on Semantic Evaluation, Proceedings of the Workshop, 2021, pp. 1032–1036.

[228] C. Liu, G. Geigle, R. Krebs, I. Gurevych, Figmemes: A dataset for figurative language identification in politically-opinionated memes, in: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, 2022, pp. 7069–7086.

[229] H. Wu, X. Li, L. Li, Q. Wang, Propaganda techniques detection in low-resource memes with multi-modal prompt tuning, in: Proceedings - IEEE International Conference on Multimedia and Expo, 2022, p. 0. doi:10.1109/icme52920.2022.9859642.

[230] X. Zhu, J. Wang, X. Zhang, Ynu-hpcc at semeval-2021 task 6: Combining albert and text-cnn for persuasion detection in texts and images, in: SemEval 2021 - 15th International Workshop on Semantic Evaluation, Proceedings of the Workshop, 2021, pp. 1045–1050.

[231] P. Chen, L. Zhao, Y. Piao, H. Ding, X. Cui, Multimodal visual-textual object graph attention network for propaganda detection in memes, Multimedia Tools and Applications 1 (2023) 1–10.

[232] J. Cui, L. Li, X. Tao, Be-or-not prompt enhanced hard negatives generating for memes category detection, in: Proceedings - IEEE International Conference on Multimedia and Expo, 2023.

[233] M. Das, S. Banerjee, A. Mukherjee, hate-alertdravidianlangtech-acl2022: Ensembling multi-modalities for tamil trollmeme classification, arXiv (2022).

[234] S. U. Hegde, A. Hande, R. Priyadharshini, B. Bharathi, B. R. Chakravarthi, Do images really do the talking? analysing the significance of images in tamil troll meme classification, arXiv (2021).

[235] R. N. Nandi, F. Alam, P. Nakov, Teamxdravidianlangtech-acl2022: A comparative analysis for troll-based meme classification, arXiv (2022).

[236] H. Zhang, H. Xu, X. Wang, Q. Zhou, S. Zhao, J. Teng, Mintrec: A new dataset for multimodal intent recognition, in: Proceedings of the 30th ACM International Conference on Multimedia, 2022, pp. 1688–1697.

[237] N. Mirzoeff, White Sight: Visual Politics and Practices of Whiteness, MIT Press, 2023.

[238] S. Sekimoto, C. Brown, Race and multimodality: An introduction to the special issue, 2023.

[239] F. Yus, Multimodality in memes: A cyberpragmatic approach, Analyzing digital disc.: New insights and future dir. (2019) 105–131.

[240] J. L. Lemke, Metamedia literacy: Transforming meanings and media, Handbook of literacy and technology: Transformations in a post-typographic world 283301 (1998).

[241] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, I. Sutskever, Learning transferable visual models from natural language supervision, in: International conference on machine learning, PMLR, 2021, pp. 8748–8763.

[242] F. Wu, B. Gao, X. Pan, L. Li, Y. Ma, S. Liu, Z. Liu, Fuser: An enhanced multimodal fusion framework with congruent reinforced perceptron for hateful memes detection, Information Processing & Management 61 (2024) 103772.

[243] W. Wu, H. Li, H. Wang, K. Q. Zhu, Probase: A probabilistic taxonomy for text understanding, in: Proceedings of the 2012 ACM SIGMOD international conference on management of data, 2012, pp. 481–492.

[244] R. Tommasini, F. Ilievski, T. Wijesiriwardene, Imkg: The internet meme knowledge graph, in: European Semantic Web Conference, Springer, 2023, pp. 354–371.

[245] D. Vrandečić, M. Krötzsch, Wikidata: a free collaborative knowledgebase, Communications of the ACM 57 (2014) 78–85.

[246] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, Z. Ives, Dbpedia: A nucleus for a web of open data, in: international semantic web conference, Springer, 2007, pp. 722–735.

[247] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, J. Taylor, Freebase: a collaboratively created graph database for structuring human knowledge, in: Proceedings of the 2008 ACM SIGMOD international conference on Management of data, 2008, pp. 1247–1250.

[248] L. Bates, P. E. Christensen, P. Nakov, I. Gurevych, A template is all you meme, arXiv preprint arXiv:2311.06649 (2023).

[249] M. K. Scheuerman, J. A. Jiang, C. Fiesler, J. R. Brubaker, A framework of severity for harmful content online, Proceedings of the ACM on Human-Computer Interaction 5 (2021) 1–33.

[250] F. Jahanbakhsh, A. X. Zhang, A. J. Berinsky, G. Pennycook, D. G. Rand, D. R. Karger, Exploring lightweight interventions at posting time to reduce the sharing of misinformation on social media, Proceedings of the ACM on Human-Computer Interaction 5 (2021) 1–42.

[251] T. Kumarage, A. Bhattacharjee, J. Garland, Harnessing artificial intelligence to combat online hate: Exploring the challenges and opportunities of large language models in hate speech detection, arXiv preprint arXiv:2403.08035 (2024).

[252] D. M. Beskow, S. Kumar, K. M. Carley, The evolution of political memes: Detecting and characterizing internet memes with multi-modal deep learning, Information Processing & Management 57 (2020) 102170.