

Do AI Systems Understand What We Mean? A Survey of Multimodal Human Creative Expression Understanding

Anonymous ACL submission

Abstract

Multimodal human creativity, such as memes, cartoons, comics, advertisements, and humorous or satirical videos, presents unique challenges for AI systems due to its non-literal, culturally grounded, and rhetorically structured nature. While Multimodal Large Language Models (MLLMs) excel at physical-world understanding, they continue to struggle with creative communication. This survey provides a systematic overview of multimodal human creativity research. We first introduce a unified taxonomy that characterizes creative content along data forms, meaning-making mechanisms, and communicative goals. We then organize existing work into a capability-oriented hierarchy, spanning recognition, interpretation, and generation. Within this framework, we audit the architectural shifts that have shaped the transition from task-specific models to MLLMs. We conclude by identifying critical benchmark trends and open socio-technical challenges, providing a roadmap toward AI systems capable of sophisticated creative understanding.

1 Introduction

Recent advances in Multimodal Large Language Models (MLLMs) have demonstrated strong capabilities in physical-world grounding, such as scene understanding, object recognition, and factual visual question answering (Yin et al., 2024; Comanici et al., 2025; Bai et al., 2025). These successes have motivated expectations that MLLMs can move beyond literal perception to comprehend more abstract and socially grounded forms of human communication. In particular, human creative expression, which dominates contemporary online communication, appears to be a natural next frontier for multimodal AI. However, growing evidence shows that even state-of-the-art MLLMs struggle to interpret such content reliably (Jia et al., 2024; Saakyan et al., 2024; Joshi, 2025).

We define **multimodal human creative expression** as the strategic synthesis of heterogeneous signals (e.g., linguistic, visual, auditory, and temporal cues) to convey meaning that is intentionally *non-literal* and *context-dependent* (Boden, 1998; Wiggins, 2006; Colton and Wiggins, 2012). Common examples include comics, memes, editorial cartoons, and satirical videos. Unlike factual communication, these artifacts use rhetorical mechanisms like metaphor, irony, and incongruity to create a “semantic gap” between the literal depiction and the intended message.

Decoding these artifacts is challenging for MLLMs because their intended meaning is rarely recoverable from surface-level perception alone. While a model may accurately recognize objects, actions, or scenes at the perceptual level, it often fails to infer their communicative function. For instance, an MLLM may accurately identify a “burning house” in an image (physical perception), yet fail to resolve its function as a satirical metaphor for “climate policy” (creative interpretation). Thus, achieving a deep understanding of such content requires more than multimodal alignment; it requires inferring implicit intent and grounding its interpretation in cultural norms and socio-historical context (Saakyan et al., 2024; Ryan et al., 2025).

Despite the increasing interest, existing surveys on human creativity understanding remain fragmented across data, task, and modeling dimensions. From a task perspective, prior reviews often focus on isolated objectives such as detection, understanding (Sharma et al., 2022b; Farabi et al., 2024), or generation (Amin and Burghardt, 2020; Loakman et al., 2023), without capturing the full capability progression from recognition to interpretive reasoning and creativity synthesis. From a data perspective, surveys that consider multimodal inputs typically concentrate on a single content type, such as memes (Farabi et al., 2024), humor (Ren et al., 2024; Kalloniatis and Adamidis, 2024), or

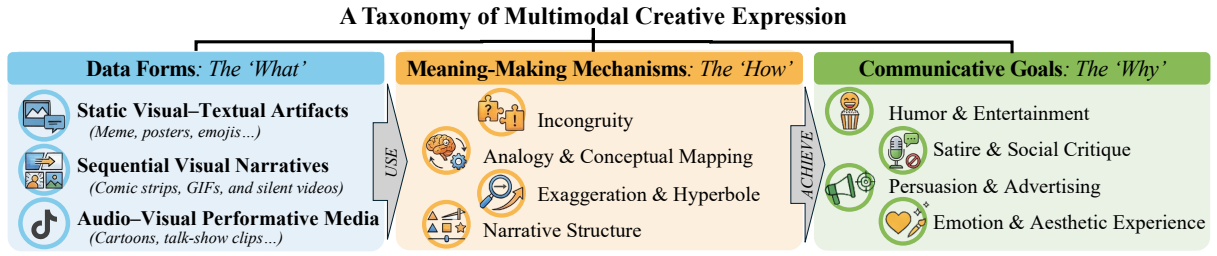


Figure 1: A unified data taxonomy of multimodal human creative expression.

text-only creativity (Amin and Burghardt, 2020; Su et al., 2025), overlooking the diversity of multimodal and **socially-grounded** inputs in real-world settings. From a modeling perspective, existing work (Nguyen and Ng, 2024; Loakman et al., 2025) tends to explore modeling techniques in a piecemeal manner, rather than providing a systematic and diverse discussion of multimodal alignment and reasoning mechanisms.

We argue that progress in multimodal creative AI is constrained not only by individual advances in data or models, but also by the lack of unified representations, capability-aware evaluation, principled synthesis across data, tasks, and models, and ethical grounding across the research pipeline. We address the limitation by systematically examining the problem from the following perspectives:

- **Unified Taxonomy:** We propose a unified taxonomy that organizes multimodal human creativity along *data forms*, *meaning-making mechanisms*, and *communicative goals*, connecting representation, interpretation, and creative intent beyond task- or model-centric views.
- **Capability-Centric Tasks:** We introduce a three-level task hierarchy—*Recognition*, *Interpretation and Reasoning*, and *Creative Generation*—that clarifies the capabilities required across datasets.
- **Capability-Guided Models:** We synthesize modeling paradigms according to the capabilities they support, highlighting how progress emerges from integrating perceptual alignment, reasoning, and grounded knowledge.
- **Benchmarks and Ethical Frontiers:** We bridge our taxonomy to existing benchmarks, and outline open socio-technical challenges to guide future creative AI development.

2 Background: A Taxonomy of Multimodal Creative Expression

To systematically characterize multimodal human creative expression, we propose a three-dimensional taxonomy organized around three fundamental questions: (1) **Data Forms**—*what* me-

dia the content is instantiated in; (2) **Meaning-Making Mechanisms**—*how* non-literal meaning is constructed; and (3) **Communicative Goals**—*why* the content is produced, as shown in Figure 1.

These dimensions capture complementary aspects of creative artifacts: their representational substrate, their underlying rhetorical and cognitive operations, and their pragmatic intent. Together, they provide a principled framework grounded in theories of multimodal communication and rhetoric (Foss, 2017; Lakoff and Johnson, 2024), while aligning naturally with contemporary AI research on multimodal understanding.

2.1 Data Forms: The “What”

Human creative expression spans a wide range of data forms, including text, images, audio, video, and their combinations. While early computational work focused primarily on textual humor and figurative language (Mihalcea and Strapparava, 2005), modern creative communication is overwhelmingly multimodal. Understanding such content requires models to integrate perceptual signals with cultural knowledge and rhetorical intent, going beyond literal semantic interpretation. We group creative data forms into three broad categories:

Static Visual-Textual Artifacts include memes, posters, and emojis. These artifacts typically consist of a single image, often paired with minimal text, to convey humor, stance, or affect. Interpretation depends on visual grounding, symbolic association, and shared socio-cultural context (Shifman, 2013; Sharma et al., 2020).

Sequential Visual Narratives encompass comic strips, GIFs, and silent videos, where meaning emerges from temporal progression across frames. Understanding these forms requires tracking entities and events, inferring causal relations, and integrating information across a visual sequence (Paval et al., 2025; Wang et al., 2025b).

Audio-Visual Performative Media include cartoons, talk-show clips, and short-form videos on platforms such as TikTok or Reels. These formats

167 incorporate speech, prosody, facial expressions,
168 music, and editing rhythm, making interpretation
169 contingent on fine-grained temporal alignment and
170 pragmatic inference (Marone, 2016).

171 Across all categories, data form shapes not only
172 perceptual requirements but also the types of rea-
173 soning and alignment needed for understanding
174 creative meaning.

175 2.2 Meaning-Making Mechanisms: The 176 “How”

177 Creative meaning often arises from rhetorical and
178 cognitive mechanisms that introduce ambiguity,
179 surprise, or abstraction. We highlight four recurrent
180 mechanisms that cut across data forms.

181 **Incongruity** refers to a mismatch between expect-
182 ation and observation and underlies many forms
183 of humor, irony, and sarcasm (Forabosco, 1992;
184 Veale, 2004). In multimodal settings, incongruity
185 often emerges from cross-modal clashes—e.g.,
186 a benign image paired with an unexpected cap-
187 tion—requiring models to reconcile conflicting
188 cues (Schifanella et al., 2016; Farabi et al., 2024).

189 **Analogy and Conceptual Mapping** construct
190 meaning by projecting structure from a familiar
191 source domain onto a target concept. This mech-
192 anism unifies *metaphor* and *symbolism* (Lakoff
193 and Johnson, 2024), where images, language, or
194 other modalities jointly instantiate cross-domain
195 correspondences that must be interpreted non-
196 literally (Refaie, 2003; Foss, 2004).

197 **Exaggeration and Hyperbole** amplify attributes
198 for comic or emphatic effect. In multimodal arti-
199 facts, exaggeration often arises through visual dis-
200 tortion, intensified motion, or prosodic emphasis,
201 with meaning emerging from cross-modal ampli-
202 fication rather than any single signal (Kreuz et al.,
203 1996; Zhang and Wan, 2024).

204 **Narrative Structure** operates primarily in sequen-
205 tial media, where meaning is shaped by setups,
206 payoffs, and causal relations across time (Genette,
207 1980; Bruner, 1991). Understanding such con-
208 tent requires identifying how events and modalities
209 jointly signal story progression (Paval et al., 2025).

210 These mechanisms explain how creative artifacts
211 encode meaning beyond surface semantics, posing
212 challenges distinct from literal understanding.

213 2.3 Communicative Goals: The “Why”

214 The third dimension concerns communicative in-
215 tent—the purpose a creative artifact serves for its

216 audience. This pragmatic layer shapes how both
217 data forms and mechanisms should be interpreted.
218 We identify four recurring goals.

219 **Humor and Entertainment** aim to amuse or
220 delight. Beyond recognizing incongruity, mod-
221 els must infer humorous intent to distinguish
222 jokes from errors or misinformation (Raskin, 1979;
223 Hasan et al., 2019a; Hu et al., 2024).

224 **Satire and Social Critique** use irony or absurdity
225 to comment on political or social issues. Correctly
226 identifying satirical intent is crucial for avoiding lit-
227 eral misinterpretation (Burfoot and Baldwin, 2009;
228 Li et al., 2020; Nandy et al., 2024).

229 **Persuasion and Advertising** seek to influence be-
230 liefs or behavior rather than entertain. Creative ads
231 and propaganda memes employ metaphor, emo-
232 tional framing, and stylization to persuade audi-
233 ences (Forceville, 2002; Kumar et al., 2023; Wang
234 et al., 2025a).

235 **Emotion and Aesthetic Experience** aim to evoke
236 affective or aesthetic responses, such as nostalgia,
237 empathy, or wonder. These goals require interpret-
238 ing how multimodal cues jointly produce emotional
239 resonance (Bruner, 1991; Christ et al., 2024; Shi
240 et al., 2025a; Padó and Thomas, 2025).

241 Although communicative goals often overlap,
242 identifying the primary intent provides essential
243 context for interpretation (e.g., distinguishing satire
244 from literal claims). In our taxonomy, communica-
245 tive goals form the pragmatic dimension that com-
246 plements data forms and meaning-making mech-
247 anisms, together offering a unified view of how
248 creative meaning is constructed and interpreted.

249 3 Task Hierarchy: Recognition, 250 Interpretation, and Generation

251 We summarize existing tasks into a three-level ca-
252 pability hierarchy to provide a structured lens on
253 current models (Figure 2); importantly, each level
254 aligns with distinct evaluation signals that reflect
255 the core requirements of the tasks.

256 3.1 Level 1: Recognition

257 Understanding human creativity begins with recog-
258 nition tasks, ranging from binary classification,
259 multi-class classification and MLLMs perception.
260 Early work focuses on **binary classification**, where
261 systems determine whether an input exhibits a
262 rhetorical or non-literal phenomenon. Prior to large
263 language models, this paradigm dominates creativ-
264 ity research across humor (Yang et al., 2015; Chan-

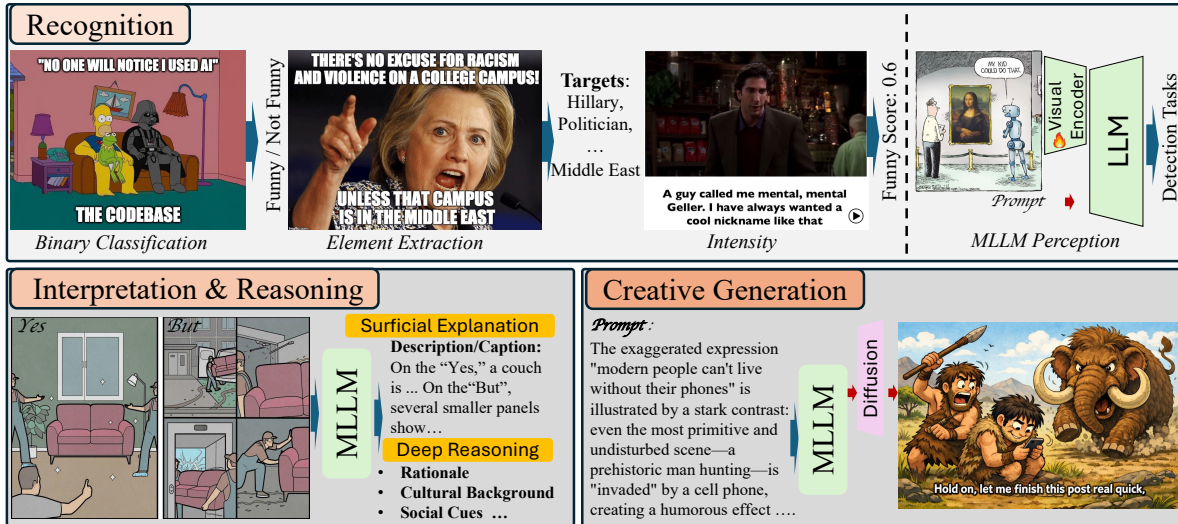


Figure 2: Capability-centric task hierarchy for multimodal human creativity.

drasekaran et al., 2016; Hasan et al., 2021), sarcasm (Schifanella et al., 2016; Liang et al., 2022; Liu et al., 2022a) etc., relying mainly on basic perceptual encoding, multimodal fusion, and surface-level cue discrimination.

Moving beyond binary decisions, recognition tasks extend to finer-grained settings (multi classification), such as **element extraction** and **intensity analysis**. Element extraction identifies constituent components of creative expressions, including persuasion techniques (Dimitrov et al., 2021b), entity roles (Sharma et al., 2022a, 2023b), sentiment and rhetorical categories in memes (Sharma et al., 2020), and punchlines or dialogue cues in multi-panel content (Kayatani et al., 2021; Liu et al., 2022b; Martínek et al., 2024). In contrast, intensity analysis quantifies the strength or perceived quality of creative effects, such as degrees of funniness (Alnajjar et al., 2022) or offensiveness (Kumari et al., 2025). These tasks require object-level localization, semantic disentanglement of overlapping signals, and judgement of creativity intensity.

MLLMs Perception exploits the language understanding and knowledge priors of LLMs to extend multimodal recognition, primarily by aligning semantic and visual representations through instruction tuning or finetune for detection-oriented tasks (Zhang and Wan, 2024; Huang et al., 2024; Gu et al., 2025). Compared to the others, these tasks still focus on recognition rather than deep reasoning, but they require stronger cross-modal semantic alignment, and language-grounded visual abstraction, enabling models to map perceptual cues to high-level labels using natural language supervision rather than handcrafted features.

For evaluation, recognition tasks are commonly

assessed using discriminative metrics such as accuracy, F1, or correlation with human ratings. While suitable for measuring perceptual alignment and cue sensitivity, these metrics fail to capture interpretive depth or causal understanding required for higher-level creative reasoning.

3.2 Level 2: Interpretation and Reasoning

Recognition reflects surface-level understanding, where models map multimodal cues to labels, while this level advances to deeper reasoning, requiring models to explain why creative effects arise. **Surficial Explanation** produces descriptive or paraprastic explanations (Hu et al., 2024) that verbalize creative content, but does not explicitly model underlying creative mechanisms or factors. This task requires natural language generation ability (Vaswani et al., 2017; Mann et al., 2020), cross-modal semantic alignment beyond label prediction (Liu et al., 2023a), and basic social and commonsense knowledge (Zhang et al., 2024b), enabling models to produce fluent, plausible descriptions by leveraging salient surface cues in a given instance without explicit multi-step reasoning.

Deep Interpretation and Reasoning requires models to explicitly reason about the mechanisms underlying creative effects through structured or theory-guided inference. It involves externalizing reasoning processes (e.g. CoT (Chen et al., 2024b), multi-stages (Tikhonov and Shtykovskiy, 2024)), abstracting non-literal concepts such as symbolism (Hussain et al., 2017; Yu et al., 2025) or metaphor (Akula et al., 2023; Saakyan et al., 2024), and integrating contextual (Wang et al., 2024b; Kumari et al., 2025) or external knowledge (Sharma et al., 2023c; Garg et al., 2025) to justify *why* a

creative effect arises rather than merely describing *what* is present.

At this level, Evaluation must go beyond fluency or surface plausibility. While superficial explanations can be assessed via semantic relevance and visual grounding (e.g., BERTScore (Zhang et al., 2019), BLEU (Papineni et al., 2002)), deep interpretation requires evaluating correct identification of rhetorical mechanisms and reasoning quality—such as correctness and consistency (Hu et al., 2024; Xiao et al., 2024)—motivating explanation-aware and evidence-based protocols (e.g., designed human evaluation and G-Eval (Liu et al., 2023b)).

3.3 Level 3: Creativity Generation

Creative expression generation spans diverse output forms—single-image (memes, humorous images, advertising posters), sequential (comics, visual narratives), and textual (captions, jokes)—with cross-modal consistency as a primary challenge: maintaining coherent creative intent and rhetorical structure across disparate representations. The field has evolved from **pattern-based approaches** that relied on predefined templates or rules, to **MLLM-guided generation** that produces flexible, context-aware content through learned multimodal representations. This shift requires models to move beyond pattern reuse toward deep interpretation and reasoning, including internalizing creative mechanisms (e.g., incongruity (Tanaka et al., 2024), metaphor (Chakrabarty et al., 2023), CoT (Zhong et al., 2024; Wang et al., 2024a)), aligning multimodal context with language generation (Shah-mohammadi et al., 2023), and synthesizing novel content grounded in learned representations and knowledge priors.

Creative generation requires evaluating both output quality and consistency with the intended creative intent. Existing automatic metrics largely assess surface similarity or visual quality and struggle to capture the diversity and subjectivity of human creativity, making human evaluation still the most reliable—yet costly—option.

4 Modeling Paradigms

Our literature review covers the decade from 2015 to 2025, capturing the paradigm shift from early neural architectures to modern Foundation Models. To ensure technical rigor, we included peer-reviewed research and high-impact pre-prints focused on non-literal multimodal intent, novel benchmarks, or creative modeling paradigms. We

excluded text-only studies and papers restricted to factual image captioning, as these lack the rhetorical or creative complexity central to this survey.

4.1 Task Specific Models: Non-MLLM Era

Prior to MLLMs, creativity understanding was dominated by task-specific discriminative architectures tightly coupled with individual tasks such as sarcasm, humor, or metaphor detection.

Early models focused on effective fusion mechanisms for heterogeneous features. Schifanella et al. (2016) first incorporated visual cues into sarcasm detection via separate encoders and concatenation, while Cai et al. (2019) showed hierarchical fusion of text, images, and attributes better captures cross-modal interactions. Zhang et al. (2021) introduced contrastive attention to explicitly model inter-modal incongruity for finer detection of cross-modal discrepancies.

Beyond fusion design, Later work emphasized implicit knowledge and commonsense. HKT (Hasan et al., 2021) injected humor-related knowledge into Transformer architecture for deeper incongruity modeling. Lee et al. (2021) and Liang et al. (2022) leveraged object-level visual representations to construct richer contextual embeddings and cross-modal graphs, facilitating localized reasoning over visual–textual conflicts. Liu et al. (2022a) incorporated external commonsense and semantic knowledge into hierarchical congruity modeling, improving implicit intent interpretation.

Overall, non-MLLM approaches relied on carefully-engineered fusion strategies and explicit knowledge injection. While effective on specific tasks, these models remained limited in scalability and generalization, as summarized in related surveys (Sharma et al., 2022b; Farabi et al., 2024).

4.2 Large Models: Current Era

With the rise of MLLMs, modeling paradigms for human creativity have fundamentally shifted. Beyond task-specific discriminative architectures, LLM-based systems bring strong linguistic priors, world knowledge, and abstraction capabilities that enable not only recognition, but also explanation and interpretation of creative phenomena. This transition reframes multimodal creativity understanding from handcrafted feature fusion toward alignment-driven representation learning, multi-step reasoning, and knowledge-grounded interpretation, laying the foundation for the three emerging paradigms in the MLLM era: perceptual represen-

tation and multimodal alignment, reasoning mechanisms, and external knowledge integration.

Perceptual Representation and Multimodal Alignment. Modeling of human creativity with MLLMs is shifting from task-specific architectures toward large-scale multimodal alignment, where perception and semantics are jointly learned. Models such as SoMeLVLM (Zhang et al., 2024b) show that instruction tuning on creativity-oriented corpora enables handling humor, memes, and implicit intent without explicit symbolic pipelines. Complementary directions explore architectural modularity (Yu et al., 2024) to decouple perception and reasoning, and text-centric representation strategies (Hasan et al., 2023; Baluja, 2025) that reduce fusion complexity while improving interpretability.

Beyond single-panel inputs, recent work extends alignment to structured and long-form creativity settings. YesBut-v2 (Liang et al., 2025) scales instruction tuning to multi-panel memes requiring cross-panel reasoning; MeSum (Khan et al., 2024) frames meme understanding as multimodal summarization by fusing image, text, and audio with vision backbones and LLMs; FunnyNet-W (Liu et al., 2024b) adapts aligned representations to video-based humor; and ComicsPAP (Vivoli et al., 2025) targets multi-panel comic understanding. Closely related, AI4VA-FG (Chen et al., 2025) improves fine-grained comic understanding by learning region-aware visual representations via reinforcement learning, showing that more precise perceptual alignment at the region level further enhances multimodal understanding in complex visual narratives.

In general, these methods offer high scalability by using broad visual-linguistic priors to interpret diverse creative content. While effective for “surface-level” tasks like meme classification, it often suffers from semantic shallowness, relying on superficial correlations rather than deep rhetoric. This frequently leads to hallucinated intent, where models invent plausible but incorrect justifications for creative artifacts they cannot truly decode.

Reasoning Mechanics. Reasoning-based approaches emphasize *explicit* intermediate inference processes for transparent, controllable interpretation rather than implicit representation. A prominent line adopts *CoT* reasoning (Wei et al., 2022) to expose intermediate semantic steps. Tanaka et al. (2024) employs incongruity theory by prompting models to reason through expectation violation be-

fore generating humor, while Gu et al. (2025) extends this at scale with human-annotated multimodal *CoT* traces for harmful memes, demonstrating that explicit reasoning substantially improves both accuracy and interpretability over pattern-based baselines.

Beyond vanilla *CoT*, recent work proposes *multi-step and theory-guided reasoning*. Humor Mechanic (Tikhonov and Shtykovskiy, 2024) demonstrates that structured multi-step reasoning grounded in humor theory generalizes better than single-pass predictions, motivating multimodal extensions. Zhang et al. (2025) explicitly encodes humor theories into a staged framework, decomposing understanding into sequential conflict detection, resolution, and affective judgment.

Other studies highlight *abstract, symbolic, and cross-modal reasoning* as essential for creativity understanding. Loakman et al. (2024) shows MLLMs can perform explicit reasoning over sound symbolism by aligning abstract auditory–visual associations, while Liu et al. (2025) probes reasoning over non-literal, abstract visual concepts. Kundu et al. (2025) and Qiu et al. (2024) further demonstrate that interpreting metaphors, diagrams, or symbolic graphics requires explicit multi-step reasoning over latent structure rather than surface perception.

Collectively, these studies demonstrate that explicit reasoning enhances transparency by breaking down rhetorical conflicts into logical steps, enabling understanding that extends beyond simple pattern matching. However, this comes at the cost of increased computational latency and the risk of reasoning drift, where models over-analyze simple humor or become trapped in “logical loops” when faced with ambiguity.

External Knowledge Integration. Understanding multimodal creative expression often requires knowledge beyond what is explicitly present in the input, motivating the integration or retrieval of external knowledge. MemeX (Sharma et al., 2023c) explicitly formulates meme understanding as an explanatory evidence retrieval task, retrieving background documents to explain implicit cultural or contextual references. Similarly, Garg et al. (2025) and Kumari et al. (2025) incorporate external commonsense and cultural knowledge to better interpret indecent or harmful memes, showing that offensiveness and intent cannot be inferred from surface cues alone. In advertising, KAFA (Jia et al., 2023) augments vision–language models with structured

ad-related knowledge to adapt visual features toward persuasion-aware semantics. Beyond static knowledge injection, Tang et al. (2024) leverages retrieval-based external knowledge to improve robustness in sarcasm detection under domain shift.

These works demonstrate that external knowledge integration via retrieval or adaptation is crucial for grounding the understanding in cultural, social, and commonsense contexts. However, it might introduce retrieval noise. Irrelevant or distractive background data can interfere with the model’s internal logic, often obscuring the core message or distracting the system from the specific visual cues that drive the creative intent.

5 Datasets and Benchmarks

In this section, we organize resources by task complexity, ranging from low-level recognition to high-level interpretation and creative generation. Detailed summaries of datasets and benchmarks are provided in Appendix Tables 1 and 2.

Recognition Resources. Early benchmarks for multimodal creativity primarily focus on existence-level recognition using simple textual or visual cues. For example, in the domain of humor, initial studies such as *Inside Jokes* (Shahaf et al., 2015) and early meme-based work ground humor detection in textual comparison and incongruity modeling (Tanaka et al., 2022). Subsequent benchmarks (e.g., UR-FUNNY (Hasan et al., 2019b), MUCH (Guo et al., 2024), and the Laughing Machine (Kayatani et al., 2021)) expand recognition to audiovisual cues, dialogue context, and temporal dynamics. More recent work further moves toward fine-grained and socially grounded signals, including laughter intensity and timing (Alnajjar et al., 2022), multimodal laughter detection (Kuznetsova and Strapparava, 2024), comic mischief (Baharlouei et al., 2024), and real-world reaction and stand-up comedy videos (Hyun et al., 2024; Barriere et al., 2025).

Overall, recognition datasets have shifted toward fine-grained, sequential annotations but remain classification-oriented and insufficient for evaluating deep reasoning or generation.

Understanding and Generation Resources. Understanding and generation are closely linked at the dataset level, as both are grounded in paired multimodal inputs (e.g., image/video–text). Thus, many datasets designed for understanding naturally support generation-oriented evaluation, since inter-

preting creativity often requires producing structured language grounded in multimodal context. In the MLLM era, datasets increasingly move beyond descriptive captioning toward reasoning-aware and knowledge-augmented evaluation. Early resources such as MemeCap (Hwang and Shwartz, 2023), EXCLAIM (Sharma et al., 2023a), and OxfordTVG-HIC (Li et al., 2023) extend captioning into explanatory understanding by requiring models to articulate intent or communicative goals, though they largely remain at a surface-level reasoning stage. More recent datasets, including Mementos (Wang et al., 2024b), MemeMind (Gu et al., 2025), and MemeGuard (Jha et al., 2024a), explicitly incorporate human-written rationales or CoT supervision to support deeper reasoning.

Complementarily, benchmarks such as the Yes-But series (Hu et al., 2024; Nandy et al., 2024; Liang et al., 2025), V-Flute (Saakyan et al., 2025), UnPIE (Chung et al., 2024), and Zhong et al. (2024) probe reasoning over incongruity, symbolism, and implicit knowledge, while video-focused datasets (e.g., Smile (Hyun et al., 2024), V-HUB (Shi et al., 2025b), EmoVid (Qiu et al., 2025)) further extend understanding and generation to temporal, affective, and social dynamics.

Overall, multimodal creativity datasets are moving beyond static recognition toward reasoning-oriented understanding across visual, temporal, and social cues, underscoring the need for fine-grained annotations and explicit reasoning supervision.

6 Challenges and Future Directions

Despite rapid progress in MLLMs on literal scene perception, a substantial gap remains between recognizing *what is depicted* and interpreting *what is meant* in human creative expression. Understanding these contents requires moving beyond physical description toward socio-cultural, rhetorical, and value-laden interpretation. This shift introduces a distinct set of technical and conceptual challenges. In this section, we synthesize the core challenges and propose future directions for the community.

6.1 The Evaluation Crisis

A central limitation in current research lies in evaluation. Most benchmarks rely on multiple-choice questions (MCQs) or binary classification tasks (Hessel et al., 2023a; Hu et al., 2024; Yang et al., 2024). While these formats enable scalable comparison, they are poorly aligned with the nature of creative interpretation because:

638 First, such benchmarks often permit **surface-**
639 **level shortcut learning** (Li et al., 2024). Models
640 may exploit lexical cues, dataset artifacts, or answer
641 priors without engaging in genuine rhetorical
642 reasoning. Correctly selecting a label for a joke
643 or meme does not guarantee that the model under-
644 stands the underlying humor, irony, or critique.

645 Second, creativity interpretation is inherently
646 **open-ended and subjective**. They admit multi-
647 ple plausible readings depending on context and
648 audience. Reducing these phenomena to a single
649 “correct” option obscures interpretive depth and dis-
650 courages explanation-driven reasoning.

651 Future evaluation must therefore move beyond
652 discriminative setups toward *explainable and gen-*
653 *erative assessments*, where models are required to
654 articulate why an artifact is humorous, ironic, or
655 persuasive, and to justify their interpretations.

656 6.2 Data Granularity and Domain Diversity

657 To move beyond surface recognition, AI requires a
658 new generation of datasets that goes beyond simple
659 labeling. The first is need for granularity, where
660 existing datasets often treat “creativity” as a mono-
661 lithic attribute. There is a pressing need for fine-
662 grained evaluation, separating the literal scene de-
663 scription from the figurative intent, the emotional
664 subtext, and the specific rhetorical device.

665 Moreover, the diverse domains of data remain
666 a problem. While social media memes are abun-
667 dant, more specialized creative forms, such as polit-
668 ical cartoons, visual satire, and culturally-specific
669 rituals, remain underrepresented. Creating high-
670 fidelity, expert-annotated datasets for these “long-
671 tail” domains is essential for robust model training.

672 6.3 Social Knowledge

673 Creative artifacts are deeply embedded in so-
674 cial context. Their meaning often depends on
675 shared background knowledge, implicit norms,
676 and assumptions about audience beliefs. Current
677 MLLMs, despite strong semantic modeling, remain
678 limited in their ability to reason about such social
679 pragmatics (Hu and Shu, 2023).

680 In particular, many failures stem from a lack of
681 **theory-of-mind-like reasoning**: models struggle
682 to infer whose perspective is being expressed, what
683 beliefs are being challenged, or how an audience
684 is expected to react. Moreover, creative expression
685 is frequently tied to **temporally-evolving events**,
686 trends, or controversies, requiring access to up-to-
687 date external knowledge.

688 Future systems must therefore integrate social
689 reasoning, explicit modeling of intent and audience,
690 and retrieval mechanisms that ground interpretation
691 in contemporary and community-specific context,
692 rather than treating creative artifacts as isolated.

693 6.4 Safety, Ethics, and Ownership

694 **Detecting Hidden Harm and Bias.** Creative me-
695 dia often embeds harmful narratives through irony,
696 sarcasm, and coded symbols, making subtle toxici-
697 ty harder to flag than overt hate speech (Sharma
698 et al., 2022b). Detection is further complicated
699 by the multimodal nature of these formats—where
700 meaning emerges only from the interplay of text
701 and imagery—and cultural subjectivity, which of-
702 ten leads to inconsistent human annotations and
703 false positives (Cao et al., 2024).

704 **Ensuring Safe and Unbiased AI Outputs.** Gen-
705 erative models frequently mirror demographic bi-
706 ases, reinforcing societal prejudices at scale. While
707 alignment and filtering are standard mitigations,
708 they present a “safety-utility” trade-off: overly con-
709 servative filters can stifle benign creative expres-
710 sion, yet remain vulnerable to brittle multimodal
711 adversarial attacks (Lee et al., 2025). Ensuring safe,
712 unbiased interaction with human-centered creative
713 media remains an open challenge requiring layered
714 alignment and robust, continuous monitoring.

715 **Copyright and Intellectual Property.** Creative
716 datasets often scrape content from social media
717 and art platforms without explicit consent from
718 creators. This raises critical concerns regarding
719 data provenance and the “Right to Style.” As mod-
720 els move toward generation, the ability to mimic
721 a specific artist’s style or a writer’s satirical tone
722 without attribution poses a threat to the economic
723 and intellectual rights of human creators.

724 6.5 Conclusion

725 Multimodal human creative expression challenges
726 AI systems as its meaning extends beyond sur-
727 face perception. Achieving robust understand-
728 ing and generation therefore requires moving to-
729 ward socially and rhetorically grounded reasoning.
730 This survey introduces a unified taxonomy and a
731 capability-centric task hierarchy that clarify how
732 creative understanding progresses from recognition
733 to interpretation and generation. We further outline
734 future directions for meaningful, reliable, and re-
735 sponsible engagement with human-created media.

736 Limitations

737 Despite growing interest in multimodal creative ex-
738 pression, this survey has several limitations. First,
739 most existing research and datasets focus on West-
740 ern, internet-centric creative forms (e.g., memes,
741 cartoons), leaving many cultural traditions and non-
742 mainstream media underexplored. Moreover, while
743 our taxonomy offers a unified organizational frame-
744 work, real-world creative artifacts often span mul-
745 tiple data forms, mechanisms, and goals, making
746 strict categorization imperfect. Finally, as multi-
747 modal models evolve rapidly, some observations
748 may not fully generalize to future architectures or
749 training paradigms.

750 Ethical considerations

751 Understanding creative expression poses distinct
752 ethical challenges. Creative media may poten-
753 tially convey harmful or sensitive content implicitly
754 through humor, irony, or symbolism, increasing
755 the risk of misinterpretation, bias amplification,
756 or over-censorship. Dataset bias and cultural im-
757 balance further threaten fairness and robustness,
758 particularly for marginalized communities. In addi-
759 tion, many creative datasets raise unresolved copy-
760 right and ownership concerns, especially as mod-
761 els increasingly transition from understanding to
762 generation and style imitation. Addressing these is-
763 sues requires context-aware evaluation, transparent
764 dataset practices, and greater emphasis on inter-
765 pretability and human oversight when deploying
766 such systems.

767 References

768 Panos Achlioptas, Maks Ovsjanikov, Kilichbek Hay-
769 darov, Mohamed Elhoseiny, and Leonidas J Guibas.
770 2021. Artemis: Affective language for visual art. In
771 *Proceedings of the IEEE/CVF Conference on Com-
772 puter Vision and Pattern Recognition*, pages 11569–
773 11579.

774 Siddhant Agarwal, Shivam Sharma, Preslav Nakov, and
775 Tanmoy Chakraborty. 2024. Mememqa: multimodal
776 question answering for memes via rationale-based
777 inferencing. *arXiv preprint arXiv:2405.11215*.

778 Arjun R Akula, Brendan Driscoll, Pradyumna Narayana,
779 Soravit Changpinyo, Zhiwei Jia, Suyash Damle,
780 Garima Pruthi, Sugato Basu, Leonidas Guibas,
781 William T Freeman, and 1 others. 2023. Metaclue:
782 Towards comprehensive visual metaphors research.
783 In *Proceedings of the IEEE/CVF conference on com-
784 puter vision and pattern recognition*, pages 23201–
785 23211.

Khalid Alnajjar, Mika Hämmäläinen, Jörg Tiedemann, Jorma Laaksonen, and Mikko Kurimo. 2022. When to laugh and how hard? a multimodal approach to detecting humor and its intensity. *arXiv preprint arXiv:2211.01889*. 786
787
788
789
790

Miriam Amin and Manuel Burghardt. 2020. A survey on approaches to computational humor generation. In *Proceedings of the 4th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 29–41. 791
792
793
794
795
796

Elaheh Baharlouei, Mahsa Shafaei, Yigeng Zhang, Hugo Jair Escalante, and Thamar Solorio. 2024. Labeling comic mischief content in online videos with a multimodal hierarchical-cross-attention model. *arXiv preprint arXiv:2406.07841*. 797
798
799
800
801

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025. *Qwen3-vl technical report*. Preprint, arXiv:2511.21631. 802
803
804
805
806
807
808

Ashwin Baluja. 2025. Text is not all you need: Multimodal prompting helps llms understand humor. In *Proceedings of the 1st Workshop on Computational Humor (CHum)*, pages 9–17. 809
810
811
812

Kate Barnes, Tiernon R. Riesenmy, Minh Duc Trinh, Eli Lleshi, Nóra Balogh, and Roland Molontay. 2020. Dank or not? analyzing and predicting the popularity of memes on reddit. *Applied Network Science*, 6. 813
814
815
816

Valentin Barriere, Nahuel Gomez, Leo Hemamou, Sofia Callejas, and Brian Ravenet. 2025. Standup4ai: A new multilingual dataset for humor detection in stand-up comedy videos. *arXiv preprint arXiv:2505.18903*. 817
818
819
820

Margaret A Boden. 1998. Creativity and artificial intelligence. *Artificial intelligence*, 103(1-2):347–356. 821
822

Digbalay Bose, Rajat Hebbar, Tiantian Feng, Krishna Somandepalli, Anfeng Xu, and Shrikanth Narayanan. 2023. Mm-au: Towards multimodal understanding of advertisement videos. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 86–95. 823
824
825
826
827
828

Jerome Bruner. 1991. The narrative construction of reality. *Critical inquiry*, 18(1):1–21. 829
830

Clint Burfoot and Timothy Baldwin. 2009. Automatic satire detection: Are you having a laugh? In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 161–164, Suntec, Singapore. Association for Computational Linguistics. 831
832
833
834
835

Yitao Cai, Huiyu Cai, and Xiaojun Wan. 2019. Multimodal sarcasm detection in twitter with hierarchical fusion model. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 2506–2515. 836
837
838
839
840

841	Jingtao Cao, Zheng Zhang, Hongru Wang, Bin Liang,	advanced reasoning, multimodality, long context, and	897
842	Hao Wang, and Kam-Fai Wong. 2024. Ospc: De-	next generation agentic capabilities. <i>arXiv preprint</i>	898
843	etecting harmful memes with large language model as	<i>arXiv:2507.06261</i> .	899
844	a catalyst. In <i>Companion Proceedings of the ACM</i>		
845	<i>Web Conference 2024</i> , pages 1892–1895.		
846	Santiago Castro, Devamanyu Hazarika, Verónica Pérez-	Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj	900
847	Rosas, Roger Zimmermann, Rada Mihalcea, and Sou-	Alam, Fabrizio Silvestri, Hamed Firooz, Preslav	901
848	janya Poria. 2019. Towards multimodal sarcasm	Nakov, and Giovanni Da San Martino. 2021a. Detect-	902
849	detection (an <code>_obviously_</code> perfect paper). <i>arXiv</i>	ing propaganda techniques in memes. In <i>Proceedings</i>	903
850	<i>preprint arXiv:1906.01815</i> .	<i>of the 59th annual meeting of the association for com-</i>	904
851	Tuhin Chakrabarty, Arkadiy Saakyan, Olivia Winn,	<i>putational linguistics and the 11th international joint</i>	905
852	Artemis Panagopoulou, Yue Yang, Marianna Apid-	<i>conference on natural language processing (volume</i>	906
853	ianaki, and Smaranda Muresan. 2023. I spy a	<i>1: long papers)</i> , pages 6603–6617.	907
854	metaphor: Large language models and diffusion		
855	models co-create visual metaphors. <i>arXiv preprint</i>	Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj	908
856	<i>arXiv:2305.14724</i> .	Alam, Fabrizio Silvestri, Hamed Firooz, Preslav	909
857	Arjun Chandrasekaran, Ashwin K Vijayakumar, Stanis-	Nakov, and Giovanni Da San Martino. 2021b.	910
858	law Antol, Mohit Bansal, Dhruv Batra, C Lawrence	Semeval-2021 task 6: Detection of persuasion tech-	911
859	Zitnick, and Devi Parikh. 2016. We are humor be-	niques in texts and images. In <i>Proceedings of the</i>	912
860	ings: Understanding and predicting visual humor. In	<i>15th international workshop on semantic evaluation</i>	913
861	<i>Proceedings of the IEEE Conference on Computer</i>	<i>(SemEval-2021)</i> , pages 70–98.	914
862	<i>Vision and Pattern Recognition</i> , pages 4603–4612.		
863	Yule Chen, Yufan Ren, and Sabine Süssstrunk. 2025.	Shafkat Farabi, Tharindu Ranasinghe, Diptesh Kanojia,	915
864	Zooming into comics: Region-aware rl improves	Yu Kong, and Marcos Zampieri. 2024. A survey	916
865	fine-grained comic understanding in vision-language	of multimodal sarcasm detection. <i>arXiv preprint</i>	917
866	models. <i>arXiv preprint arXiv:2511.06490</i> .	<i>arXiv:2410.18882</i> .	918
867	Yuyan Chen, Songzhou Yan, Zhihong Zhu, Zhixu Li,	Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rah-	919
868	and Yanghua Xiao. 2024a. Xmecap: Meme caption	wan, and Sune Lehmann. 2017. Using millions of	920
869	generation with sub-image adaptability. In <i>Proceed-</i>	emoji occurrences to learn any-domain representa-	921
870	<i>ings of the 32nd ACM International Conference on</i>	tions for detecting sentiment, emotion and sarcasm.	922
871	<i>Multimedia</i> , pages 3352–3361.	<i>arXiv preprint arXiv:1708.00524</i> .	923
872	Yuyan Chen, Yichen Yuan, Panjun Liu, Dayiheng Liu,	Giovanantonio Forabosco. 1992. Cognitive aspects of	924
873	Qinghao Guan, Mengfei Guo, Haiming Peng, Bang	the humor process: The concept of incongruity.	925
874	Liu, Zhixu Li, and Yanghua Xiao. 2024b. Talk		
875	funny! a large-scale humor response dataset with	Charles Forceville. 2002. <i>Pictorial metaphor in adver-</i>	926
876	chain-of-humor interpretation. In <i>Proceedings of</i>	<i>tising</i> . Routledge.	927
877	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	Sonja K Foss. 2004. Theory of visual rhetoric. In	928
878	ume 38, pages 17826–17834.	<i>Handbook of visual communication</i> , pages 163–174.	929
879	Lukas Christ, Shahin Amiriparian, Manuel Milling, Il-	Routledge.	930
880	han Aslan, and Björn Schuller. 2024. Modeling emo-	Sonja K Foss. 2017. <i>Rhetorical criticism: Exploration</i>	931
881	tional trajectories in written stories utilizing trans-	<i>and practice</i> . Waveland Press.	932
882	formers and weakly-supervised learning. In <i>Find-</i>	Rahul Garg, Trilok Padhi, Hemang Jain, Ugur Kursuncu,	933
883	<i>ings of the Association for Computational Linguis-</i>	and Ponnurangam Kumaraguru. 2025. Just kiddin’:	934
884	<i>tics: ACL 2024</i> , pages 7144–7159.	Knowledge infusion and distillation for detection of	935
885	Jiwan Chung, Seungwon Lim, Jaehyun Jeon, Seung-	indecent memes. In <i>Findings of the Association for</i>	936
886	been Lee, and Youngjae Yu. 2024. Can visual lan-	<i>Computational Linguistics: ACL 2025</i> , pages 23067–	937
887	guage models resolve textual ambiguity with visual	23086.	938
888	cues? let visual puns tell you! <i>arXiv preprint</i>	Gérard Genette. 1980. <i>Narrative discourse: An essay</i>	939
889	<i>arXiv:2410.01023</i> .	<i>in method</i> , volume 3. Cornell University Press.	940
890	Simon Colton and Geraint A Wiggins. 2012. Computa-	Hexiang Gu, Qifan Yu, Saihui Hou, Zhiqin Fang, Huijia	941
891	tional creativity: The final frontier? In <i>ECAI 2012</i> ,	Wu, and Zhaofeng He. 2025. Mememind: A large-	942
892	pages 21–26. IOS Press.	scale multimodal dataset with chain-of-thought rea-	943
893	Gheorghe Comanici, Eric Bieber, Mike Schaekermann,	soning for harmful meme detection. <i>arXiv preprint</i>	944
894	Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-	<i>arXiv:2506.18919</i> .	945
895	cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and	Diandian Guo, Cong Cao, Fangfang Yuan, Yanbing	946
896	1 others. 2025. Gemini 2.5: Pushing the frontier with	Liu, Guangjie Zeng, Xiaoyan Yu, Hao Peng, and	947
		Philip S Yu. 2025. Multi-view incongruity learning	948
		for multimodal sarcasm detection. In <i>Proceedings of</i>	949
		<i>the 31st International Conference on Computational</i>	950
		<i>Linguistics</i> , pages 1754–1766.	951

952	Hongyu Guo, Wenbo Shang, Xueyao Zhang, Shubo Zhang, Xu Han, and Binyang Li. 2024. Much: A multimodal corpus construction for conversational humor recognition based on chinese sitcom. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 11692–11698.	1010
953		1011
954		1012
955		1013
956		
957		1014
958		1015
959		1016
960	SI Harini, Somesh Singh, Yaman K Singla, Aanisha Bhattacharyya, Veeky Baths, Changyou Chen, Rajiv Ratn Shah, and Balaji Krishnamurthy. 2025. Long-term ad memorability: Understanding & generating memorable ads. In <i>2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pages 5707–5718. IEEE.	1017
961		1018
962		1019
963		1020
964		
965		1021
966		1022
967	Md Kamrul Hasan, Md Saiful Islam, Sangwu Lee, Wasifur Rahman, Iftekhar Naim, Mohammed Ibrahim Khan, and Ehsan Hoque. 2023. Textmi: Textualize multimodal information for integrating non-verbal cues in pre-trained language models. <i>arXiv preprint arXiv:2303.15430</i> .	1023
968		1024
969		1025
970		1026
971		
972		1027
973	Md Kamrul Hasan, Sangwu Lee, Wasifur Rahman, Amir Zadeh, Rada Mihalcea, Louis-Philippe Morency, and Ehsan Hoque. 2021. Humor knowledge enriched transformer for understanding multimodal humor. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pages 12972–12980.	1028
974		1029
975		1030
976		1031
977		
978		1032
979		1033
980	Md Kamrul Hasan, Wasifur Rahman, AmirAli Bagher Zadeh, Jianyuan Zhong, Md Iftekhar Tanveer, Louis-Philippe Morency, and Mohammed (Ehsan) Hoque. 2019a. UR-FUNNY: A multimodal language dataset for understanding humor . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2046–2056, Hong Kong, China. Association for Computational Linguistics.	1034
981		1035
982		
983		1036
984		1037
985		1038
986		1039
987		1040
988		1041
989		
990	Md Kamrul Hasan, Wasifur Rahman, AmirAli Bagher Zadeh, Jianyuan Zhong, Md Iftekhar Tanveer, Louis-Philippe Morency, and Mohammed Ehsan Hoque. 2019b. Ur-funny: A multimodal language dataset for understanding humor . In <i>Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)</i> , pages 2046–2056.	1042
991		1043
992		1044
993		1045
994		1046
995		1047
996		
997		1048
998		1049
999	Ming Shan Hee, Wen-Haw Chong, and Ka-Wei Roy Lee. 2023. Decoding the underlying meaning of multimodal hateful memes. In <i>32nd International Joint Conference on Artificial Intelligence (IJCAI 2023)</i> . International Joint Conferences on Artificial Intelligence (IJCAI).	1050
1000		1051
1001		1052
1002		1053
1003		1054
1004		1055
1005	Jack Hessel, Ana Marasovic, Jena D. Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff, and Yejin Choi. 2023a. Do androids laugh at electric sheep? humor “understanding” benchmarks from the new yorker caption contest . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 688–714, Toronto, Canada. Association for Computational Linguistics.	1056
1006		1057
1007		1058
1008		1059
1009		1060
		1061
		1062
		1063
		1014
		1015
		1016
		1017
		1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
		1054
		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063

1064	Prince Jha, Raghav Jain, Konika Mandal, Aman Chadha, Sriparna Saha, and Pushpak Bhattacharyya. 2024a. Memeguard: An llm and vlm-based framework for advancing content moderation via meme intervention. <i>arXiv preprint arXiv:2406.05344</i> .	Roger J Kreuz, Richard M Roberts, Brenda K Johnson, and Eugenie L Bertus. 1996. Figurative language occurrence and co-occurrence in contemporary literature. <i>Advances in Discourse Processes</i> , 52:83–98.	1120 1121 1122 1123
1065			
1066			
1067			
1068			
1069	Prince Jha, Krishanu Maity, Raghav Jain, Apoorv Verma, Sriparna Saha, and Pushpak Bhattacharyya. 2024b. Meme-ingful analysis: Enhanced understanding of cyberbullying in memes through multimodal explanations. <i>arXiv preprint arXiv:2401.09899</i> .	Yaman Kumar, Rajat Jha, Arunim Gupta, Milan Aggarwal, Aditya Garg, Tushar Malyan, Ayush Bhargava, Rajiv Ratn Shah, Balaji Krishnamurthy, and Changyou Chen. 2023. Persuasion strategies in advertisements. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 37, pages 57–66.	1124 1125 1126 1127 1128 1129
1070			
1071			
1072			
1073			
1074	Mengzhao Jia, Can Xie, and Liqiang Jing. 2024. Debiasing multimodal sarcasm detection with contrastive learning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 38, pages 18354–18362.	Gitanjali Kumari, Jitendra Solanki, and Asif Ekbal. 2025. Memedetoxnet: Balancing toxicity reduction and context preservation. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 25076–25098.	1130 1131 1132 1133 1134
1075			
1076			
1077			
1078			
1079	Zhiwei Jia, Pradyumna Narayana, Arjun Akula, Garima Pruthi, Hao Su, Sugato Basu, and Varun Jampani. 2023. Kafa: Rethinking image ad understanding with knowledge-augmented feature adaptation of vision-language models. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)</i> , pages 772–785.	Manishit Kundu, Sumit Shekhar, and Pushpak Bhattacharyya. 2025. Looking beyond the pixels: Evaluating visual metaphor understanding in vlms. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 23137–23158.	1135 1136 1137 1138 1139
1080			
1081			
1082			
1083			
1084			
1085			
1086	Narendra Nath Joshi. 2025. Evaluating human perception and bias in ai-generated humor. In <i>Proceedings of the 1st Workshop on Computational Humor (CHum)</i> , pages 79–87.	Anna Kuznetsova and Carlo Strapparava. 2024. Multimodal and multilingual laughter detection in stand-up comedy videos. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 11884–11889.	1140 1141 1142 1143 1144 1145
1087			
1088			
1089			
1090	Antonios Kalloniatis and Panagiotis Adamidis. 2024. Computational humor recognition: a systematic literature review. <i>Artificial Intelligence Review</i> , 58(2):43.	George Lakoff and Mark Johnson. 2024. <i>Metaphors we live by</i> . University of Chicago press.	1146 1147
1091			
1092			
1093	Sai Kartheek Reddy Kasu, Mohammad Zia Ur Rehman, Shahid Shafi Dar, Rishi Bharat Junghare, Dhanvin Sanjay Namboodiri, and Nagendra Kumar. 2025. D-humor: Dark humor understanding via multimodal open-ended reasoning—a benchmark dataset and method. <i>arXiv preprint arXiv:2509.06771</i> .	DongGeon Lee, Joonwon Jang, Jihae Jeong, and Hwanjo Yu. 2025. Are vision-language models safe in the wild? a meme-based benchmark study. <i>arXiv preprint arXiv:2505.15389</i> .	1148 1149 1150 1151
1094			
1095			
1096			
1097			
1098			
1099	Yuta Kayatani, Zekun Yang, Mayu Otani, Noa Garcia, Chenhui Chu, Yuta Nakashima, and Haruo Takemura. 2021. The laughing machine: Predicting humor in video. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 2073–2082.	Roy Ka-Wei Lee, Rui Cao, Ziqing Fan, Jing Jiang, and Wen-Haw Chong. 2021. Disentangling hate in online memes. In <i>Proceedings of the 29th ACM international conference on multimedia</i> , pages 5138–5147.	1152 1153 1154 1155
1100			
1101			
1102			
1103			
1104			
1105	Anas Anwarul Haq Khan, Tanik Saikh, Arpan Phukan, and Asif Ekbal. 2024. Hope ‘the paragraph guy’ explains the rest: Introducing mesum, the meme summarizer. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 6654–6668.	Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. 2024. Naturalbench: Evaluating vision-language models on natural adversarial samples. <i>Advances in Neural Information Processing Systems</i> , 37:17044–17068.	1156 1157 1158 1159 1160 1161 1162
1106			
1107			
1108			
1109			
1110			
1111	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. <i>Advances in neural information processing systems</i> , 33:2611–2624.	Lily Li, Or Levi, Pedram Hosseini, and David Broniatowski. 2020. A multi-modal method for satire detection using textual and visual cues. In <i>Proceedings of the 3rd NLP4IF Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda</i> , pages 33–38, Barcelona, Spain (Online). International Committee on Computational Linguistics (ICCL).	1163 1164 1165 1166 1167 1168 1169 1170
1112			
1113			
1114			
1115			
1116			
1117	Dayoon Ko, Sangho Lee, and Gunhee Kim. 2023. Can language models laugh at youtube short-form videos? <i>arXiv preprint arXiv:2310.14159</i> .	Runjia Li, Shuyang Sun, Mohamed Elhoseiny, and Philip Torr. 2023. Oxfordtv-gic: Can machine make humorous captions from images? In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 20293–20303.	1171 1172 1173 1174 1175
1118			
1119			

1286	Reuben Narad, Siddharth Suresh, Jiayi Chen, Pine SL	Chengjuan Ren, Ziyu Guo, Ping Zhang, and Yuhan Gao.	1339
1287	Dysart-Bricken, Bob Mankoff, Robert Nowak, Ji-	2024. Humor detection using deep learning in 10	1340
1288	fan Zhang, and Lalit Jain. 2025. Which llms get	years: A survey. <i>Métodos numéricos para cálculo y</i>	1341
1289	the joke? probing non-stem reasoning abilities with	<i>diseño en ingeniería: Revista internacional</i> , 40(1):1–	1342
1290	humorbench. <i>arXiv preprint arXiv:2507.21476</i> .	13.	1343
1291	Khoi PN Nguyen, Terrence Li, Derek Lou Zhou,	Yuriel Ryan, Rui Yang Tan, Kenny Tsu Wei Choo, and	1344
1292	Gabriel Xiong, Pranav Balu, Nandhan Alahari, Alan	Roy Ka-Wei Lee. 2025. Humor in pixels: Bench-	1345
1293	Huang, Tanush Chauhan, Harshavardhan Bala, Emre	marking large multimodal models understanding of	1346
1294	Guzelordu, and 1 others. 2025. Memeqa: Holistic	online comics. In <i>Findings of the Association for</i>	1347
1295	evaluation for meme understanding. In <i>Proceedings</i>	<i>Computational Linguistics: EMNLP 2025</i> , pages	1348
1296	<i>of the 63rd Annual Meeting of the Association for</i>	14024–14050.	1349
1297	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Arkadiy Saakyan, Shreyas Kulkarni, Tuhin Chakrabarty,	1350
1298	pages 18926–18946.	and Smaranda Muresan. 2024. V-flute: Visual figu-	1351
1299	Khoi PN Nguyen and Vincent Ng. 2024. Computational	rative language understanding with textual explana-	1352
1300	meme understanding: A survey. In <i>Proceedings of</i>	tions. <i>CoRR</i> .	1353
1301	<i>the 2024 Conference on Empirical Methods in Natu-</i>	Arkadiy Saakyan, Shreyas Kulkarni, Tuhin Chakrabarty,	1354
1302	<i>ral Language Processing</i> , pages 21251–21267.	and Smaranda Muresan. 2025. Understanding figura-	1355
1303	Xuan Ouyang, Senan Wang, Bouzhou Wang, Siyuan	rative meaning through explainable visual entailment.	1356
1304	Xiahou, Jinrong Zhou, and Yuekang Li. 2025. Laugh,	In <i>Proceedings of the 2025 Conference of the Na-</i>	1357
1305	relate, engage: Stylized comment generation for short	<i>tions of the Americas Chapter of the Association for</i>	1358
1306	videos. <i>arXiv preprint arXiv:2511.03757</i> .	<i>Computational Linguistics: Human Language Tech-</i>	1359
1307	Sebastian Padó and Kerstin Thomas. 2025. Artwork in-	<i>nologies (Volume 1: Long Papers)</i> , pages 1–23.	1360
1308	terpretation with vision language models: A case	Rossano Schifanella, Paloma De Juan, Joel Tetreault,	1361
1309	study on emotions and emotion symbols. <i>arXiv</i>	and Liangliang Cao. 2016. Detecting sarcasm in	1362
1310	<i>preprint arXiv:2511.22929</i> .	multimodal social platforms. In <i>Proceedings of the</i>	1363
1311	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	<i>24th ACM international conference on Multimedia</i> ,	1364
1312	Jing Zhu. 2002. Bleu: a method for automatic evalua-	pages 1136–1145.	1365
1313	tion of machine translation. In <i>Proceedings of the</i>	Christoph Schuhmann, Romain Beaumont, Richard	1366
1314	<i>40th annual meeting of the Association for Computa-</i>	Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti,	1367
1315	<i>tional Linguistics</i> , pages 311–318.	Theo Coombes, Aarush Katta, Clayton Mullis,	1368
1316	Sandro Paval, Pascal Meißner, and Ivan P. Yamshchikov.	Mitchell Wortsman, and 1 others. 2022. Laion-5b:	1369
1317	2025. ComicScene154: A scene dataset for comic	An open large-scale dataset for training next genera-	1370
1318	analysis . In <i>Proceedings of the 2025 Conference on</i>	tion image-text models. <i>Advances in neural informa-</i>	1371
1319	<i>Empirical Methods in Natural Language Processing</i> ,	<i>tion processing systems</i> , 35:25278–25294.	1372
1320	pages 31562–31568, Suzhou, China. Association for	Siddhant Bikram Shah, Shuvam Shiwakoti, Maheep	1373
1321	Computational Linguistics.	Chaudhary, and Haohan Wang. 2024. Meme-	1374
1322	Zeju Qiu, Weiyang Liu, Haiwen Feng, Zhen Liu,	clip: Leveraging clip representations for multi-	1375
1323	Tim Z Xiao, Katherine M Collins, Joshua B Tenen-	modal meme classification. <i>arXiv preprint</i>	1376
1324	baum, Adrian Weller, Michael J Black, and Bernhard	<i>arXiv:2409.14703</i> .	1377
1325	Schölkopf. 2024. Can large language models under-	Dafna Shahaf, Eric Horvitz, and Robert Mankoff. 2015.	1378
1326	stand symbolic graphics programs? <i>arXiv preprint</i>	Inside jokes: Identifying humorous cartoon captions.	1379
1327	<i>arXiv:2408.08313</i> .	In <i>Proceedings of the 21th ACM SIGKDD interna-</i>	1380
1328	Zongyang Qiu, Bingyuan Wang, Xingbei Chen,	<i>tional conference on knowledge discovery and data</i>	1381
1329	Yingqing He, and Zeyu Wang. 2025. Emovid: A mul-	<i>mining</i> , pages 1065–1074.	1382
1330	timodal emotion video dataset for emotion-centric	Hassan Shahmohammadi, Adhiraj Ghosh, and Hendrik	1383
1331	video understanding and generation. <i>arXiv preprint</i>	Lensch. 2023. Vipe: Visualise pretty-much every-	1384
1332	<i>arXiv:2511.11002</i> .	thing. <i>arXiv preprint arXiv:2310.10543</i> .	1385
1333	Victor Raskin. 1979. Semantic mechanisms of humor.	Chhavi Sharma, Deepesh Bhageria, William Scott,	1386
1334	In <i>Annual Meeting of the Berkeley Linguistics Society</i> ,	Srinivas Pykl, Amitava Das, Tanmoy Chakraborty,	1387
1335	pages 325–335.	Viswanath Pulabaigari, and Bjorn Gambäck. 2020.	1388
1336	Elisabeth El Refaie. 2003. Understanding visual	Semeval-2020 task 8: Memotion analysis—the visuo-	1389
1337	metaphor: The example of newspaper cartoons. <i>Vi-</i>	lingual metaphor! <i>arXiv preprint arXiv:2008.03781</i> .	1390
1338	<i>sual communication</i> , 2(1):75–95.	Shivam Sharma, Siddhant Agarwal, Tharun Suresh,	1391
		Preslav Nakov, Md Shad Akhtar, and Tanmoy	1392
		Chakraborty. 2023a. What do you meme? gener-	1393
		ating explanations for visual semantic role labelling	1394

1395	in memes. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pages 9763–9771.	
1396		
1397	Shivam Sharma, Md Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2022a. Disarm: Detecting the victims targeted by harmful memes. <i>arXiv preprint arXiv:2205.05738</i> .	
1398		
1399		
1400		
1401	Shivam Sharma, Firoj Alam, Md Shad Akhtar, Dimitar Dimitrov, Giovanni Da San Martino, Hamed Firooz, Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and Tanmoy Chakraborty. 2022b. Detecting and understanding harmful memes: A survey. <i>arXiv preprint arXiv:2205.04274</i> .	
1402		
1403		
1404		
1405		
1406		
1407	Shivam Sharma, Atharva Kulkarni, Tharun Suresh, Himanshi Mathur, Preslav Nakov, Md Shad Akhtar, and Tanmoy Chakraborty. 2023b. Characterizing the entities in harmful memes: Who is the hero, the villain, the victim? <i>arXiv preprint arXiv:2301.11219</i> .	
1408		
1409		
1410		
1411		
1412	Shivam Sharma, S Rameswaran, Udit Arora, Md Shad Akhtar, and Tanmoy Chakraborty. 2023c. Memex: Detecting explanatory evidence for memes via knowledge-enriched contextualization. In <i>Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 5272–5290.	
1413		
1414		
1415		
1416		
1417		
1418		
1419	Yi Shi, Wenlong Meng, Zhenyuan Guo, Chengkun Wei, and Wenzhi Chen. 2025a. Enhancing meme emotion understanding with multi-level modality enhancement and dual-stage modal fusion. <i>arXiv preprint arXiv:2511.11126</i> .	
1420		
1421		
1422		
1423		
1424	Zhengpeng Shi, Hengli Li, Yanpeng Zhao, Jianqun Zhou, Yuxuan Wang, Qinrong Cui, Wei Bi, Songchun Zhu, Bo Zhao, and Zilong Zheng. 2025b. V-hub: A visual-centric humor understanding benchmark for video llms. <i>arXiv preprint arXiv:2509.25773</i> .	
1425		
1426		
1427		
1428		
1429	Limor Shifman. 2013. <i>Memes in digital culture</i> . MIT press.	
1430		
1431	Yuchen Su, Yonghua Zhu, Ruofan Wang, Zijian Huang, Diana Benavides-Prado, and Michael Witbrock. 2025. A survey of pun generation: Datasets, evaluations and methodologies. <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 7375–7395.	
1432		
1433		
1434		
1435		
1436		
1437	Kohtaro Tanaka, Kohei Uehara, Lin Gu, Yusuke Mukuta, and Tatsuya Harada. 2024. Content-specific humorous image captioning using incongruity resolution chain-of-thought . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 2348–2367, Mexico City, Mexico. Association for Computational Linguistics.	
1438		
1439		
1440		
1441		
1442		
1443		
1444	Kohtaro Tanaka, Hiroaki Yamane, Yusuke Mori, Yusuke Mukuta, and Tatsuya Harada. 2022. Learning to evaluate humor in memes based on the incongruity theory. In <i>Proceedings of the Second Workshop on When Creative AI Meets Conversational AI</i> , pages 81–93.	
1445		
1446		
1447		
1448		
1449		
	Binghao Tang, Boda Lin, Haolong Yan, and Si Li. 2024. Leveraging generative large language models with visual instruction and demonstration retrieval for multimodal sarcasm detection. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1732–1742.	1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
	Li Tao, Shunsuke Nakamura, Xueting Wang, Tatsuya Kawahara, Gen Tamura, and Toshihiko Yamasaki. 2024. A large-scale television advertising dataset for detailed impression analysis. <i>Multimedia Tools and Applications</i> , 83(7):18779–18802.	1458
		1459
		1460
		1461
		1462
	Alexey Tikhonov and Pavel Shtykovskiy. 2024. Humor mechanics: Advancing humor generation with multi-step reasoning. <i>arXiv preprint arXiv:2405.07280</i> .	1463
		1464
		1465
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	1466
		1467
		1468
		1469
		1470
	Tony Veale. 2004. Incongruity in humor: Root cause or epiphenomenon?	1471
		1472
	Emanuele Vivoli, Artemis Llabrés, Mohamed Ali Souibgui, Marco Bertini, Ernest Valveny Llobet, and Dimosthenis Karatzas. 2025. Comicspap: understanding comic strips by picking the correct panel. In <i>International Conference on Document Analysis and Recognition</i> , pages 337–350. Springer.	1473
		1474
		1475
		1476
		1477
		1478
	Han Wang, Yilin Zhao, Dian Li, Xiaohan Wang, Gang Liu, Xuguang Lan, and Hui Wang. 2024a. Innovative thinking, infinite humor: Humor research of large language models through structured thought leaps. <i>arXiv preprint arXiv:2410.10370</i> .	1479
		1480
		1481
		1482
		1483
	Jiquan Wang, Lin Sun, Yi Liu, Meizhi Shao, and Zengwei Zheng. 2022. Multimodal sarcasm target identification in tweets. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8164–8175.	1484
		1485
		1486
		1487
		1488
	Wenshuo Wang, Ziyong Jiang, Junjie Wang, Mingyang Li, Jie Huang, Yuekai Huang, Zhiyuan Chang, Feiyan Duan, and Qing Wang. 2025a. Learning from mistakes: Enhancing harmful meme detection via misjudgment risk patterns. <i>arXiv preprint arXiv:2510.15946</i> .	1489
		1490
		1491
		1492
		1493
		1494
	Xiaochen Wang, Heming Xia, Jialin Song, Longyu Guan, Qingxiu Dong, Rui Li, Yixin Yang, Yifan Pu, Weiyao Luo, Yiru Wang, Xiangdi Meng, Wenjie Li, and Zhifang Sui. 2025b. Beyond single frames: Can LLMs comprehend implicit narratives in comic strip? In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 6436–6452, Suzhou, China. Association for Computational Linguistics.	1495
		1496
		1497
		1498
		1499
		1500
		1501
		1502
		1503
	Xiyao Wang, Yuhang Zhou, Xiaoyu Liu, Hongjin Lu, Yuancheng Xu, Feihong He, Jaehong Yoon, Taixi Lu,	1504
		1505

1506	Fuxiao Liu, Gedas Bertasius, and 1 others. 2024b.	<i>International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 8652–8661.	1561
1507	Mementos: A comprehensive benchmark for multi-		1562
1508	modal large language model reasoning over image		1563
1509	sequences. In <i>Proceedings of the 62nd Annual Meet-</i>		
1510	<i>ing of the Association for Computational Linguistics</i>	Jiajun Zhang, Shijia Luo, Ruikang Zhang, and Qi Su.	1564
1511	<i>(Volume 1: Long Papers)</i> , pages 416–442.	2025. Humorchain: Theory-guided multi-stage rea-	1565
		soning for interpretable multimodal humor genera-	1566
		tion. <i>arXiv preprint arXiv:2511.21732</i> .	1567
1512	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Jifan Zhang, Lalit Jain, Yang Guo, Jiayi Chen, Kuan	1568
1513	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	Zhou, Siddharth Suresh, Andrew Wagenmaker, Scott	1569
1514	and 1 others. 2022. Chain-of-thought prompting elic-	Sievert, Timothy T Rogers, Kevin G Jamieson, and	1570
1515	its reasoning in large language models. <i>Advances</i>	1 others. 2024a. Humor in ai: Massive scale crowd-	1571
1516	<i>in neural information processing systems</i> , 35:24824–	sourced preferences and benchmarks for cartoon cap-	1572
1517	24837.	tioning. <i>Advances in Neural Information Processing</i>	1573
		<i>Systems</i> , 37:125264–125286.	1574
1518	Geraint A Wiggins. 2006. A preliminary framework	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q	1575
1519	for description, analysis and comparison of creative	Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-	1576
1520	systems. <i>Knowledge-based systems</i> , 19(7):449–458.	uating text generation with bert. <i>arXiv preprint</i>	1577
		<i>arXiv:1904.09675</i> .	1578
1521	Ruiyu Xiao, Lei Wu, Yuhang Gou, Weinan Zhang, and	Xiaoqiang Zhang, Ying Chen, and Guangyuan Li. 2021.	1579
1522	Ting Liu. 2024. Prove your point!: Bringing proof-	Multi-modal sarcasm detection based on contrastive	1580
1523	enhancement principles to argumentative essay gen-	attention mechanism. In <i>CCF International Confer-</i>	1581
1524	eration. <i>arXiv preprint arXiv:2410.22642</i> .	<i>ence on Natural Language Processing and Chinese</i>	1582
		<i>Computing</i> , pages 822–833. Springer.	1583
1525	Shweta Yadav, Cornelia Caragea, Chenye Zhao, Naincy	Xinnong Zhang, Haoyu Kuang, Xinyi Mou, Hanjia Lyu,	1584
1526	Kumari, Marvin Solberg, and Tanmay Sharma. 2023.	Kun Wu, Siming Chen, Jiebo Luo, Xuan-Jing Huang,	1585
1527	Towards identifying fine-grained depression symp-	and Zhongyu Wei. 2024b. Somelvlm: A large vi-	1586
1528	toms from memes. In <i>Proceedings of the 61st Annual</i>	sion language model for social media processing. In	1587
1529	<i>Meeting of the Association for Computational Lin-</i>	<i>Findings of the Association for Computational Lin-</i>	1588
1530	<i>guistics (Volume 1: Long Papers)</i> , pages 8890–8905.	<i>guistics: ACL 2024</i> , pages 2366–2389.	1589
1531	Diyi Yang, Alon Lavie, Chris Dyer, and Eduard Hovy.	Zhengyi Zhao, Shubo Zhang, Yuxi Zhang, Yanxi Zhao,	1590
1532	2015. Humor recognition and humor anchor extrac-	Yifan Zhang, Zezhong Wang, Huimin Wang, Yutian	1591
1533	tion. In <i>Proceedings of the 2015 conference on empir-</i>	Zhao, Bin Liang, Yefeng Zheng, and 1 others. 2025.	1592
1534	<i>ical methods in natural language processing</i> , pages	Memereacon: Probing contextual meme understand-	1593
1535	2367–2376.	ing in large vision-language models. <i>arXiv preprint</i>	1594
		<i>arXiv:2505.17433</i> .	1595
1536	Yixin Yang, Zheng Li, Qingxiu Dong, Heming Xia, and	Shanshan Zhong, Zhongzhan Huang, Shanghua Gao,	1596
1537	Zhifang Sui. 2024. Can large multimodal models un-	Wushao Wen, Liang Lin, Marinka Zitnik, and Pan	1597
1538	cover deep semantics behind images? In <i>Findings of</i>	Zhou. 2024. Let’s think outside the box: Explor-	1598
1539	<i>the Association for Computational Linguistics: ACL</i>	ing leap-of-thought in large language models with	1599
1540	<i>2024</i> , pages 1898–1912, Bangkok, Thailand. Associ-	creative humor generation. In <i>Proceedings of the</i>	1600
1541	ation for Computational Linguistics.	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	1601
		<i>tern Recognition</i> , pages 13246–13257.	1602
1542	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing	Naitian Zhou, David Jurgens, and David Bamman. 2024.	1603
1543	Sun, Tong Xu, and Enhong Chen. 2024. A survey on	Social meme-ing: Measuring linguistic variation in	1604
1544	multimodal large language models. <i>National Science</i>	memes. In <i>Proceedings of the 2024 Conference of</i>	1605
1545	<i>Review</i> , 11(12):nwae403.	<i>the North American Chapter of the Association for</i>	1606
		<i>Computational Linguistics: Human Language Tech-</i>	1607
1546	Haofei Yu, Zhengyang Qi, Lawrence Keunho Jang, Russ	<i>nologies (Volume 1: Long Papers)</i> , pages 3005–3024.	1608
1547	Salakhutdinov, Louis-Philippe Morency, and Paul Pu		
1548	Liang. 2024. Mmoe: Enhancing multimodal mod-		
1549	els with mixtures of multimodal interaction experts.		
1550	In <i>Proceedings of the 2024 Conference on Empiri-</i>		
1551	<i>cal Methods in Natural Language Processing</i> , pages		
1552	10006–10030.		
1553	Haorui Yu, Yang Zhao, Yijia Chu, and Qiufeng Yi. 2025.		
1554	Seeing symbols, missing cultures: Probing vision-		
1555	language models’ reasoning on fire imagery and cul-		
1556	tural meaning. In <i>Proceedings of the 9th Widening</i>		
1557	<i>NLP Workshop</i> , pages 1–8.		
1558	Huixuan Zhang and Xiaojun Wan. 2024. Image matters:		
1559	A new dataset and empirical study for multimodal hy-		
1560	perbole detection. In <i>Proceedings of the 2024 Joint</i>		

A Overview of Datasets and Benchmarks

Here we provide a comprehensive tabular overview of datasets and benchmarks studied in this survey, organized according to the unified taxonomy and capability hierarchy introduced in the main paper.

Table 1: Overview of multimodal human creativity datasets focused on **Recognition** tasks. Recognition datasets typically use binary labels (e.g., humorous vs. non-humorous), element-level labels (e.g., persuasion techniques, punchlines, or rhetorical roles), and intensity scores (e.g., degree of funniness or offensiveness) to capture different granularities of creative signals. In Data Forms, StaVT denotes Static Visual–Textual Artifacts; SeqVN denotes Sequential Visual Narratives; PM denotes Audio–Visual Performative Media. In Mechanism, Multi denotes this dataset contain multi rhetorical mechanics we define in Sec. 2. The Availability column provides links to publicly accessible datasets and "N/A" indicate unpublished datasets.

Dataset	Venue	Data Forms	Mechanism	Communicative Goal	Size	Availability
D-HUMOR (Kasu et al., 2025)	ICDM'25	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	4,379	Link
MemeMind (Gu et al., 2025)	Arxiv'25	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	43,223	N/A
TOXICN MM (Lu et al., 2024)	NeurIPS'24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	12K	Link
PrideMM (Shah et al., 2024)	EMNLP'24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	5,063	Link
BHM (Hossain et al., 2024)	ACL'24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	7,148	N/A
Ext-Harm-P (Sharma et al., 2022a)	NAACL'23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	4,446	Link
HVVMemes (Sharma et al., 2023b)	EACL'23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	6,933	Link
RESTORE (Yadav et al., 2023)	ACL'23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	4,664	Link
Dank or not? (Barnes et al., 2020)	App. Net. Sci.'21	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	70K	N/A
The Hateful Memes Challenge Set (Kiela et al., 2020)	NeurIPS'20	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	10K	Link
Propaganda-techniques-in-memes (Dimitrov et al., 2021a)	ACL'21	StaVT: <i>Meme</i>	Multi	Persuasion & Advertising	950	Link
SemEval-2021 Task 6 (Dimitrov et al., 2021b)	ACL'21	StaVT: <i>Meme</i>	Multi	Persuasion & Advertising	950	Link
SemEval-2020 Task 8 (Sharma et al., 2020)	ACL'20	StaVT: <i>Meme</i>	Multi	Emotion & Aesthetic Experience	10K	N/A
RedEval (Tang et al., 2024)	NAACL'24	StaVT: <i>Sarcastic Image</i>	Sarcasm	Satire and Social Critique	1,004	Link
SPMSD (Guo et al., 2025)	COLING'24	StaVT: <i>Sarcastic Image</i>	Sarcasm	Satire and Social Critique	1K	N/A
MSTI dataset (Wang et al., 2022)	ACL'22	StaVT: <i>Sarcastic Image</i>	Sarcasm	Satire and Social Critique	5,015	Link
Multi-Modal Sarcasm Detection in Twitter (Cai et al., 2019)	ACL'19	StaVT: <i>Sarcastic Image</i>	Sarcasm	Satire and Social Critique	24,635	N/A
Sarcasm in Multimodal Social Platforms (Schifanella et al., 2016)	ACM MM'16	StaVT: <i>Sarcastic Image</i>	Sarcasm	Satire and Social Critique	10K	N/A
ArtELingo (Mohamed et al., 2024)	EMNLP'22	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	80k	Link
WikiArt Emotions (Mohammad and Kiritchenko, 2018)	LREC'18	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	4,105	Link
AVA (Murray et al., 2012)	CVPR'12	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	250K	Link
HumorDB (Jain et al., 2025)	ICCV'25	StaVT: <i>Humorous Image</i>	Multi	Humor & Entertainment	3,542	Link
Image Matters (Zhang and Wan, 2024)	COLING'24	StaVT: <i>Hyperbole image</i>	Hyperbole	Emotion & Aesthetic Experience	2,160	Link
Persuasion Strategies in Advertisements (Kumar et al., 2023)	AAAI'23	StaVT: <i>Poster</i>	Multi	Persuasion & Advertising	64,832	Link
DeepMoji (Felbo et al., 2017)	EMNLP'17	StaVT: <i>Emoji</i>	Multi	Satire and Social Critique, Emotion & Aesthetic Experience	1246M	Link
Inside Jokes (Shahaf et al., 2015)	KDD'15	StaVT: <i>Cartoon</i>	Multi	Humor & Entertainment	76,928	N/A

(Continued with Table 1)

Dataset	Venue	Data Forms	Mechanism	Communicative Goal	Size	Availability
COMICORDA (Martínek et al., 2024)	COLING'24	SeqVN: <i>Comic Strip</i>	Narrative	Humor & Entertainment	1,438	N/A
AVH & FOR (Chandrasekaran et al., 2016)	CVPR'16	SeqVN: <i>Comic Strip</i>	Narrative	Humor & Entertainment	7,150	Link
StandUp4AI (Barriere et al., 2025)	EMNLP-Findings'25	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	330.00 hours	Link
Multimodal Comic Mischief Dataset (Baharlouei et al., 2024)	COLING'24	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	74.63 hours	Link
DY11k (Liu et al., 2024a)	ICMR'24	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	35.38 hours	Link
SMILE (Hyun et al., 2024)	NAACL-Findings'24	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	14.41 hours	Link
Friends (Alnajjar et al., 2022)	COLING'22	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	78.60 hours	N/A
MUCH (Guo et al., 2024)	COLING'24	PM: <i>Video</i>	Multi	Humor & Entertainment	62.00 hours	Link
Multimodal and Multilingual Laughter Detection (Kuznetsova and Strapparava, 2024)	COLING'24	PM: <i>Video</i>	Multi	Humor & Entertainment	17.00 hours	Link
The Laughing Machine (Kayatani et al., 2021)	WACV'21	PM: <i>Video</i>	Multi	Humor & Entertainment	77.42 hours	N/A
MUSStARD (Castro et al., 2019)	ACL'20	PM: <i>Video</i>	Sarcasm	Satire and Social Critique	9.31 hours	Link
UR-FUNNY (Hasan et al., 2019b)	EMNLP'19	PM: <i>Video</i>	Multi	Humor & Entertainment	90.23 hours	Link

Table 2: A collection of datasets for multimodal creative **Understanding** and **Generation**. Understanding and generation datasets typically pair multimodal inputs with explanations or rationales—often augmented with contextual or external knowledge—to justify why creativity arises, alongside target creative outputs that support coherent and controllable generation across modalities (e.g., Memes, Comics or Videos + Caption/Rationales). In Data Forms, StaVT denotes Static Visual–Textual Artifacts; SeqVN denotes Sequential Visual Narratives; PM denotes Audio–Visual Performative Media. In Mechanism, Multi denotes this dataset contain multi rhetorical mechanics we define in Sec. 2. The Availability column provides links to publicly accessible datasets, and "N/A" indicates unpublished datasets.

Dataset	Venue	Data Forms	Mechanism	Communicative Goal	Size	Availability
MemeReaCon (Zhao et al., 2025)	EMNLP’25	StaVT: <i>Meme</i>	Multi	Satire & Social Critique	1,565	N/A
MEMESAFETY-BENCH (Lee et al., 2025)	EMNLP’25	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	50,430	Link
MemeMind (Gu et al., 2025)	Arxiv’25	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	43,223	N/A
MemeQA (Nguyen et al., 2025)	ACL’25	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	9K	Link
SEMANTICMEMES (Zhou et al., 2024)	NAACL’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment, Satire & Social Critique	3.8M	Link
MMD (Khan et al., 2024)	EMNLP-Findings’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	13,494	Link
MultiBully-Ex (Jha et al., 2024b)	EACL’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	5,854	Link
Oogiri-GO (Zhong et al., 2024)	CVPR’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	130K	Link
MemeMQACorpus (Agarwal et al., 2024)	ACL-Findings’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	1,880	N/A
ICMM (Jha et al., 2024a)	ACL’24	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	1K	Link
OxfordTVG-HIC (Li et al., 2023)	ICCV’23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	2.9M	Link
MEMECAP (Hwang and Shwartz, 2023)	EMNLP’23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	6,300	Link
MCC (Sharma et al., 2023c)	ACL’23	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	3,4K	Link
HatReD (Hee et al., 2023)	IJCAI’24	StaVT: <i>Meme</i>	Multi	Satire & Social Critique	3,228	Link
ExHVV (Sharma et al., 2023a)	AAAI’22	StaVT: <i>Meme</i>	Multi	Humor & Entertainment	3K	Link
HAIVMet (Chakrabarty et al., 2023)	ACL’23	StaVT: <i>Artwork</i>	Metaphor	Humor & Entertainment, Emotion & Aesthetic Experience	6,476	Link
ArtELingo (Mohamed et al., 2024)	EMNLP’22	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	80k	Link
LAION-Aesthetics (Schuhmann et al., 2022)	NeurIPS’22	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	51.9M	Link
ArtEmis (Achlioptas et al., 2021)	CoRR’21	StaVT: <i>Artwork</i>	Multi	Emotion & Aesthetic Experience	81,446	Link
HumorBench (Narad et al., 2025)	Arxiv’25	StaVT: <i>Cartoon</i>	Multi	Humor & Entertainment	300	N/A
Humor in AI (Zhang et al., 2024a)	NeurIPS’24	StaVT: <i>Cartoon</i>	Incongruity	Humor & Entertainment	2.2M	Link
Do Androids Laugh at Electric Sheep? (Hessel et al., 2023b)	ACL’23	StaVT: <i>Cartoon</i>	Multi	Humor & Entertainment	24,048	N/A
V-FLUTE (Saakyan et al., 2025)	NAACL’25	StaVT: <i>Social Media Image</i>	Multi	Humor & Entertainment, Satire & Social Critique	6,027	Link
SoMeLVLM (Zhang et al., 2024b)	ACL’24	StaVT: <i>Social Media Image</i>	Multi	Humor & Entertainment, Satire & Social Critique, Emotion & Aesthetic Experience	653.8K	Link
HumorDB (Jain et al., 2025)	ICCV’25	StaVT: <i>Humorous Image</i>	Multi	Humor & Entertainment	3,542	Link
VisualPun_UNPIE (Chung et al., 2024)	EMNLP’24	StaVT: <i>Humorous Image</i>	Incongruity	Humor & Entertainment	1K	Link

(Continued with Table 2)

Dataset	Venue	Data Forms	Mechanism	Communicative Goal	Size	Availability
ImageMet (Kundu et al., 2025)	EMNLP-Findings'25	StaVT: <i>Metaphor Image</i>	Metaphor	Emotion & Aesthetic Experience	2,527	N/A
MangaUB (Ikuta et al., 2025)	IEEE MM'25	SeqVN: <i>Comic Strip</i>	Narrative	Humor & Entertainment	18,179	Link
AI4VA-FG (Chen et al., 2025)	Arxiv'25	SeqVN: <i>Comic Strip</i>	Narrative	Humor & Entertainment	16,264	N/A
PixelHumor (Ryan et al., 2025)	EMNLP-Findings'25	SeqVN: <i>Comic Strip</i>	Multi	Humor & Entertainment	2.8K	Link
YesBut-v2 (Liang et al., 2025)	Arxiv'25	SeqVN: <i>Comic Strip</i>	Multi	Humor & Entertainment	1,262	Link
YesBut (Hu et al., 2024)	NeurIPS'24	SeqVN: <i>Comic Strip</i>	Multi	Humor & Entertainment	348	Link
YesBut (synthetic 3D stick) (Nandy et al., 2024)	EMNLP'24	SeqVN: <i>Comic Strip</i>	Multi	Humor & Entertainment	2,547 (synthetic)	Link
ComicsPAP (Vivoli et al., 2025)	Arxiv'25	SeqVN: <i>Comic Strip</i>	Narrative	Humor & Entertainment	103,933	Link
XMeCap (Chen et al., 2024a)	ACM MM'24	SeqVN: <i>Meme</i>	Multi	Humor & Entertainment	12,320	N/A
Laugh, Relate, Engage (Ouyang et al., 2025)	Arxiv'25	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	1K	N/A
V-HUB (Shi et al., 2025b)	Arxiv'25	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	4.00 hours	N/A
EmoVid (Qiu et al., 2025)	Arxiv'25	PM: <i>Short-form Video</i>	Symbolism	Emotion & Aesthetic Experience	39.00 hours	N/A
ExFunTube (Ko et al., 2023)	EMNLP'23	PM: <i>Short-form Video</i>	Multi	Humor & Entertainment	84.00 hours	Link
MUStARD (Castro et al., 2019)	ACL'20	PM: <i>Short-form Video</i>	Sarcasm	Satire & Social Critique	9.31 hours	Link
AdsVQA (Long et al., 2025)	ICCV'25	PM: <i>Ads Video</i>	Multi	Persuasion & Advertising	1,544 (22.70 hours)	Link
LAMDBA (Harini et al., 2025)	WACV'25	PM: <i>Ads Video</i>	Multi	Persuasion & Advertising	2,205	Link
TV Ads (Tao et al., 2024)	Multimedia Tools App.'24	PM: <i>Ads Video</i>	Multi	Persuasion & Advertising	14,490	N/A
MM-AU (Bose et al., 2023)	ACM MM'23	PM: <i>Ads Video</i>	Multi	Persuasion & Advertising	147.00 hours	N/A
Automatic Understanding of Image and Video Advertisements (Hussain et al., 2017)	CVPR'17	StaVT & PM: <i>Ads Poster & Video</i>	Multi	Persuasion & Advertising	64,832 & 3,477	Link