

What Makes a Meme a Meme? Identifying Memes for Memetics-Aware Dataset Creation

Muzhaffar Hazman¹, Susan McKeever², Josephine Griffith²

¹University of Galway, Galway Ireland.

²TU Dublin, Dublin Ireland.

m.hazman1@universityofgalway.ie, susan.mckeever@tudublin.ie, josephine.griffith@universityofgalway.ie

Abstract

WARNING: This paper contains memes that may be offensive to some readers.

Multimodal Internet Memes are now a ubiquitous fixture in online discourse. One strand of meme-based research is the classification of memes according to various affects, such as sentiment and hate, supported by manually compiled meme datasets. Understanding the unique characteristics of memes is crucial for meme classification. Unlike other user-generated content, memes spread via memetics, i.e. the process by which memes are imitated and transformed into symbols used to create new memes. In effect, there exists an ever-evolving pool of visual and linguistic symbols that underpin meme culture and are crucial to interpreting the meaning of individual memes. The current approach of training supervised learning models on static datasets, without taking memetics into account, limits the depth and accuracy of meme interpretation. We argue that meme datasets must contain genuine memes, as defined via memetics, so that effective meme classifiers can be built. In this work, we develop a meme identification protocol which distinguishes meme from non-memetic content by recognising the memetics within it. We apply our protocol to random samplings of the leading 7 meme classification datasets and observe that more than half (50.4%) of the evaluated samples were found to contain no signs of memetics. Our work also provides a meme typology grounded in memetics, providing the basis for more effective approaches to the interpretation of memes and the creation of meme datasets.

1 Introduction

Internet Memes have become a staple of expression and interaction within digital communities. Although commonly recognised as multimodal humorous jokes, memes present a medium to participate in a global digital culture (Shifman 2014b), spread ideologies (Zannettou et al. 2018), debate political issues (Pearce and Hajizada 2014), and even distribute propaganda (Dimitrov et al. 2021). The expressive nature of memes has led to recent interest in classifying them by the various affects they convey including sentiment, hate, humour, sarcasm and offensiveness. The goals underlying meme classification are varied, e.g.: identifying hateful memes to help curb the proliferation of hateful rhetoric (Sharma et al. 2022), detecting signs of cyberbullying (Shang et al. 2021b),

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

recognising trolling behaviour (Suryawanshi et al. 2020), and discerning the sentiment conveyed by memes (Sharma et al. 2020).

Meme classification datasets provide memes labelled with various affects to aid in the development of meme classifiers via supervised learning. The current state-of-the-art in collecting and labelling sample memes is almost a completely manual process. Since not all content on the Internet are memes, the authors of these datasets have indicated the need to identify memes from non-memetic content during data collection. However, none of the 12 publicly available peer-reviewed meme classification datasets we surveyed disclosed any standards or criteria used to identify memes from non-memetic content. For clarity, we use *meme classification* to describe the task of classifying the various affects conveyed by memes as posed by these datasets, while we refer to *meme identification* as the task of selecting memes from non-memetic content (non-memes).

Beyond meme classification datasets, work in both the fields of computational and philological research identify, characterise, and interpret memes through the lens of *memetics* (Courtois and Frissen 2022; Segev et al. 2015), the process by which memes are derived from previous memes through imitation and variation (Shifman 2014b). Through memetics, memes form groups with other memes via shared *memetic elements*, which are components or traits common between them. These memetic elements not only indicate that a piece of content has been created from prior memes, but also often encapsulate the meaning and context of a meme (Wiggins and Bowers 2015). The repeated imitation and reuse of memetic elements in millions of memes gives each element its meaning (Dubey et al. 2018; Qu et al. 2022), adding to the constantly evolving collection of visual symbols shared between memes (Zannettou et al. 2018).

Interpreting memetic elements is an additional challenge necessary in classifying memes that does not exist when classifying other multimodal content. Understanding the meaning conveyed by these elements is often subtle, and requires familiarity with how each element has been used in preceding memes (Wiggins and Bowers 2015). Consider, for example, the memes in Fig. 1, each using a distinct but similar *meme template* (which Nissenbaum and Shifman (2017) describes as a recurring image with a text overlay that typically addresses a particular subject or circumstance). The *Socially*



Figure 1: Example memes using the popular meme templates: (a) Socially Awesome Penguin and (b) Socially Awkward Penguin

Awesome Penguin meme template, as shown in Fig. 1(a) subverts an awkward social scenario with a charismatic response. This penguin and red background has come to symbolise “social awesomeness” and is an inversion of the *Socially Awkward Penguin* meme template¹, see Fig. 1(b). The two templates are nearly identical; whether the meme conveys awkwardness or charisma depends solely on which direction the penguin is facing and the colour of its background. The relationship between penguins and either awkwardness or charisma in social situations is obscure. Like many memetic elements, the specific meaning conveyed by these penguins arises solely through their reuse and variation across millions of memes. Similarly to how new memes are derived from previous memes, memetic elements can be derivations of previously established memetic elements, such as how Fig. 1(a) was derived from Fig. 1(b). The new element incorporates some of the meaning of the original element, but is given added meaning through visual editing (Qu et al. 2022).

In contrast, non-memetic content conveys meaning without relying on memetic elements or on the intended reader being familiar with the larger meme culture. Combining memetic and non-memetic content into meme datasets hinders the development of meme classifiers, as it undermines meme-specific challenges in multimodal affective classification. Current approaches to meme classification are not adequate given that the datasets on which they trained contain both meme and non-meme content. To enable research that addresses challenges unique to meme classification, datasets of genuine verifiable memes are required. To create such datasets, a clear definition of memes and a repeatable methodology for meme identification are needed.

The main contributions of this work are to provide a theoretically grounded definition of memes and repeatable meme identification protocol based on the concept of memetics. We propose that this definition should be used in the creation of meme datasets that are created for the task of affective classification by supervised machine learning models. By surveying the data collection approaches that were used for existing meme classification datasets, we show that none reported any replicable methodology for meme identification or for the filtration of non-memetic content. We present a protocol that identifies a piece of multimodal user-generated content either

¹<https://knowyourmeme.com/memes/socially-awkward-penguin>

as a meme or non-meme by verifying its memetic nature. By applying our protocol to samples from the surveyed datasets, we show that current meme classification datasets contain both memes and non-memetic content alike.

Our paper is organised as follows: Section 2 provides an overview of how memes are conceptualised in the literature in two key ways. S2.1 explores how memetic elements signify the memetic nature of every meme. S2.2 illustrates how memetic elements are also used to computationally characterise large collections of memes. Section 3 introduces our novel typology and protocol for meme identification. In Section 4, we present a survey of 12 contemporary meme datasets, revealing that while the authors of these datasets regard memes as being distinct from non-memes, none reported a repeatable protocol for meme identification. Section 5 applies our protocol to seven of these datasets to ascertain whether they solely comprise memes or include non-memetic content. We also analyse instances of non-memes found within meme datasets and discuss the implications of their inclusion in meme datasets. Section 6 positions our work within the context of meme research, outlining its limitations and future research directions.

2 Background

2.1 Existing Definitions of Memes

A meme distinguishes itself from non-memetic content by the way it is created. Unlike other digital content, a meme remains identifiable as a meme and as a member of a larger group of memes due to the observable characteristics it shares with other memes (Wiggins and Bowers 2015). Although each meme is an individual unit of expression, this *sharedness* exhibits how each meme is derived from other memes, through *memetics* (Cannizzaro 2016; Wiggins 2019).

Memetics in this context entails a large-scale collective process of imitation and variation on popular cultural media. This process is part of a larger participatory culture online (Wiggins and Bowers 2015) where popular media is “remixed” to create a meme that expresses the perspectives of its creator (Wiggins 2019). Cannizzaro (2016) describes this process as a collective *habit* of translating elements from previous memes to create new interpretations, becoming new cultural resources (Nissenbaum and Shifman 2017) that represent shared collective identities (Gal, Shifman, and Kampf 2016).

Inspired by several eminent perspectives on memetics, Shifman (2013) reconciled these into the modern digital sense of memes. In doing so, they proposed what has become an influential definition of Internet memes (hereafter, we refer to as “**Shifman’s definition**”):

An internet meme [is defined] as a group of digital items that: (a) share common characteristics of content, form, and/or stance; (b) are created with awareness of each other; and (c) are circulated, imitated, and transformed via the Internet by multiple users.

First, note how Shifman spells out the role memetics play in defining Internet memes: every meme is created from “imitating and transforming” other memes. The derivation of memes from other memes is central to how memes are created

and spread (Shifman 2014b). Shifman uses this derivative process to contrast memes against *virals*: singular units of cultural media that, while popular, do not inspire derivatives. Whereas related memes are distinct from one another. Similarly, Wiggins and Bowers (2015) suggest that digital media only become a meme once they generate “imitations, remixes, and iterations”. This imitation and transformation relates to how memes are created in groups and that a meme must share “common characteristics” with other memes within its group. For example, consider the group of memes in Fig.2 which all share the same background.

Second, the memetics process creates groups of memes that are related by some shared characteristics in one of three “dimensions” (Shifman 2013): *Content*, *Form*, and *Stance*. We refer to these components, which are shared and reused between related but distinct memes, as *memetic elements*, such as the penguins in Fig. 1 and the reused background in Fig. 2. The presence of memetic elements within a piece of digital content indicates that the content was created memetically from a preceding meme, and thus can be verifiably identified as a meme.

Third, we note that modality, in terms of using or combining visual, textual, or verbal elements, does not, in itself, define memes. Although memes are often multimodal, a meme can also be, in its entirety, textual (He et al. 2019) or visual (Shifman 2014a). Thus, memetics create memes that could manifest as unimodal or multimodal content; conversely, not all multimodal content is memetic.

2.2 Memetic Elements in Previous Works

In this section, we highlight several approaches leveraging various types of memetic elements to computationally infer the relationship between one memes and another. Most works that seek to analyse multimodal Internet memes at scale have chosen to focus on features along the *Form* dimension of Shifman’s definition to identify shared common characteristics to link related memes. Memetics through this dimension are grounded in a shared visual format or linguistic syntax, i.e. the memes share observable visual, audible, or textual features.

To assess the evolution and influence of memes across various online communities, Zannettou et al. (2018) presents an early example of using shared observable characteristics to link related memes. First, they represented memes based on overall visual similarity. They created vector representations of memes collected from various online communities. They applied a distance metric based on this representation to cluster memes based on visual similarity to a separate set of annotated memes they had collected from meme annotation sites, such as KnowYourMeme.com. Their work also shows how visually similar memes originate from certain online communities and spread through other communities.

To track the evolution of memes across political conversations, Beskow, Kumar, and Carley (2020) expanded on characterising memes based on overall similarity by incorporating multiple modalities into a deep learning-based representation. They initially trained a classifier that extracts textual and facial features along with visual features to classify memes from non-memetic multimodal content. The authors went on



Figure 2: Example group of related memes as captured by Beskow, Kumar, and Carley (2020).

to use the meme representations trained in their classifier to cluster memes based on similar overall appearances. Using this approach, the authors were able to track the evolution of memes in political conversations. While their classifier may appear to solve the task of identifying memes from non-memes, the definition applied by the authors is based solely on the appearance of the text overlays and does not consider memetics. They defined a meme as an image that is either superimposed by text in the Impact font or has text placed in some whitespace over it. Although many memes are consistent with these criteria, this definition strictly limits memes based on the visual appearance of its textual content and does not consider the memetics of individual elements. In fact, within their own illustrated examples, Beskow, Kumar, and Carley (2020) published several images which they described as memes, clearly showing memetic imitation via the visual modality, which do not fit either stated criteria (see Fig. 2).

Dubey et al. (2018) employed sparse matching of meme templates to establish links between memes. To isolate the background of every meme, they removed overlays and superimposed elements from every meme and compared it to a set of popular meme templates. Clustering memes by their backgrounds allowed the authors to identify memes that shared templates, the memetic elements considered by this approach. The authors were also able to identify memes that used templates that had been remixed with visual alterations and overlays. Although template-based memes are often considered the most common form of memes (Knobel and Lankshear 2007), such a template-based approach limits the possible shared memetic elements to the background of a meme, ignoring how memetic transfer can be observed in foreground elements.

Courtois and Frissen (2022) proposed grouping memes based on local visual features as opposed to the overall appearance. This approach recognises visual memetic elements that may be in the template background of one meme but is used as a foreground element in another meme (see example in Table 1(d) where the character from the *Disaster Girl* meme template² is used as a foreground element). They established the relationship between two memes by counting the number of matching features and the difference in position of each matched feature. Much like other works discussed here, this method groups memes by memetic elements. Additionally, this approach acknowledges that visual elements can be replicated as whole-image templates or as individual elements superimposed on another background image, similar to examples discussed in Shifman (2014a).

From these works, we can note two considerations when

²<https://knowyourmeme.com/memes/disaster-girl>

Meme Type	(a) Character Macro	(b) Format Macro	(c) Memetic Images	(d) Transferred Symbols	(e) Memetic Trend
Memetic Element (Location)	<i>Imagination Spongebob</i> (Background)	<i>Drake Template</i> (Background)	<i>Jim Halpert Smiling</i> (Image)	<i>Disaster Girl</i> (Superimposed)	<i>Better love story than Twilight</i> (Text)
Novel Element	Text caption	Text clusters	Text caption	Background image	Background image

Table 1: Example of each meme type from our meme typology. Each labelled with the *memetic element* (and its location) each shares with, and the *novel elements* that distinguish it from, its *related but distinct* memes (not shown here).

applying the *Form* dimension of Shifman’s definition at scale: First, the relatedness of a meme to another can be determined by observing shared visual elements between them. Thus, the memetics of a meme can be established by identifying the memetic elements within it that are shared with other memes. Second, these works illustrate three common types of visual memetic elements: *meme templates* as used by Dubey et al. (2018) in which the background is memetically transferred between memes, *foreground visual features* transferred from other memes (Courtois and Frissen 2022), and *memetic images* appended with text placed in a whitespace (Beskow, Kumar, and Carley 2020).

A limitation of these approaches is that they each rely on a limited pool of memetic elements that was collected from reference memes (Dubey et al. 2018; Zannettou et al. 2018) or within a collected corpus (Courtois and Frissen 2022). Given the dynamic remixing of memetic elements (Qu et al. 2022), a method is needed to recognise memetic elements based on whether they are used contemporarily in the wider online culture, rather than limited to those from a finite set of reference memes.

3 Proposed Meme Identification Model

Building upon the characteristics of memes as we discussed in the preceding sections, we present the following **definition for multimodal Internet Memes** which serves as a basis for our meme identification model.

A multimodal³ user-generated piece of digital content that contains at least one visual or textual component that can be identified as a *memetic element* that is shared with other *related but distinct* memes.

Where we define the following:

Memetic Element: A textual or visual element that is observably reused in multiple distinct memes; e.g. the background shared between multiple memes in Fig. 4(a).

Related but Distinct: Pairs or groups of memes that share common memetic elements, but differ from one another by some *novel element*, such as the distinct text captions in Fig. 4(a) and (b).

³We discuss limitations on modalities in Section 3.3.

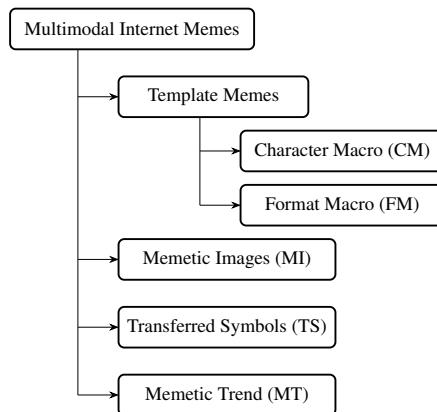


Figure 3: Proposed typology of multimodal Internet memes.

To identify a piece of user-generated content (a “**candidate meme**”) as a meme, one would need to observe at least one memetic element within it. For an element in a candidate meme to be considered memetic, we must search for and find *related but distinct* memes using that same element. To recognise such elements, we present a typology of multimodal memes based on the type of memetic element it holds and a protocol to systematically identify such elements.

3.1 Meme Typology

To assist with identifying memetic elements within candidate memes, we present a meme typology based on which component is found to be memetic. By combining multiple memetic elements, a meme can fit multiple types simultaneously. In the context of verifying whether a sample is a meme, identifying memetics in at least one of these categories would establish that it is a meme. Our typology consists of:

Template Memes encompass memes that use identical visual backgrounds as templates. These backgrounds are reused across related but distinct memes, often with only the textual content being different. Template Memes are commonly found on *meme generation sites*⁴. Among Template Memes,

⁴e.g.: <https://imgflip.com/>

two common categories exist: *Character Macro* and *Format Macro*, see (a) and (b) in Table 1, respectively.

Character Macro (CM) memes are established visual templates centred around characters that uses superimposed text captions to signify some “stereotypical behaviour” (Shifman 2014a). Lists of such characters can be found on meme encyclopaedia sites⁵.

Format Macro (FM) memes use visual templates where the meaning is conveyed through the relative position of visual and textual elements, as described by Hazman, McKeever, and Griffith (2023). The visual template and the intentional positioning of the superimposed text create metaphors (Shang et al. 2021b), which are then reused by editing the textual content (e.g. Table 1(b)).

Memetic Images (MI) describe memes in which an image is appended with a whitespace containing a text caption as described by Beskow, Kumar, and Carley (2020), e.g.: Table 1(c). The memetics of such images are realised by applying multiple different text captions to the same image. Similar to template memes, MI memes are visually very similar to their related but distinct memes e.g.: Fig. 4(c).

Transferred Symbols (TS) memes rely on conventionalised visual elements but do not make use of memetic backgrounds. Instead, they use symbols that have gained meaning in the larger meme culture; transferring both the symbol and its adopted meaning onto a new and otherwise non-memetic background. These transferred symbols are usually observed as either a superimposed visual element e.g. in Table 1(d) or as a visual segment e.g. Fig 4(b).

Memetic Trend (MT) memes are perhaps the most difficult to identify. These memes usually share similar text but may not use visually memetic elements. These are created in response to a participative trend or topic that inspired large amounts of visually dissimilar memes, such as the *Better Love Story than Twilight* memetic trend shown in Table 1(e).

3.2 Meme Identification Protocol

We present our protocol for multimodal meme identification below. First, we define C as a candidate multimodal meme: an image file that contains visual and textual elements that we would like to identify as either a meme or as non-memetic. Further, we define $IS(image)$ as a *reverse image search* where image file, $image$, is submitted as a query to the search engine, which returns N visually similar image files. $TS(text)$ represents a traditional image search function that returns N image files that correspond to, or contains, a given text input, $text$. When reviewing the results of these search functions, the objective is to find at least one meme that is *Related to but Distinct* from C , i.e.: image files that share at least one common memetic element with C while being distinct from each other by some novel element(s). For both IS and TS , we reviewed all returned search results, which proved to be manageably few. Our protocol consists of the following steps:

1. Retrieve images that are visually similar to C , $IS(C)$ and review these for memes that share the same background

⁵e.g.: <https://knowyourmeme.com/>

⁶<https://images.google.com/>

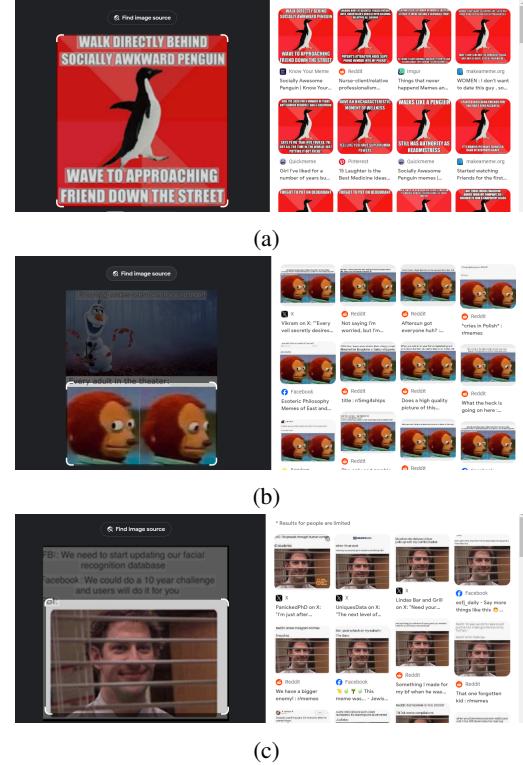


Figure 4: Google Image Search⁶ results for (a) a Character Macro meme showing related memes using the same template, (b) Transferred Symbols meme, where the search input is cropped to the bottom segment, showing related memes using the *Monkey Puppet* memetic element and (c) Memetic Image meme, where the search input is cropped to remove the textual element, showing related memes that reuse the same image.

as C but with different foreground text or visual content. If at least one is found, C is a *Template Meme*. Else: go to Step 2.

2. Consider the overall visual format of C . If C comprises a character in the centre with a text caption overlay (e.g. Fig 4(a)), go to Step 3. If C comprises of multiple segments, each containing a separate visual content (e.g. Fig. 4(b)), go to Step 4. Otherwise, if C comprises an image and a separately appended whitespace (e.g. Fig. 4(c)), go to Step 5. Otherwise, go to Step 6.
3. Crop C to remove textual content, creating c_{crop} . Review $IS(c_{crop})$ for any meme that shares the same background as C but with different foreground text or visual content. If found, C is a *Character Macro Meme*. Otherwise, go to Step 6.
4. Crop C into its visual segments, c_m , where m ranges from 1 to the number of segments in C . Review $IS(c_m)$ for each m , and search for memes that pair segment m with other visual segments or different text captions. If found for any c_m , C is a *Transferred Symbols* meme. If none is found for all c_m , go to Step 6.

5. Crop C to remove the white space to create c_{img} . Review $IS(c_{img})$ for memes that contain c_{img} but have different textual content. If found for any c_{img} , C is a *Memetic Image* meme. Otherwise, go to Step 6
6. Crop each distinct visual elements particularly those that appear to have been superimposed onto a background image, to create c_{elem} . Review $IS(c_{elem})$ for other memes that include c_{elem} . If found for any c_{elem} , C is a *Transferred Symbols* meme. Otherwise, go to Step 7.
7. Extract the textual elements from within C , c_{text} . Review $TS(c_{text})$ for distinct memes that contain textual content that is, wholly or in most parts, identical to that in the candidate meme. If c_{text} is shared across numerous visually dissimilar memes, C is a *Memetic Trend* meme. Otherwise, go to Step 8.
8. If the candidate meme does not contain any identifiable memetic elements, as per the steps above, it is considered non-memetic.

3.3 Addressing Current Challenges in Meme Identification

Basing meme identification on the existence of *related but distinct* memes, as outlined in our protocol, removes the reliance on personal judgement or familiarity with meme culture from the process of recognising memetic elements as was used in some recent datasets (Miliani et al. 2020; Suryawanshi et al. 2020). Due to the extremely large number of visual and textual elements that may possibly be memetic, no annotator could be expected to recognise all such elements from personal knowledge alone. Other alternatives to meme identification have attempted to leverage reference memes sourced from community-maintained encyclopaedias of popular memes (Tommasini, Ilievski, and Wijesiriwardene 2023) or randomly sampled from social media platforms (Sherratt, Pimbblet, and Dethlefs 2023). However, since memetic elements are constantly added to meme parlance (Wiggins and Bowers 2015), it is likely impossible to prepare a set of reference memes that is large enough to encompass the entirety of meme culture.

Arbitrary collections of reference memes, whether formally collected or within the personal familiarity of an annotator, impose artificial constraints on the potential diversity of memetic elements. Both approaches rely on memetic elements achieving some degree of global popularity and are likely to overlook two significant categories of memes: non-English and newly emerging memes. Due to being more culturally specific, memes that are rooted in non-English cultures or use non-English text are less likely to achieve the global popularity required to be included into a meme encyclopaedia or to be represented within a random sampling of a meme-sharing communities which tend to use English as their lingua franca. Furthermore, since a meme only requires one memetic derivative to qualify as a meme, these approaches are likely to overlook newly emerging memes that have yet to achieve significant popularity.

To allow us to recognise personally unfamiliar, culturally specific, as well as emerging memetic elements, our protocol instead employs visual search engines. Specifically, we

use Google Image Search to retrieve visually similar images and/or multimodal memes, which are then assessed for memetic elements shared with the candidate meme. If any of the results share memetic element(s) with the candidate meme but is otherwise distinct, it qualifies as a *related but distinct* meme. Candidate memes for which the search only returns results identical to itself, even from multiple different sources, indicate that the candidate meme is likely a viral piece, rather than a memetic one, as described by Shifman (2014b).

Our protocol caters specifically for multimodal memes as they are represented within contemporary multimodal meme datasets: where each sample comprises an image file that contains both visual and textual content. We acknowledge that this representation does not encompass all forms of memes or multimodality within memes. For example, a *Memetic Image* meme could manifest as a memetically reused image accompanied by novel text captions within the body of a social media post; instead of textual content within the image file. Despite differences in format, it combines the novel and memetic elements in the same way that a Memetic Image meme (e.g. Fig. 5(c)) does and is both multimodal and memetic. However, since current meme datasets do not capture textual content outside of image files, neither does our protocol.

4 Analysis of Current Meme Datasets

To assess the current practice of meme identification, we examine how current meme datasets handle the task of meme identification from two perspectives. Each dataset reportedly consists exclusively of memes, each having been manually selected and annotated with affective classification label(s). First, we present a survey of current approaches to meme identification taken by the authors of current meme classification datasets. Here, we attempt to identify (1) whether current datasets regard memes as distinct from non-memes and (2) the criteria and methods used to identify the memes selected for each dataset. Specifically, we examine publications that present memes as samples in a dataset, each annotated with one or more affective classification label(s). We excluded publications presenting datasets that are re-annotations or extensions of another dataset or those that lack scientific peer review. Second, we apply our protocol to samples taken from the datasets. This allows us to determine whether current meme datasets include both memes and non-memes.

4.1 Meme Identification Approaches in Existing Meme Classification Datasets

We reviewed the publications that accompany 12 meme datasets. The purpose of examining these datasets is to establish how memes are currently being defined in the field of meme classification. We sought information on any steps taken to either identify memes or to remove non-memetic content during data collection. Then, we extracted any criteria reportedly used during data collection, focussing on whether each criterion helps distinguish memetics within memes. Table 2 summarises our findings on: (1) the presence of a meme identification step, and (2) any criteria used to qualify memes

Dataset	Meme Identification Reported?: Approach	Reported Filtration Criteria
Miliani et al. (2020)	Yes: Annotated by labelers based on: “formal aspects (layout, multimodality and manipulation) as well as content, e.g. ironic intent”.	– Relevance to selected topic. – Italian-only text.
Sharma et al. (2020)	Yes: Manual selection by authors.	– Contains “clear background picture”, and “embedded textual content”. – English-only text
Patwa et al. (2022)	Yes: Manual selection by authors and “extensive cleaning” of scraped content by authors.	– Collected from “public domains”
Kiela et al. (2020)	Yes: Manual selection by authors.	– The original image in the meme must be replacable while preserving the meme’s meaning.
Pramanick et al. (2021a)	Yes: Keyword search appended with “memes” followed by manual selection by authors.	– English-only text. – Text is readable. – Contains both text and visual content. – Does not contain <i>cartoon</i> visuals.
Dimitrov et al. (2021)	Yes: Manual selection by authors.	– Does not contain “diagrams/graphs/tables” or <i>cartoon</i> visuals – Contains both text and visual content.
Suryawanshi et al. (2020)	Yes: Volunteers were requested to submit memes specifically.	– Includes Tamil-only text.
Prakash, Hee, and Lee (2023)	Yes: Keyword search appended with “memes”, followed by filtration by annotators.	– Minimum image resolution. – Maximum text length. – Contains text.
Fersini et al. (2022)	Yes: Manual selection by authors including from meme creation sites.	– N/A
Mishra et al. (2023)	Yes: Manual selection by authors.	– Text is in code-mixed Hindi-English.
Hossain, Sharif, and Hoque (2022a)	Yes: Content search appended with keyword “memes” followed by manual selection by authors.	– Bengali-only text. – Text is readable. – Contains both text and visual content. – Does not contain <i>cartoon</i> visuals.
Hossain, Sharif, and Hoque (2022b)	Yes: Manual selection by authors from content search appended with keyword “memes” and “popular meme pages”.	– Text is in code-mixed Bengali-English. – Contains both text and visual content. – Text and visual content are clear. – Does not contain <i>cartoon</i> visuals.

Table 2: Meme Identification Approaches and Reported Filtration Criteria used by current meme classification datasets. For the classification task and the size of each dataset, see Appendix.

for inclusion in each dataset.

Memes as Identifiable Content As shown in Table 2, all 12 dataset publications treat memes as a distinct type of content and have reported some additional effort to filter memes from a large collection of image files. All reviewed works opted for manual approaches to meme identification. Most frequently, the authors of each dataset reported manually selecting memes from image search results, web scraping, and/or browsing online social groups (Sharma et al. 2020; Patwa et al. 2022; Kiela et al. 2020; Pramanick et al. 2021a,b; Fersini et al. 2022). However, none of these works reported criteria for how memes were chosen during this manual selection process. Patwa et al. (2022) also reportedly performed an “extensive cleaning” step following their data collection via web-crawling but neither the purpose nor the criteria used

for this step were reported.

Furthermore, Pramanick et al. (2021a,b); Hossain, Sharif, and Hoque (2022a); Prakash, Hee, and Lee (2023) used search functions to assist in filtering out non-memes by adding the keyword “meme” to their image search input. This delegation of the meme identification task to search engines is perhaps best demonstrated by Sabat, Ferrer, and Giro i Nieto (2019), who instead of publishing the data used to train their classifiers, uploaded code that sends “meme”-appended queries⁷ to Google Image Search and downloads the returned results.

Three publications delegated meme identification to annotators: Prakash, Hee, and Lee (2023) asked their annotators if they considered each sample to be a meme. Miliani et al.

⁷Queries used: *racist meme*, *jew meme*, and *muslim meme*.

(2020) included Meme versus notMeme as a classification label, where the authors reported that these labels were applied based on “formal aspects (such as layout, multimodality and manipulation) as well as content (e.g. ironic intent)”⁸. Suryawanshi et al. (2020) prompted public participants to submit memes. However, none of these three works published any annotator guidelines on what constitutes a meme nor any protocol to recognise indications of memetics.

Filtration Criteria The rightmost column in Table 2 shows the filtration criteria that were applied during the creation of each dataset. Although each criterion characterises the samples included in each meme dataset, none of the criteria directly relate to memetics nor do they appear to outline what constitutes a meme. In lieu of memetics, we expected to identify the authors of each dataset to determine whether a sample constitutes a meme. Instead, none of the criteria reported reveal how the authors of these datasets identified memes from non-memes. Furthermore, the criteria reported varied significantly between the datasets, and do not collectively point towards a common set of characteristics used to identify memes.

Most frequently, dataset authors restricted their collection according to multimodality by excluding samples that included only image or text content (Miliani et al. 2020; Sharma et al. 2020; Patwa et al. 2022; Pramanick et al. 2021a,b; Hossain, Sharif, and Hoque 2022b,a). Across all these datasets, multimodality is represented as text that appears alongside visual elements within a single image file. Although most of these works assert that multimodality is a common characteristic of Internet memes (and restricted their datasets to multimodal samples out of consistency), Dimitrov et al. (2021); Pramanick et al. (2021b) go further by incorporating multimodality into their definitions of memes; both defining memes specifically as an image with appended/superimposed text. As we have noted previously, multimodality alone, without considering memetics, does not make a meme (Shifman 2013). While other works analyse unimodal text-only memes (He et al. 2019) and image-only memes (Shifman 2014a), our analysis is entirely focused on multimodal meme data sets.

Furthermore, several data sets also restricted their data collection based on language (Miliani et al. 2020; Sharma et al. 2020; Kiela et al. 2020; Pramanick et al. 2021a; Hossain, Sharif, and Hoque 2022a,b; Mishra et al. 2023). However, this did not appear to be related to how the authors identified memes from non-memetic content. Some chose to exclude samples that contained “cartoon” visual elements (Pramanick et al. 2021a; Dimitrov et al. 2021; Hossain, Sharif, and Hoque 2022a,b). This decision seems to originate from the assumption that machine learning models cannot effectively handle cartoon visual elements (Pramanick et al. 2021a).

Summary The authors of current meme datasets refer to memes as a distinct and identifiable form of digital content, incorporating additional effort to filter memes from non-memetic content during data collection. However, none reported a replicable and theoretically grounded methodology to do so. Although these works report various filtration

criteria used during data collection, none seem to aid in distinguishing memes from non-memetic content. In fact, these works do not appear to consider memetics during meme collection. Furthermore, none of the works discussed the role of memetics during the labelling process or how annotators were asked to handle the interpretation of unfamiliar memetic elements. Current approaches rely heavily on manual labour and subjective personal judgments. Therefore, the current approaches and criteria used to create current meme datasets do not offer any alternatives to our meme identification model and identification protocol. As long as data collection methods for creating meme datasets do not apply a clear standard for meme identification, the performance of models built on such datasets may reflect the state of the art in multimodal classification, but not necessarily that of meme classification.

4.2 Identifying Non-Memes in Current Meme Datasets using Our Protocol

This section assesses whether the lack of criteria for identifying memes (as highlighted in the previous section) leads to the existence of non-memes in current datasets. To do this, we manually apply our protocol to samples from each of the meme datasets that include only English text. Since our protocol requires recognising *Memetic Trend* memes via textual elements of memes, we excluded datasets that were specifically filtered for memes with text in languages other than English (Miliani et al. 2020), (Suryawanshi et al. 2020), (Hossain, Sharif, and Hoque 2022a), (Hossain, Sharif, and Hoque 2022b), (Mishra et al. 2023).

Since we performed this evaluation manually, it was only viable for us to take a small sample from each dataset. We note that for some of the datasets, this sample size represents a smaller portion than for others (see Appendix for a list each dataset size). We randomly sampled 200 memes from each dataset using a Python script (see Appendix). For datasets that included splits or multiple subsets, each subset was sampled with the same probability. Our sampling did not take into account any classification labels or contents of the image files as all the samples are equally asserted to be memes by their respective authors regardless of label. As current datasets report to already exclude non-memes, one would expect to only find memes in such a small sampling. Thus, we argue that a substantial presence of non-memes in our sample is indicative of a wider presence of non-memes in the full dataset.

Applying our protocol to each sample chosen, we identify it as one of the following:

- non-memetic image-text multimodal content (as **nMIT** in Table 3);
- non-multimodal text-only or visual-only (as **nMM** in Table 3); or
- for samples that are identified as memes, we record its type per our meme typology.

Results Table 3 presents the proportions of samples from each dataset identified either as a meme or as a non-meme (as nMIT or nMM). On average, 50.4% of the content sampled from each dataset was found not to contain identifiable

⁸No further details were provided.

Dataset	Memes - by type (% of total sample)						Non-Memes (% of total sample)		
	CM	FM	MI	TS	MT	Total	nMIT	nMM	Total
Sharma et al. (2020)	33	7	7	6	2	55	39	6	45
Patwa et al. (2022)	27	19	20	9	2	77	21	2	23
Kiela et al. (2020)	5	0	7	3	1	16	80	4	84
Pramanick et al. (2021a)	14	6	18	6	3	48	48	4	52
Dimitrov et al. (2021)	15	7	9	7	3	41	58	2	60
Prakash, Hee, and Lee (2023)	10	16	10	9	4	49	33	18	51
Fersini et al. (2022)	20	12	18	7	3	62	38	0	38
Average						49.6			50.4

Table 3: Proportion of samples (200 per dataset) found to contain memetic elements by type: CM: Character Macro memes, FM: Format Macro memes, MI: Memetic Images, TS: Transferred Symbols memes, MT: Memetic trend memes; and proportion of samples not containing any memetic elements: nMIT: Non-Memetic Image-Text multimodal content, nMM: non-memetic non-multimodal content. For each dataset’s classification task, size and distribution licence, see Appendix.

memetic elements. This proportion varies widely between the different datasets, with a minimum non-meme proportion of 23%. These findings show that all the evaluated datasets contain a sizable number of samples without memetic elements and therefore are not memes. Combined with the lack of reported criteria used in relation to meme identification, it appears that current meme datasets do not systematically distinguish memes from non-memes. Further, it is unclear whether memetics was considered at all during the creation of these meme datasets, and if so, what definition of memetics was used and applied.

Discussion Consider the sample found in the (Sharma et al. 2020) dataset: Fig. 5(a). Although this sample fulfils all the criteria stated by the authors of that dataset viz. contains a “clear background image” and “embedded textual content”, which is “wholly in English”, it does not contain any memetic elements. When evaluated through our protocol, the reverse image search returned numerous identical copies, indicating that it had been widely propagated across the Internet in its original form. As we discussed previously, this fits into Shifman’s Shifman (2014b) definition for a “viral”, but does not contain any elements that have been reused memetically prior or since. Thus, understanding the meaning conveyed by this sample requires some cultural context but does not require any interpretation of memetic elements. In contrast, the same three people are included in Fig. 5(b), a Template Meme for which we found only two related but distinct memes, including the one shown in Fig. 5(c). The existence of these memes is proof that the background is a memetic element. Furthermore, this exemplifies how our protocol can identify memetic elements that have not achieved significant popularity.

Other than the inclusion of memetic elements, we noted that memetic and non-memetic samples appear superficially very similar. Like many multimodal memes, the sample shown in Fig. 5(a) utilise visual segments to create comic-like frames, similar to memes shown in Fig. 4(b) and Table 1(b). Also, many non-memetic samples use visual formats used in Character Macro and Memetic Images memes, but with images that have not inspired derivative memes. These similarities highlight the challenge of identifying memes from non-memetic content without systematic steps to verify the



Figure 5: (a) Non-meme sampled from (Sharma et al. 2020); (b) & (c) Example of related but distinct meme pair that share common albeit unpopular visual background.

memetic nature of the contents within a meme, and referring to large collections of memetic elements.

Further, we also observed that all samples, memetic and otherwise, almost universally include cultural references. As such, cultural contextual information is needed to classify the affect of both types of content, such as that proposed by Pramanick et al. (2021b). However, within non-memetic content, these references can be contextualised by recognising what Shang et al. (2021a) describes as “commonsense knowledge” that underlie the symbolic meanings of elements within multimodal content. Crucially, *commonsense* here describes facts that an average person is expected to know (Shang et al. 2021a). In contrast, the memetic process itself serves as the cultural context for memetic elements (Cannizzaro 2016), necessitating familiarity with the meme culture to be able to interpret such elements (Wiggins and Bowers 2015). Since memes combine references to both ongoing cultural events and meme culture in particular, interpreting memes requires context and knowledge from beyond *commonsense* alone, requiring knowledge of the memetic elements within the meme (Dubey et al. 2018).

Furthermore, consider the sample shown in Fig. 6(a), sampled from the Kiela et al. (2020) dataset. This example illustrates the nuanced semantics conveyed by memetic elements.

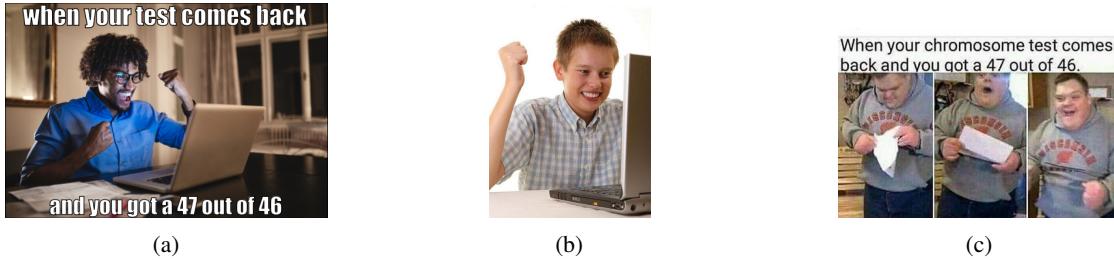


Figure 6: (a) Synthetically recreated meme (©Getty Images) sampled from (Kiela et al. 2020); (b) *First Day on the Internet Kid* Character Macro template; and (c) Example meme from the derogatory “47 out of 46” Memetic Trend.

Citing concerns over copyright infringement, the authors of this dataset tasked annotators with recreating memes using “near-identical” visual content; often removing visual memetic elements in the process. This sample depicts a person excited to receive a result of “47 out of 46” for some undefined test. It is ostensibly joyful and is labelled as non-hateful in the dataset. However, we believe that two memetic elements were removed during the recreation process, resulting in a loss of meaning of the original meme. First, we believe that the original meme used the Character Macro template *First Day on the Internet Kid*, see Fig. 6b. Unlike the background of the recreated sample, this character has come to symbolise naive excitement for a false, deceitful, or otherwise unexciting occurrence online. This suggests a more subtle (and possibly negative) meaning of the test results mentioned here. Second, applying Step 7 of our proposed protocol reveals that the original meme was likely part of a Memetic Trend surrounding memes with captions including “47 out of 46”. These numbers refer to chromosomes and, as part of a memetic trend, were often appended derogatorily to images depicting persons with Trisomy or Down Syndrome; see example in Fig. 6c. Without recognising these memetic elements, which were possibly obfuscated by the recreation process, this meme takes on a harmless and joyful appearance. Only when viewed through the relevant memetic contexts can we clearly recognise the derogatory “humour” being conveyed.

Since current datasets contain both memetic and non-memetic content, they may be well suited for training supervised learning models to classify multimodal content in general, but not specifically memes. The lack of clear delineation in terms of memetics or any reported criteria used to differentiate memes from non-memes in these datasets begs the question: what qualifies these datasets as “meme datasets”? If our protocol had been applied during the creation of these datasets, each sample would be systematically and verifiably identified as a meme. This, in turn, would allow the distinction between meme classification and multimodal content classification to be clearly delineated; where meme classification specifically requires the interpretation of memetic elements.

5 Limitations & Future Works

In line with current practices in preparing meme datasets, our protocol is designed to be performed manually. However, this

severely limited the number of samples we could evaluate per dataset. Although this allowed us to observe that current datasets do not exclude non-memetic content, the proportions of memes to non-memes we reported in Table 3 may not hold across the entirety of each dataset, and we caution against generalising our observations as such.

Although this work does not directly propose an automated meme identification solution, our protocol illustrates three core capabilities required by any such solutions. Firstly, the solution must recognise memetic elements both globally and locally (Courtois and Frissen 2022). Secondly, it should access a vast repository of digital artefacts to recognise memetic elements within emerging and culturally specific memes. To achieve this, Image Search APIs, such as Google Image Search, could be used instead of our manual interaction with Google Image Search. Lastly, an automated solution cannot rely on matching candidate memes to other memes on similarity alone, but must also be able to recognise both novelty between related-but-distinct memes. Additionally, to build machine learning-based solution, a dataset of systematically labelled memes and non-memes is needed. Our protocol provides a repeatable and verifiable approach to collecting such data.

We limited our analysis to meme datasets that are reported to contain only monolingual English memes. However, since our protocol is based on language-agnostic characteristics shared by all multimodal media (e.g., shared visual features, differences in textual content), we posit that our approach can be applied to identify non-English or multilingual memes given annotators who are fluent in the relevant language(s). Additionally, we foresee two linguistic relationships between candidate memes that are absent from a monolingual meme corpus: (1) translation pairs, where two visually identical memes contain semantically identical textual content in different languages, and (2) transliteration pairs, where two visually identical memes contain semantically and linguistically identical textual content in different scripts. Consider the translation pair in Fig. 7 (the English version was sampled from the Hossain, Sharif, and Hoque (2022a) dataset). One would need to determine whether semantically identical memes in different languages qualify as distinct from each other or not.

An issue that stems from current datasets containing both memetic and non-memetic content is that the role memetic elements play in how memes convey meaning remains un-

When the light turned green .00001 seconds ago and someone beeps the horn behind you..



Cuando hace 0,00001 segundos que el semáforo se ha puesto en verde y alguien te pita atrás.



Figure 7: Example Memetic Image meme *Translation Pair* with textual content in English and Spanish, respectively.

recognised. Since these elements gain meaning through the memetic process, dedicated knowledge mining and representation approaches are needed to incorporate knowledge of memetic elements into meme classifiers. Recent works that extract and represent large collections of memetic elements (Tommasini, Ilievski, and Wijesiriwardene 2023; Sherratt 2022) suggest that knowledge of these elements would likely not be gleaned from a finite dataset of sample memes.

Throughout this paper, we argue that meme classification is a distinct and more complex task than the classification of non-memes. However, we could not demonstrate this difference because the existing meme datasets are not composed exclusively of memes (see Table 3). As such, we view meme identification as a key step in creating meme classification datasets and advancing meme affective classification approaches. Using our proposed protocol, we plan to develop a dataset that consists only of memes. With such a dataset, we aim to quantitatively assess the difference in performance of state-of-the-art approaches to meme classification versus classification of multimodal content in general.

We also plan to leverage the inherent nature of memes to form genres (Wiggins and Bowers 2015) to collect entire groups of visually similar but distinct memes during the *reverse image search* step within our protocol. Within such groups, our aim is to collect memes that serve as contrastive examples of one another: memes that are very similar but contain small differences that allow each to express a different affect. This enables the training of meme classifiers using a supervised contrastive learning approach (Pinitas et al. 2022) to capture how the interaction between memetic and novel elements create memes that convey different affects.

6 Conclusion

In this work, we provided a definition for multimodal Internet Memes based on well-established perspectives on memetics from both philological and computational research. We proposed a meme identification model to distinguish memes from non-memetic content. As part of this model, we developed a typology of memes to identify which component of a meme is memetic. We extended this identification model into a protocol that allows meme dataset authors to identify memes in a repeatable and verifiable manner, without relying on personal judgement or familiarity with meme culture. While current meme datasets incorporate some additional effort to separate memes from non-memes, none reported any methodology or criteria used for meme identification.

Applying our protocol to 7 meme classification datasets revealed that current meme datasets include both memetic and non-memetic content. Finally, we argued for the creation of meme datasets that exclusively contain memes as a key step towards much needed advancements in meme classification tasks.

References

- Beskow, D. M.; Kumar, S.; and Carley, K. M. 2020. The evolution of political memes: Detecting and characterizing internet memes with multi-modal deep learning. *Information Processing & Management*, 57(2): 102170.
- Cannizzaro, S. 2016. Internet memes as internet signs: A semiotic view of digital culture. *Sign Systems Studies*, 44: 562.
- Courtois, C.; and Frissen, T. 2022. Computer Vision and Internet Meme Genealogy: An Evaluation of Image Feature Matching as a Technique for Pattern Detection. *Communication Methods and Measures*, 17: 17 – 39.
- Dimitrov, D.; Bin Ali, B.; Shaar, S.; Alam, F.; Silvestri, F.; Firooz, H.; Nakov, P.; and Da San Martino, G. 2021. Detecting Propaganda Techniques in Memes. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 6603–6617. Online: Association for Computational Linguistics.
- Dubey, A.; Moro, E.; Cebrian, M.; and Rahwan, I. 2018. Meme-sequencer: Sparse matching for embedding image macros. In *Proceedings of the 2018 world wide web conference*, 1225–1235.
- Fersini, E.; Gasparini, F.; Rizzi, G.; Saibene, A.; Chulvi, B.; Rosso, P.; Lees, A.; and Sorensen, J. 2022. SemEval-2022 Task 5: Multimedia Automatic Misogyny Identification. In Emerson, G.; Schluter, N.; Stanovsky, G.; Kumar, R.; Palmer, A.; Schneider, N.; Singh, S.; and Ratan, S., eds., *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, 533–549. Seattle, United States: Association for Computational Linguistics.
- Gal, N.; Shifman, L.; and Kampf, Z. 2016. “It Gets Better”: Internet memes and the construction of collective identity. *New Media & Society*, 18(8): 1698–1714.
- Hazman, M.; McKeever, S.; and Griffith, J. 2023. Unimodal Intermediate Training for Multimodal Meme Sentiment Classification. In Mitkov, R.; and Angelova, G., eds., *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, 494–506. Varna, Bulgaria: INCOMA Ltd., Shoumen, Bulgaria.
- He, S.; Yang, H.; Zheng, X.; Wang, B.; Zhou, Y.; Xiong, Y.; and Zeng, D. 2019. Massive Meme Identification and Popularity Analysis in Geopolitics. In *2019 IEEE International Conference on Intelligence and Security Informatics (ISI)*, 116–121.
- Hossain, E.; Sharif, O.; and Hoque, M. M. 2022a. MemoSen: A Multimodal Dataset for Sentiment Analysis of Memes. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, 1542–1554. Marseille, France: European Language Resources Association.
- Hossain, E.; Sharif, O.; and Hoque, M. M. 2022b. MUTE: A Multimodal Dataset for Detecting Hateful Memes. In Hanqi, Y.; Zonghan, Y.; Ruder, S.; and Xiaojun, W., eds., *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing: Student Research Workshop*, 32–39. Online: Association for Computational Linguistics.

- Kiela, D.; Firooz, H.; Mohan, A.; Goswami, V.; Singh, A.; Ringshia, P.; et al. 2020. The Hateful Memes Challenge: Detecting hate speech in multimodal memes. *Advances in Neural Information Processing Systems*, 33: 2611–2624.
- Knobel, M.; and Lankshear, C. 2007. *Online memes, affinities, and cultural production*. New literacies and digital epistemologies. New York, NY, US: Peter Lang Publishing. ISBN 978-0-8204-9523-1.
- Miliani, M.; Giorgi, G.; Rama, I.; Anselmi, G.; and Lebani, G. E. 2020. DANKMEMES @ EVALITA 2020: The Memeing of Life: Memes, Multimodality and Politics. In *International Workshop on Evaluation of Natural Language and Speech Tools for Italian*.
- Mishra, S.; Suryavardan, S.; Patwa, P.; Chakraborty, M.; Rani, A.; Reganti, A.; Chadha, A.; Das, A.; Sheth, A. P.; Chinnakotla, M.; Ekbal, A.; and Kumar, S. 2023. Memotion 3: Dataset on sentiment and emotion analysis of codemixed Hindi-English Memes. In Das, A.; Sheth, A. P.; and Ekbal, A., eds., *Proceedings of De-Factify 2: 2nd Workshop on Multimodal Fact Checking and Hate Speech Detection, co-located with AAAI 2023, Washington DC, USA, February 14, 2023*, volume 3555 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Nissenbaum, A.; and Shifman, L. 2017. Internet memes as contested cultural capital: The case of 4chan's /b/ board. *New Media & Society*, 19(4): 483–501.
- Patwa, P.; Ramamoorthy, S.; Gunti, N.; Mishra, S.; Suryavardan, S.; Reganti, A.; et al. 2022. Findings of Memotion 2: Sentiment and Emotion Analysis of Memes. In *De-Factify @ AAAI 2022. First Workshop on Multimodal Fact-Checking and Hate Speech Detection*, CEUR Workshop Proceedings. AAAI.
- Pearce, K.; and Hajizada, A. 2014. No laughing matter humor as a means of dissent in the digital Era: The case of Authoritarian Azerbaijan. *Demokratizatsiya*, 22: 67–85.
- Pinitas, K.; Makantasis, K.; Liapis, A.; and Yannakakis, G. N. 2022. Supervised Contrastive Learning for Affect Modelling. In *Proceedings of the 2022 International Conference on Multimodal Interaction*, ICMI '22, 531–539. New York, NY, USA: Association for Computing Machinery. ISBN 9781450393904.
- Prakash, N.; Hee, M. S.; and Lee, R. K.-W. 2023. TotalDefMeme: A Multi-Attribute Meme Dataset on Total Defence in Singapore. In *Proceedings of the 14th Conference on ACM Multimedia Systems*, MMSys '23, 369–375. New York, NY, USA: Association for Computing Machinery. ISBN 9798400701481.
- Pramanick, S.; Dimitrov, D.; Mukherjee, R.; Sharma, S.; Akhtar, M. S.; Nakov, P.; and Chakraborty, T. 2021a. Detecting Harmful Memes and Their Targets. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2783–2796. Online: Association for Computational Linguistics.
- Pramanick, S.; Sharma, S.; Dimitrov, D.; Akhtar, M. S.; Nakov, P.; and Chakraborty, T. 2021b. MOMENTA: A Multimodal Framework for Detecting Harmful Memes and Their Targets. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 4439–4455. Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Qu, Y. Q.; He, X.; Pierson, S.; Backes, M.; Zhang, Y.; and Zannettou, S. 2022. On the Evolution of (Hateful) Memes by Means of Multimodal Contrastive Learning. *2023 IEEE Symposium on Security and Privacy (SP)*, 293–310.
- Sabat, B. O.; Ferrer, C. C.; and Giro i Nieto, X. 2019. Hate Speech in Pixels: Detection of Offensive Memes towards Automatic Moderation. In *NeurIPS 2019 Workshop on AI for Social Good*. Vancouver, Canada.
- Segev, E.; Nissenbaum, A.; Stolero, N.; and Shifman, L. 2015. Families and Networks of Internet Memes: The Relationship Between Cohesiveness, Uniqueness, and Quiddity Concreteness. *Journal of Computer-Mediated Communication*, 20(4): 417–433.
- Shang, L.; Youn, C.; Zha, Y.; Zhang, Y.; and Wang, D. 2021a. KnowMeme: A Knowledge-enriched Graph Neural Network Solution to Offensive Meme Detection. In *2021 IEEE 17th International Conference on eScience (eScience)*, 186–195.
- Shang, L.; Zhang, Y.; Zha, Y.; Chen, Y.; Youn, C.; and Wang, D. 2021b. AOMD: An Analogy-Aware Approach to Offensive Meme Detection on Social Media. *Inf. Process. Manage.*, 58(5).
- Sharma, C.; Bhageria, D.; Scott, W.; PYKL, S.; Das, A.; Chakraborty, T.; et al. 2020. SemEval-2020 Task 8: Memotion Analysis- the Visuo-Lingual Metaphor! In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, 759–773. Barcelona (online): International Committee for Computational Linguistics.
- Sharma, S.; Alam, F.; Akhtar, M. S.; Dimitrov, D.; Da San Martino, G.; Firooz, H.; Halevy, A.; Silvestri, F.; Nakov, P.; and Chakraborty, T. 2022. Detecting and Understanding Harmful Memes: A Survey. In Raedt, L. D., ed., *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, 5597–5606. International Joint Conferences on Artificial Intelligence Organization. Survey Track.
- Sherratt, V. 2022. Towards Contextually Sensitive Analysis of Memes: Meme Genealogy and Knowledge Base. In Raedt, L. D., ed., *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, 5871–5872. International Joint Conferences on Artificial Intelligence Organization. Doctoral Consortium.
- Sherratt, V.; Pimbblet, K.; and Dethlefs, N. 2023. Multi-Channel Convolutional Neural Network for Precise Meme Classification. In *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*, ICMR '23, 190–198. New York, NY, USA: Association for Computing Machinery. ISBN 9798400701788.
- Shifman, L. 2013. Memes in a Digital World: Reconciling with a Conceptual Troublemaker. *Journal of Computer-Mediated Communication*, 18(3): 362–377.
- Shifman, L. 2014a. The Cultural Logic of Photo-Based Meme Genres. *Journal of Visual Culture*, 13(3): 340–358.
- Shifman, L. 2014b. *Memes versus Virals*. The MIT Press. ISBN 9780262525435.
- Suryawanshi, S.; Chakravarthi, B. R.; Verma, P.; Arcan, M.; McCrae, J. P.; and Buitelaar, P. 2020. A Dataset for Troll Classification of TamilMemes. In *Proceedings of the WILDRE5– 5th Workshop on Indian Language Data: Resources and Evaluation*, 7–13. Marseille, France: European Language Resources Association (ELRA). ISBN 979-10-95546-67-2.
- Tommasini, R.; Ilievski, F.; and Wijesiriwardene, T. 2023. IMKG: The Internet Meme Knowledge Graph. In Pesquita, C.; Jimenez-Ruiz, E.; McCusker, J.; Faria, D.; Dragoni, M.; Dimou, A.; Troncy, R.; and Hertling, S., eds., *The Semantic Web*, 354–371. Cham: Springer Nature Switzerland. ISBN 978-3-031-33459-9.
- Wiggins, B. 2019. *The Discursive Power of Memes in Digital Culture: Ideology, Semiotics, and Intertextuality*. Routledge. ISBN 9780429492303.
- Wiggins, B. E.; and Bowers, G. B. 2015. Memes as genre: A structural analysis of the memescape. *New Media & Society*, 17(11): 1886–1906.
- Zannettou, S.; Caulfield, T.; Blackburn, J.; De Cristofaro, E.; Sirivianos, M.; Stringhini, G.; and Suarez-Tangil, G. 2018. On the Origins of Memes by Means of Fringe Web Communities. In *Proceedings of the Internet Measurement Conference 2018*, IMC '18, 188–202. New York, NY, USA: Association for Computing Machinery. ISBN 9781450356190.

A File Sampling Python Script

The following are Python functions written to sample candidate memes from each f the dataset analysed in Table 3.

```
# Functions to randomly sample candidate
memes from folders containing meme
datasets.
# Assumes each datasets to comprise
either a folder containing image
files or a folder of subfolders of
image files.
from pathlib import Path
from os import mkdir, listdir
from os.path import join, exists, isdir,
basename
import shutil
from typing import List

# Filter non-image files from list of
filepaths
def filter_images(files: List[Path]):
    from PIL import Image
    for f in files:
        try:
            Image.open(f).verify()
        except Exception:
            files.remove(f)
    return files

# Randomly select N number of files
from list of filepaths
def random_file_selection(files: List[
    Path], n: int):
    import random
    return [f for f in random.choices(
        files, k=n)]

# Randomly select N samples from a
dataset
def sample_dataset(input_dir: Path,
output_dir: Path, n_samples: int =
200):
    if not exists(output_dir):
        mkdir(output_dir)
    dires = listdir(input_dir)
    # Checks if dataset is split into
    multiple subfolders.
    if all([isdir(join(input_dir, dire)) for
dire in dires]):
        # Collates files from all
        subfolders
        all_files = []
        for dire in dires:
            for file in listdir(join(
                input_dir, dire)):
                all_files.append(join(
                    input_dir, dire, file))
    # Filter for image files only
    img_files = filter_images(
        all_files)
    # Sample with equal probability
    # regardless of file source
```

```
samples = random_file_selection(
    img_files, n_samples)
for f in samples:
    shutil.move(f, join(
        output_dir, basename(f)))
else:
    assert all([not isdir(dire) for
        dire in dires])
    files = [join(input_dir, f) for f
        in dires]
    # Filter for image files only
    img_files = filter_images(files)
    # Sample N files
    samples = random_file_selection(
        img_files, n_samples)
    for f in samples:
        shutil.move(f, join(
            output_dir, basename(f)))
```

B License Types of Evaluated Datasets

The following lists the publications and license for each of the dataset used in our Evaluation section.

Dataset	License Type
Kiela et al. (2020)	Proprietary by Facebook, Inc.
Fersini et al. (2022)	CC BY-NC-SA 4.0
Pramanick et al. (2021a)	BSD
Prakash, Hee, and Lee (2023)	Unspecified but mentions “non-commercial research purposes”
Dimitrov et al. (2021)	BSD
Sharma et al. (2020)	CC BY 4.0
Patwa et al. (2022)	CC BY 4.0

Table 4: Distribution license of each Meme Classification Dataset that we sampled for evaluation.

C Additional Dataset Details

The following lists the Classification Labels, Size and Meme Sources for each of the Dataset assessed in this work.

Dataset	Classification Task	# Memes
Miliani et al. (2020)	<ul style="list-style-type: none"> - Meme vs NotMeme - HateSpeech - Event Clustering 	2,361
Sharma et al. (2020)	<ul style="list-style-type: none"> - Sentiment - Sarcasm - Humour - Offensiveness - Motivational 	10,000
Patwa et al. (2022)	<ul style="list-style-type: none"> - Sentiment - Sarcasm - Humour - Offensiveness - Motivational 	10,000
Kiela et al. (2020)	Hate speech	10,000
Pramanick et al. (2021a)	<ul style="list-style-type: none"> - Harmfulness - Harm Target 	3,544
Dimitrov et al. (2021)	<ul style="list-style-type: none"> - Propaganda Technique 	2,488
Suryawanshi et al. (2020)	- Trolling	2,669
Prakash, Hee, and Lee (2023)	<ul style="list-style-type: none"> - Defence Topic - Stance 	5,301
Fersini et al. (2022)	<ul style="list-style-type: none"> - Misogyny - Misogyny type 	5,500
Mishra et al. (2023)	<ul style="list-style-type: none"> - Sentiment - Sarcasm - Humour - Offensiveness - Motivational 	10,000
Hossain, Sharif, and Hoque (2022a)	- Sentiment	4,417
Hossain, Sharif, and Hoque (2022b)	- Hate speech	4,158

Table 5: Classification tasks and size of the Meme Classification Datasets Evaluated.