

---

# Is AI fun? HumorDB: a curated dataset and benchmark to investigate graphical humor

---

**Vedaant Jain**

University of Illinois Urbana-Champaign  
vvjain3@illinois.edu

**Felipe dos Santos Alves Feitosa**

University of São Paulo  
felipefeitosa@usp.br

**Gabriel Kreiman**

Children’s Hospital, Harvard Medical School  
gabriel.kreiman@tch.harvard.edu

## Abstract

Despite significant advancements in computer vision, understanding complex scenes, particularly those involving humor, remains a substantial challenge. This paper introduces HumorDB, a novel image-only dataset specifically designed to advance visual humor understanding. HumorDB consists of meticulously curated image pairs with contrasting humor ratings, emphasizing subtle visual cues that trigger humor and mitigating potential biases. The dataset enables evaluation through binary classification (Funny or Not Funny), range regression (funniness on a scale from 1 to 10), and pairwise comparison tasks (Which Image is Funnier?), effectively capturing the subjective nature of humor perception. Initial experiments reveal that while vision-only models struggle, vision-language models, particularly those leveraging large language models, show promising results. HumorDB also shows potential as a valuable zero-shot benchmark for powerful large multimodal models. We open-source both the dataset and code under the CC BY 4.0 license.

## 1 Introduction

The last decade has seen remarkable strides in the ability to label images, segment objects, write captions, and detect objects in complex scenes. Despite these successes, the problem of scene understanding remains challenging. Understanding a scene often requires interpreting the relationship between objects and their positions, the intention of agents, and linking visual information with prior knowledge. As a paradigmatic example of scene understanding, we focus here on the ability to assess whether an image is funny or not. Graphical humor understanding demands a high level of cognitive abstraction, as it requires context awareness and the identification of incongruities [4]. Consider **Fig. 1**, left. The viewer needs to detect locations (surgical setting), agents (patient, medical practitioners), and objects (cell phone, hand). The grasping in the hand induces us to think that the medical experts have “excised” the phone from the hand. Given the prominent role of cell phones in our culture, many people have jokingly stated that the phone is physically attached to the hand, and the image plays with this idea. Of note, most likely, the reader has never seen this particular image or any similar image before. Upon first exposure to this image, readers can rapidly interpret what is going on and may consider the image to be somewhat humorous (83.3% of participants indicated that this image is funny). In stark contrast, consider **Fig. 1**, right. The two images are identical except that the cell phone was removed. Despite the strong similarity between the two images, the one on the right is no longer humorous (85.7% of participants indicated that the image on the right is not funny).

To systematically evaluate how well current machine learning algorithms can assess graphical humor, here we introduce a comprehensive and controlled image dataset, **HumorDB**, and accompanying human evaluations based on three different metrics. Inspired by previous work in action recognition



Figure 1: **Example image pair.** Left: image rated as funny (83.3% of participants). Right: modified image rated as not funny (85.7 % indicated not funny). Focus on the phone in the surgeon’s hand in the left image.

(e.g., [17]), we mitigate biases by using slightly differing images and meticulous curation to avoid potential confounding factors, incorporating pairs of similar images with contrasting humor ratings, thereby honing the focus on the nuanced modifications that elicit humor (e.g., **Fig. 1**). HumorDB further includes English descriptions for a portion of its images, providing valuable context for training and evaluating vision-language models. We anticipate that HumorDB will serve as a valuable benchmark for developing intelligent systems capable of not only recognizing visual elements, but also understanding their nuanced relationships and interpreting them within a humorous context, ultimately pushing the boundaries of scene understanding.

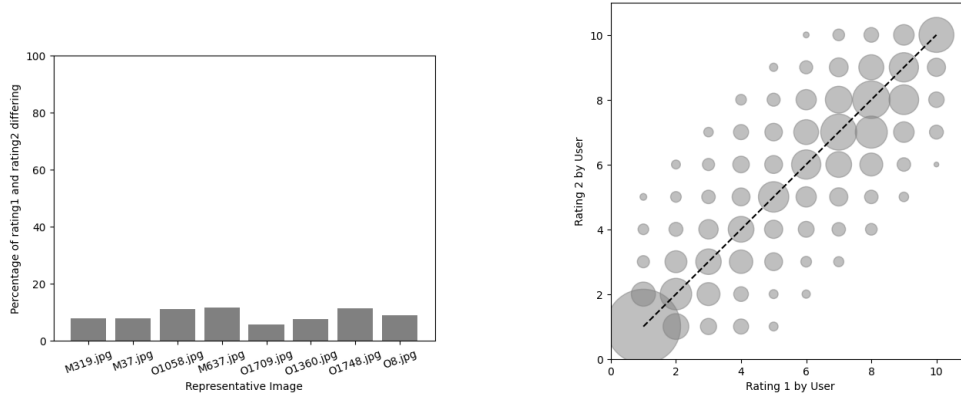
We start by outlining our methodology for data collection, including the creation of visually similar images that vary in humor ratings and the quantification of the subjective nature of humor through crowdsourcing. Next, we introduce three evaluation tasks—Binary Classification (Funny or Not Funny) and Regression (Funniness rated on a 1-10 scale), and Humor Comparison (which of two images is funnier?). Finally, we apply state-of-the-art computer vision and vision-language models as baselines to demonstrate the challenges that humor understanding poses to current technologies. We open-source the dataset and the code: [https://anonymous.4open.science/r/HumorDB\\_-51AF/](https://anonymous.4open.science/r/HumorDB_-51AF/)

## 2 Related Work

In cognitive psychology, “The best known theory of humor is probably the incongruity theory” [3]. This theory posits that humor arises from the unexpected turn of events that subverts existing expectations, which are constructed through contextual clues. The intricacy of modeling humor stems from its contextual dependencies and creative complexities, as humor often draws from the well of human creativity [14]. While recent efforts have been made to create computational humor understanding datasets [6, 13] and methods [5, 12], the majority of these projects focus on multi-modal data where inferring contextual cues might be easier from one modality rather than another. Our work aims to focus on images, presenting a unique challenge for the computer vision community.

Datasets for humor understanding have been emerging and have spanned a range of tasks: humor detection in multi-modal context [6, 13, 34], funny captioning of images [22], matching captions to cartoons [16], or explaining why a text joke is funny with a paired image [16]. These works have enabled development of powerful techniques to attempt humor understanding [23] and have triggered an increased interest in humor understanding [9, 35]. The approaches to humor understanding have involved using attention-based Bi-LSTM for identifying humor in text on social media [19], generation of negative examples [40], the use of transformers [12], and CRF-RNN-CNN [5].

Most of the humor datasets and approaches have focused on either natural language [7, 8, 19, 30] or multimodal contexts (especially videos [18, 27, 34]). To the best of our knowledge, there is a lack of datasets that concentrate only on images which are an important form of communication and may also help in developing more visually intelligent systems in the future.



(a) Participants showed high self-reliability (Humor Comparison). For each one of 10 images (x-axis), the y-axis indicates the percentage of times the second rating from the same user differed from the first rating. A value of 0 would indicate perfect self-reliability.

(b) Participants showed high self-reliability (Range Ratings). Repeated images were used to assess consistency, revealing a strong correlation ( $\rho = 0.89$ ) between first and second ratings (1,800 rating pairs; circle sizes indicate the number of ratings per pair).

Figure 2: Comparison of participant self-reliability in different contexts.

### 3 Building HumorDB

#### 3.1 Creating the Image repository

We gathered a diverse set of images from a variety of sources, including online image repositories, comics, social media platforms, and text-to-image generators like DALL-E [31] and MidJourney [29] (Table 1). This multi-source approach ensured a broad spectrum of comedic styles and content. We filtered the images to maintain diversity, eliminate potentially offensive content, and remove cases where humor depended on text within the image. We preserved text that did not contribute to humor, such as artist names, to maintain image authenticity. A key aspect of HumorDB is the creation of slightly modified image pairs designed to elicit contrasting humor responses (e.g., Fig. 1). This critical control addresses the challenge of models learning spurious correlations from biases in online images. These pairs, created using traditional and AI-powered image manipulation techniques, maintain similar visual structure while exhibiting different humor levels, highlighting the need for models to understand humor nuances rather than relying on spurious correlations. To mitigate potential biases introduced by the modifications themselves, subtle alterations were also applied to the original images, preserving their humor and preventing algorithms from simply identifying manipulation discrepancies.

#### 3.2 Human Evaluations

We conducted online psychophysics experiments with 550 participants through Amazon Mechanical Turk (MTurk) and Prolific to gather human evaluations of our curated image sets. Each participant provided ratings for 100 images without any explicit definition of "humor" in three different tasks:

**Binary Rating.** Participants classified each image as either humorous or not humorous.

**Range Rating.** Participants rated images on a scale from 1 (not at all humorous) to 10 (extremely humorous). Participants showed high self-reliability in their ratings (Fig. 2b) and rated original images as funnier than their modified counterparts (Fig. 3).

**Humor Comparison.** Participants were presented with two images at a time and had to indicate which was funnier and write a word/phrase about the funnier image. To assess the consistency of human ratings across participants, we selected 10 images and considered all  $\binom{10}{2}$  pairwise combinations. The results revealed a noticeable level of agreement across human raters (Fig. 4).

To address the impracticality of obtaining pairwise comparisons for every image, we implemented a stratified sampling approach. We first grouped images into eight strata based on their mean binary ratings. From each stratum, five representative images were selected, totaling 40 images. These

images underwent an additional round of binary ratings (40 ratings) to confirm their categorization. Finally, one image from each stratum was selected, resulting in a subset of eight images against which all other images were compared. With respect to the eight selected images, the ratings showed noticeable agreement within human raters (**Fig. 2a**).

**Controls** We took several precautions to ensure the reliability and consistency of the ratings.

(1) We required participants to spend at least 500 milliseconds before submitting a response for an image. Violating this requirement four times resulted in task termination.

(2) Unbeknownst to the participants, we repeated a randomly selected 10% of the images to assess self-consistency. In the Binary and Comparison tasks, if a participant’s ratings for  $> 3$  of the replicated images diverged, the ratings garnered from that participant were considered incongruous and discarded. For the range task, if a participant’s ratings displayed deviations exceeding 4 for any of the replicated images or witnessed variations exceeding 2 on  $> 4$  instances, the ratings were excluded from the analyses. We observed a high correlation ( $\rho = 0.89$ ) between the ratings a user gave to the repeated images in Range task (**Fig. 2b**). Subjects were also self-consistent in the Binary and Comparison tasks. In Binary Task the participants were consistent for an average of  $84.2 \pm 13.3\%$  for the repeated images and in comparison task, participants showed self-reliability at an average of  $91.3 \pm 14.8\%$  for the repeated images. Range Ratings for modified images were mostly lower than those for the original images (**Fig. 3**).

(3) We worried that labels coming from online participants may be unreliable. To address this potential concern, we collected psychophysics data (both Binary and Range ratings) from in-lab participants who were expected to provide reliable evaluations for a randomly selected subset of 200 images. For the Binary ratings, the degree of consistency between in-lab and online participants was 0.75. The Pearson correlation coefficient between the Range ratings of the in-lab participants and the online participants was 0.68.

(4) Humor is in the eye of the beholder, and there can be variability across participants. For the current effort, we focused on the *common* judgments across participants. To avoid large subjective deviations from individual participants, we implemented a criterion for identifying outliers in humor assessment. We regarded a participant’s rating as an outlier if it deviated significantly from the collective average. If the z-score of a participant’s rating fell outside the range of -1.96 to 1.96 (corresponding to a 95% confidence interval), that rating was considered an outlier. This approach, which focuses on the central tendency of the data, helped exclude extreme ratings that do not represent the common perception of humor. Data from participants whose responses met this outlier criterion for  $> 10\%$  of their responses were removed from our analyses.

(5) Participants also provided textual descriptions highlighting humorous elements in the funnier image, ensuring attentive viewing and enabling identification of common humor triggers. Using these

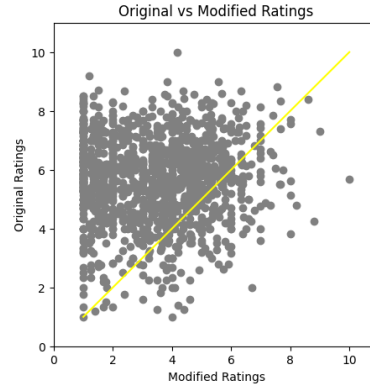


Figure 3: **Modifications rendered images less humorous.** Each point compares the rating of image pairs (y-axis: original, x-axis: modified pair; total 1,271 pairs; line = identity). For the majority of images (86.4%), the ratings for the original images were higher.

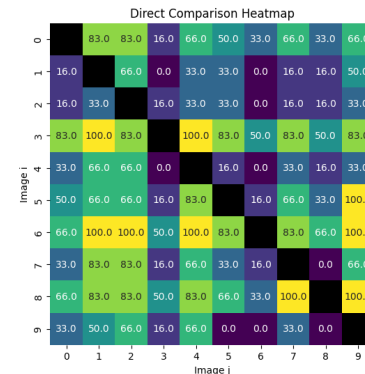


Figure 4: Each cell(i, j) represents the percentage of times when image i was rated funnier than image j. Users tend to agree which image is funnier, showing a dominance of images 6 and 3 being rated more funny than others most times, indicating contesting humor levels.

Dataset	Total	Pairs	Funny images			Not Funny		
			N	Binary	Range	N	Binary	Range
Training set	2,136	698	1,068	0.79 ± 0.18	5.75 ± 1.42	1,068	0.12 ± 0.16	3.62 ± 1.75
Validation set	703	273	351	0.78 ± 0.19	5.68 ± 1.40	351	0.12 ± 0.17	3.65 ± 1.65
testAll Set	706	300	352	0.77 ± 0.19	5.60 ± 1.37	352	0.14 ± 0.17	3.39 ± 1.68
test Only-Pairs	600	300	300	0.77 ± 0.19	5.62 ± 1.36	300	0.13 ± 0.16	3.30 ± 1.68
<b>Total</b>	<b>3,545</b>	<b>1,271</b>	<b>1,771</b>	<b>0.79 ± 0.20</b>	<b>5.70 ± 1.40</b>	<b>1,771</b>	<b>0.13 ± 0.17</b>	<b>3.58 ± 1.72</b>

Table 1: Dataset-Summary. Total is without the testOnlyPairs as it is a subset of testAllSet

textual descriptions, we identified common words or phrases that appeared in at least 30% of the responses for each image, aiding in the identification of common humorous features.

### 3.3 Final Dataset

For the Binary task, we labeled an image as humorous if the mean rating across participants was  $\geq 0.5$  (Table 1). To ensure a balanced dataset, given that we had more non-funny images, we removed those non-funny images that showed the highest standard deviation in the ratings across participants to end up with an equal number of funny and non-funny images. To separate the data into training, validation, and test splits, we ensured that each training/validation/test set contained both the original and modified versions for a given image, further avoiding potential biases across images. We conducted two separate evaluations: one on the entire test set (testAllSet) and the other on the subset of the test set consisting of only the original/modified pairs of images in the test set (testOnlyPairs).

## 4 Computational experiments

### 4.1 Models and Tasks

We evaluated state-of-the-art visual architectures, including vision-only and vision-language models, using both pretrained and trained-from-scratch settings. Pretrained vision models were pretrained with either ImageNet [36] or LAION-2B [37]. BLIP and LLaVA on the other hand use a combination of various dataset and non-trivial data-filtering methods. GPT-4o and Gemini-Flash do not open-source their training data.

Each model was trained in two ways: (i) either using the pretrained weights and fine-tuning, or (ii) training with random initialization of weights. We report results for each of these settings. We report results for the testAllSet and the testOnlyPairs.

For both the Binary and Range tasks with vision-only models, we added a final single fully connected layer at the end of each model to make predictions. For vision-only models in the Comparison task, we created a network with a single backbone that extracts features from before the classification layer. We then computed the difference between the feature values of the two images and passed this difference through a fully connected layer to classify which image was funnier. We tested two schemes: a canonical fully connected layer and one incorporating dropout; however, we did not observe significant differences between the two settings.

To evaluate the vision-language models, we framed our task as a visual question answering (VQA) problem with two variants of the question to avoid biases: (i) "Is the image funny?", and (ii) "Is the image not funny?". We report the average accuracy achieved for both questions, as performance was comparable between the two. For the Range rating, our prompt was "What is the degree of funniness of this image in the range from 1 (not funny) to 10 (extremely funny)?" For the Comparison task, we used vision-language models in a zero-shot setting with the prompt: "Given these two images, answer which is funnier by responding with 'first' or 'second' and then explain succinctly." Additionally,

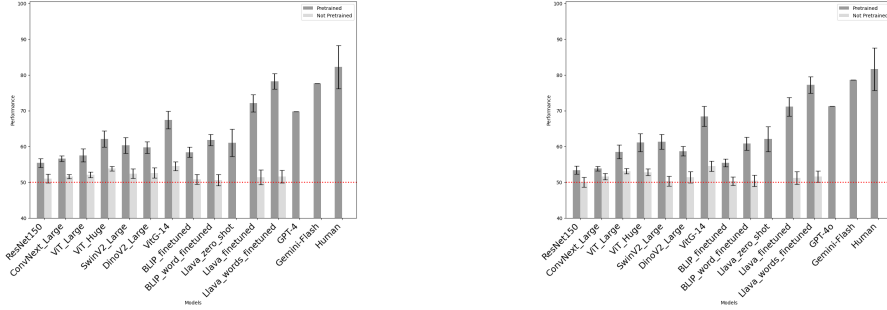


Figure 5: **Binary task results for (a) testAllSet and (b) testOnlyPairs.** The dark grey bars show the results for pretrained models whereas the light grey bars show the results from models that are not pretrained. The dotted line is random chance. The last column shows human performance. The references for models are in **Table 2**.

as we collect words about funnier image in the comparison task from users, we were able to select common words from the responses for each image. Our criteria for common words is that the word appears atleast 30% of times in the responses about the image. Then for fine-tuning with words we modify the prompt by using a prefix: “The prominent features of image are: {common words[image]}.”.

Due to the uneven distribution of range ratings, we employed a sampling strategy that grouped images into bins according to their funniness ratings. This allowed us to randomly select a balanced number of images from each bin for every training epoch, ensuring a uniform distribution of sample images across all ratings in the training set. Further training details are shown in Appendix A.1. Comparison ratings also contained slightly uneven distribution; therefore, we applied the same sampling strategy as above for training. For the Comparison task, we discarded comparisons where multiple users gave conflicting ratings and there was no majority.

## 5 Results

### 5.1 Binary task

Despite the complexity of the task, surprisingly, multiple models achieved above chance performance (**Fig.5**). As expected, all the pretrained models (dark gray) achieved higher performance than the non-pretrained models (light gray). Non-pretrained models were near chance performance, emphasizing the crucial role of pretraining. Performance was slightly higher in the testAllSet (**Fig.5** left), compared to the testOnlyPairs set (**Fig.5** right), reflecting some of the intrinsic biases that are at least partially controlled for in the modified images. Large models (LLaVA, GPT-4o, ViTG-14, Gemini-Flash) exhibit consistent performance on both the testAllSet and testOnlyPairs datasets, indicating robustness to image modifications. Conversely, smaller models (with less than a billion parameters) experience a drop in performance on testOnlyPairs, highlighting the importance of these pairs for stress-testing models and uncovering potential biases stemming from pre-trained data.

The best performance was reached for the Gemini-Flash and Llava\_words\_finetuned models. The results demonstrate the superiority of vision-language models (Llava, GPT-4, Gemini-Flash) over vision-only models (ViT\_Huge, SwinV2\_Large, DinoV2\_Large). Notably, vision-language models(LLaVA, BLIP) trained with supporting words for a portion of the images further enhanced performance. Zero-shot evaluations using LLM backbones like LLaVA, GPT-4o, and Gemini-Flash reveal promising results, in some cases reaching almost the consistency of Human responses(Gemini-Flash), suggesting HumorDB as a valuable benchmark for assessing these powerful models on a challenging task.

### 5.2 Range task

The models also performed surprisingly well in the Range task (**Table 2**). Even models that were only evaluated in a zero-shot fashion like GPT-4o or Gemini-Flash showed adequate performance.

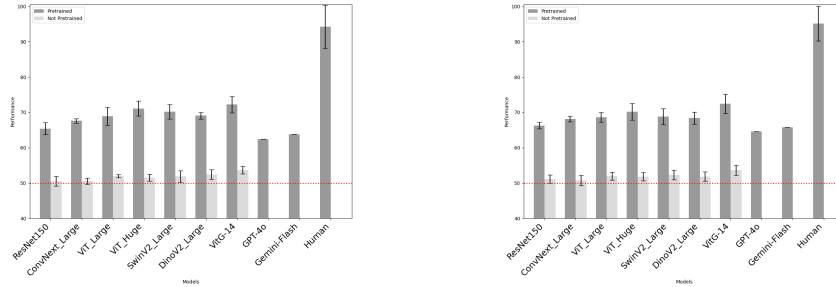


Figure 6: **Comparison task results for (a) testAllSet and (b) testOnlyPairs.** The dark grey bars show the results for pretrained models whereas the light grey bars show the results from models that are not pretrained

Model Name	testAllSet RMSE $\pm SD$	testOnlyPairs RMSE $\pm SD$
dinov2 large [33]	1.98 $\pm$ 0.08	1.96 $\pm$ 0.09
vit huge [10]	1.88 $\pm$ 0.09	1.92 $\pm$ 0.07
swin2 large [26]	1.96 $\pm$ 0.12	1.98 $\pm$ 0.11
convnext large [28]	2.10 $\pm$ 0.06	2.15 $\pm$ 0.07
vitg 14 [20, 38]	1.76 $\pm$ 0.08	1.80 $\pm$ 0.06
resnet152 [15]	2.11 $\pm$ 0.09	2.09 $\pm$ 0.11
LLaVA(Zero-Shot) [24, 25]	2.95 $\pm$ 0.33	2.98 $\pm$ 0.34
LLaVA(fine-tuned) [24, 25]	1.70 $\pm$ 0.22	1.69 $\pm$ 0.22
LLaVA(words fine-tuned) [24, 25]	<b>1.68 <math>\pm</math> 0.28</b>	<b>1.66 <math>\pm</math> 0.31</b>
BLIP (fine-tuned) [21]	1.94 $\pm$ 0.05	1.96 $\pm$ 0.05
BLIP (words fine-tuned) [21]	1.92 $\pm$ 0.06	1.95 $\pm$ 0.04
GPT-4o(Zero-Shot) [32]	2.61	2.63
Gemini-Flash [39]	2.06	2.11
Humans	2.72 $\pm$ 0.88	2.71 $\pm$ 0.86
Chance(from distribution)	2.69 $\pm$ 0.05	2.50 $\pm$ 0.06
Chance(for zero-shot)	3.58 $\pm$ 0.07	3.33 $\pm$ 0.07

Table 2: Range Regression Results on testAllSet and testOnlyPairs. Bold indicates the best model result. The chance(from distribution) scores are calculated by random sampling from the distribution of ratings in the dataset after rounding the mean range ratings to the nearest integer. The chance(zero-shot) scores are calculated after randomly sampling an integer from 1 to 10.

### 5.3 Comparison task

This task revealed larger differences between human and model performance (**Fig. 6**). As expected and as reported for the Binary task, pretrained models (dark grey) performed better than non-pretrained models (light gray), and performance was slightly higher in the testAllSet task. Although most of the models performed significantly better than chance, all models were still rather far from human performance. Here zero-shot multimodal LLMs performed worse than fine-tuned vision-only models. We do not report results for BLIP and LLaVA as they do not support multi-image inputs in a way that lets them compare images.

### 5.4 Modified Images are necessary to stress test models

Experiments revealed a significant performance increase for several models (Vit-huge, VitG-14, ResNet150, GPT-4o) when evaluated only on the test set without modified images for the binary task. Notably, GPT-4o shows a surprising 10% improvement.(See Appendix A.6) This observation reinforces the necessity of modified images for evaluating large models and preventing biases caused by already-seen data.

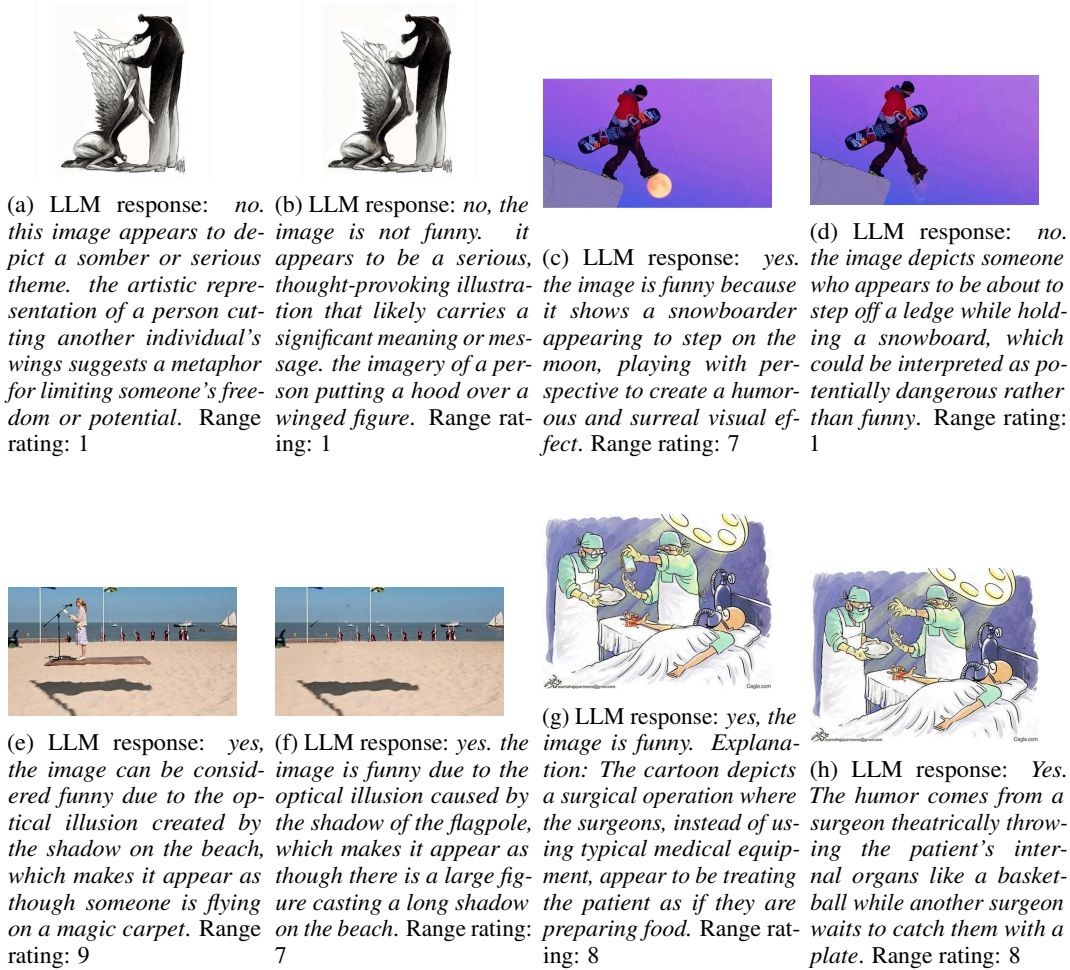


Figure 7: **Example GPT-4 binary classification, range ratings, and explanations.** The human binary and range ratings are as follows for the images: (a) 66%, 4.20, (b) 0%, 5.6, (c) 71%, 6.70, (d) 0%, 3.60, (e) 60%, 2.80, (f) 0%, 5.57, (g) 83%, 6.75, (h) 14%, 5.70

## 5.5 Multimodal LLM answer explanations

We submitted the images to GPT-4o using prompts to obtain binary classification, and also an explanation (Fig. 7). Several examples were correctly classified by GPT-4o (e.g., Figure 7b, c, d, e, g) and many others were incorrectly classified (e.g., 7a, f, h). For several examples, we show both the original image and its modified counterpart to further illustrate the importance of paired images in this dataset.

Fig. 7g-h follow the example in Fig. 1. GPT-4o labeled both images as funny. Examination of the explanation provided shows that GPT-4o failed to grasp the elements critical for humor. A similar situation is seen in Fig. 7e-f. Fig. 7a shows a case where the model did recognize the cutting of feathers as something odd; however, it failed to connect that humor. Fig. 7c-d shows a case where the model properly classified the image pairs. The explanation is remarkably eloquent and reasonable.

These examples illustrate the diversity of humorous images in the dataset and the facets of humor the models must comprehend to accurately classify, rank, and explain the images. More examples are provided in the supplement with answer explanations from Gemini as well(See A.5).



## 6 Discussion

For the Binary task, all the models evaluated here exhibited a significantly lower performance compared to humans, despite being fine-tuned on average ratings. Humans rated images in a “zero-shot” manner, bringing their inherent knowledge into play. For the Range task, the models reached a performance that was comparable to humans. For the Comparison task, the models performed better than chance but worse than humans. As in most other vision/language tasks, it is extremely difficult to compare the training regimes of humans and machine algorithms.

As the field moves from labeling individual objects to scene interpretation, we will have to contend with human variation in labels. The vast majority of people may label a picture of a chair as a chair, but not everybody will find the same image equally funny. Here we focus on the majority votes (Binary and Comparison tasks), or average scores (Range task). In the future, it will be interesting to assess whether models can capture individual assessments of humor.

A typical concern in machine learning is the amount of training data. The dataset size is slightly smaller but of the same order of magnitude in terms of images per class compared to other benchmarks in computer vision. For example, imageNet [36] has approximately 1,000 images per class, which is similar to HumorDB (**Table 1**). It is conceivable that if one were able to expand the dataset by several orders of magnitude, models could do better. It is not trivial to expand HumorDB by orders of magnitude given the multiple controls, image modifications, and curation involved plus the paucity of humorous images. However, we speculate that the main reason underlying the discrepancy between current models and human performance in the binary and the comparison task is *not* the dataset size. Rather, these differences underscore the current limitations of models in terms of effectively reasoning about abstract concepts. This is evident in the GPT-4o answer explanations (**Fig. 7**, Gemini answer explanations (Appendix 11) as well as attention maps shown in Appendix A.3.

The construction of the dataset highlights the need for controls in computer vision and how easy it is to fall into biases that can help classification and artificially inflate performance. In particular, in HumorDB, the nuanced modifications made to the images to create similarly looking images belonging to different classes (**Fig. 1**) go a long way towards avoiding such biases.

HumorDB aims to begin to illustrate our lacunae in scene understanding and provide a dataset and benchmark for future efforts to deepen our understanding of abstract concepts such as humor. Developing AI systems that truly capture human humor could have multiple beneficial applications (perhaps entertainment, therapy, and furthering our understanding of human cognition). There is also a risk that AI systems might misuse humor, potentially leading to offensive outputs and potentially relying on subtle visual cues that may unintentionally perpetuate stereotypes or biases. At the same time, improving models’ understanding of abstract concepts like humor may lead to improved systems of content moderation.

## References

- [1] S. Abnar and W. Zuidema. *Quantifying Attention Flow in Transformers*. 2020. arXiv: 2005.00928 [cs.LG].
- [2] J. Ansel et al. “PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation”. In: ASPLOS ’24. <conf-loc>, <city>La Jolla</city>, <state>CA</state>, <country>USA</country>, </conf-loc>: Association for Computing Machinery, 2024, pp. 929–947. DOI: 10.1145/3620665.3640366.
- [3] S. Attardo. “57Theories of humor and their levels”. In: *The Linguistics of Humor: An Introduction*. Oxford University Press, June 2020. ISBN: 9780198791270. DOI: 10.1093/oso/9780198791270.003.0003. eprint: <https://academic.oup.com/book/0/chapter/349333277/chapter-pdf/43339517/oso-9780198791270-chapter-3.pdf>.
- [4] S. Aykan and E. Nalçacı. “Assessing Theory of Mind by Humor: The Humor Comprehension and Appreciation Test (ToM-HCAT)”. In: *Frontiers in Psychology* 9 (2018). ISSN: 1664-1078. DOI: 10.3389/fpsyg.2018.01470.
- [5] D. Bertero and P. Fung. “Deep Learning of Audio and Language Features for Humor Prediction”. In: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*. Portorož, Slovenia: European Language Resources Association (ELRA), May 2016, pp. 496–501.
- [6] D. S. Chauhan et al. *M2H2: A Multimodal Multiparty Hindi Dataset For Humor Recognition in Conversations*. 2021. arXiv: 2108.01260 [cs.CL].
- [7] L. Chen and C. M. Lee. “Predicting Audience’s Laughter During Presentations Using Convolutional Neural Network”. In: *Proceedings of the 12th Workshop on Innovative Use of NLP for Building Educational Applications*. Copenhagen, Denmark: Association for Computational Linguistics, Sept. 2017, pp. 86–90. DOI: 10.18653/v1/W17-5009.
- [8] P.-Y. Chen and V.-W. Soo. “Humor Recognition Using Deep Learning”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, June 2018, pp. 113–117. DOI: 10.18653/v1/N18-2018.
- [9] L. Christ et al. “The MuSe 2022 Multimodal Sentiment Analysis Challenge: Humor, Emotional Reactions, and Stress”. In: *MuSe’ 22*. Lisboa, Portugal: Association for Computing Machinery, 2022, pp. 5–14. ISBN: 9781450394840. DOI: 10.1145/3551876.3554817.
- [10] M. Dehghani et al. *Scaling Vision Transformers to 22 Billion Parameters*. 2023.
- [11] S. Gugger et al. *Accelerate: Training and inference at scale made simple, efficient and adaptable*. <https://github.com/huggingface/accelerate>. 2022.
- [12] M. K. Hasan et al. “Humor Knowledge Enriched Transformer for Understanding Multimodal Humor”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 35.14 (May 2021), pp. 12972–12980. DOI: 10.1609/aaai.v35i14.17534.
- [13] M. K. Hasan et al. “UR-FUNNY: A Multimodal Language Dataset for Understanding Humor”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, 2019. DOI: 10.18653/v1/d19-1211.
- [14] W. E. Hauck and J. W. Thomas. “The Relationship of Humor to Intelligence, Creativity, and Intentional and Incidental Learning”. In: *The Journal of Experimental Education* 40.4 (1972), pp. 52–55. ISSN: 00220973, 19400683.
- [15] K. He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs.CV].
- [16] J. Hessel et al. *Do Androids Laugh at Electric Sheep? Humor "Understanding" Benchmarks from The New Yorker Caption Contest*. 2023. arXiv: 2209.06293 [cs.CL].
- [17] V. Jacquot, Z. Ying, and G. Kreiman. “Can Deep Learning Recognize Subtle Human Activities?” In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2020.
- [18] Y. Kayatani et al. “The Laughing Machine: Predicting Humor in Video”. In: *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2021, pp. 2072–2081. DOI: 10.1109/WACV48630.2021.00212.

- [19] D. Li et al. “HEMOS: A novel deep learning-based fine-grained humor detecting method for sentiment analysis of social media”. In: *Information Processing & Management* 57.6 (2020), p. 102290. ISSN: 0306-4573. DOI: <https://doi.org/10.1016/j.ipm.2020.102290>.
- [20] J. Li et al. *BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models*. 2023. arXiv: 2301.12597 [cs.CV].
- [21] J. Li et al. *BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation*. 2022. arXiv: 2201.12086 [cs.CV].
- [22] R. Li et al. *OxfordTVG-HIC: Can Machine Make Humorous Captions from Images?* 2023. arXiv: 2307.11636 [cs.CV].
- [23] P. P. Liang et al. *Multimodal Learning Without Labeled Multimodal Data: Guarantees and Applications*. 2023. arXiv: 2306.04539 [cs.LG].
- [24] H. Liu et al. *Improved Baselines with Visual Instruction Tuning*. 2023.
- [25] H. Liu et al. *Visual Instruction Tuning*. 2023.
- [26] Z. Liu et al. *Swin Transformer V2: Scaling Up Capacity and Resolution*. 2022. arXiv: 2111.09883 [cs.CV].
- [27] Z.-S. Liu, R. Courant, and V. Kalogeiton. “FunnyNet: Audiovisual Learning of Funny Moments in Videos”. In: *Computer Vision – ACCV 2022*. Ed. by L. Wang et al. Cham: Springer Nature Switzerland, 2023, pp. 433–450.
- [28] Z. Liu et al. *A ConvNet for the 2020s*. 2022. arXiv: 2201.03545 [cs.CV].
- [29] Midjourney Inc. *Midjourney*. <https://www.midjourney.com/>. 2023.
- [30] R. Mihalcea and C. Strapparava. “Making Computers Laugh: Investigations in Automatic Humor Recognition”. In: *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*. Vancouver, British Columbia, Canada: Association for Computational Linguistics, Oct. 2005, pp. 531–538.
- [31] OpenAI. *DALL-E 2*. <https://openai.com/dall-e-2/>. 2023.
- [32] OpenAI. *GPT-4 Technical Report*. 2023. arXiv: 2303.08774 [cs.CL].
- [33] M. Oquab et al. *DINOv2: Learning Robust Visual Features without Supervision*. 2023. arXiv: 2304.07193 [cs.CV].
- [34] B. N. Patro et al. “Multimodal Humor Dataset: Predicting Laughter tracks for Sitcoms”. In: *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2021, pp. 576–585. DOI: 10.1109/WACV48630.2021.00062.
- [35] S. Poria et al. “Beneath the Tip of the Iceberg: Current Challenges and New Directions in Sentiment Analysis Research”. In: *IEEE Transactions on Affective Computing* 14.1 (2023), pp. 108–132. DOI: 10.1109/TAFFC.2020.3038167.
- [36] O. Russakovsky et al. *ImageNet Large Scale Visual Recognition Challenge*. 2015. arXiv: 1409.0575 [cs.CV].
- [37] C. Schuhmann et al. *LAION-5B: An open large-scale dataset for training next generation image-text models*. 2022. arXiv: 2210.08402 [cs.CV].
- [38] Q. Sun et al. *EVA-CLIP: Improved Training Techniques for CLIP at Scale*. 2023. arXiv: 2303.15389 [cs.CV].
- [39] G. Team and et al. *Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context*. 2024. arXiv: 2403.05530 [cs.CL].
- [40] T. Winters and P. Delobelle. *Dutch Humor Detection by Generating Negative Examples*. 2020. arXiv: 2010.13652 [cs.CL].
- [41] T. Wolf et al. *HuggingFace’s Transformers: State-of-the-art Natural Language Processing*. 2020. arXiv: 1910.03771 [cs.CL].

## A Appendix

### A.1 Training details

Models were trained using the Adam optimization algorithm with weight decay. A hyperparameter grid search was conducted across learning rates in the set  $\{0.01, 0.001, 0.0001, 0.00001\}$ , batch sizes in the set  $\{4, 8, 16\}$ , and weight decay parameters in the set  $\{0.1, 0.01, 0.001\}$ . Model training proceeded for a fixed number of 10 epochs, with periodic checkpoints. For the final evaluation on the unseen test set, we used the model iteration exhibiting optimal performance on the validation cohort. We used cross-entropy loss for the Binary classification and Comparison tasks while mean square loss was used for the Regression task. For all architectures except GPT-4o and Gemini-Flash, we fine-tuned the models (Pretraining details are mentioned in Section 4.1). For LLaVA we did lora fine-tuning instead of full fine-tuning. To ensure statistical robustness, each experiment was conducted 5 times for all the experiments except for GPT-4o and Gemini-Flash which were run only one time.

Most of the experiments were run on 4 Nvidia GeForce RTX 2080 Ti GPUs which were part of an internal cluster. However, for training LLaVA and some models for the Comparison task, we used an Nvidia A100 GPU provided by the Illinois Computes project which is supported by the University of Illinois Urbana-Champaign and the University of Illinois System.

### A.2 External assets used

We utilized the following assets: The LLaVA repository (Apache-2.0 license) [24, 25], PyTorch [2], huggingface transformers (Apache-2.0 license) [41], and huggingface accelerate (Apache-2.0 license) [11].

Additionally, for the images collected from the internet we provide reference links in the repository. Link: <https://github.com/kreimanlab/HumorDB/>

### A.3 Attention maps

We examined the attention maps using the attention rollout technique [1] on the ViT-Huge model [10]. This helped us understand whether the models focused on the actual humorous parts of images or other biases in the dataset. The attention maps may help to better understand how the models classify the images and identify potential shortcomings (**Fig. 8**).

As an example, consider the case of **Fig. 1**. The attention maps for the vit huge model are shown in **Fig.8**. The model fails to pay attention to the most humorous part of the image (the phone, black rectangle), which is critical to assess whether the image is funny or not. Therefore the model is not able to correctly classify both images.

### A.4 Crowdsourcing details

There were 850 participants: 200 for binary task, 215 for the range task and 435 for the comparison task. The interfaces used by the participants for the three tasks are shown in **Fig. 9**. The generic instructions given for all tasks were:

- **Binary Task:** Please rate if the image is funny or not.
- **Range Task:** Please rate the degree of funniness of the image on a scale from 1 (not funny) to 10 (very funny).
- **Comparison Task:** Please indicate which of the two images is funnier.
- For funny images write a word that makes the image funny, for not funny images, write a word about the most prominent feature of the image.
- The time required to rate all the images is approximately 9-11 minutes
- Only click on the rating buttons once, and wait till the next image loads (maximum 1 second), a message will show you when the next image is being processed.
- Please do not refresh the page. You will lose progress and will have to start again.

- There are 100-120 images in this survey.
- At the end of the survey, we will provide you a code, please store it and use it appropriately to get the reward.
- Click the button below to begin.

Some participants were discarded due to reliability and the amount of outliers in their responses as detailed in Section 3.2. All participants were required to submit all questions in the survey for the response to count.

### A.5 Multimodal models’ answer explanations

For evaluating zero-shot performance of the large multimodal models we test them on testAllSet and testAllPairs sets. We do this to compare the performance of these models with the other fine tuned models on the same test set. The two variants of prompts in binary task as mentioned in section 4.1 were: (i) “Is the image funny?”, and (ii) “Is the image not funny?”. The performance on both prompts were similar so we reported the average for the results. In addition for succinct explanations in a particular format for the figures like Fig. 7, we add a suffix ‘start answer with yes/no then explain’. In this section we present similar figures to Fig. 7 for Gemini-Flash and LLaVA on zero-shot prompting. We also mention the range ratings these models give for the images and the range rating prompt is mentioned in section 4.1. Gemini-Flash answer explanations are presented in Fig. 11. The answer explanation from LLaVA zero-shot are presented in Fig. 12.

### A.6 Evaluating without modified images

This section investigates the impact of excluding modified images from the test set on model performance in our binary classification task. By evaluating our models solely on the original images from the testAllSet, we aim to highlight the significance of incorporating modified images in the dataset. Figure 10 showcases the results of this analysis. The data reveals a substantial performance enhancement for ResNet150, ViT-G/14, ViT-huge, and GPT-4o. Notably, GPT-4o experiences a dramatic 10% increase in accuracy. This dramatic improvement suggests that relying solely on original images, primarily sourced from online repositories, might artificially inflate performance for models that may have encountered these images during training. Our findings demonstrate that HumorDB provides a valuable benchmark for zero-shot evaluation of large multimodal models due to its inclusion of modified images. These modified images represent a challenge not encountered during the training of these models, providing a more realistic assessment of their generalization abilities.

### A.7 Dataset Documentation and metadata

1. Dataset documentation and intended uses. Recommended documentation frameworks include datasheets for datasets, dataset nutrition labels, data statements for NLP, and accountability frameworks: *Available on the github repository here: <https://github.com/kreimanlab/HumorDB/>*
2. URL to website/platform where the dataset/benchmark can be viewed and downloaded by the reviewers: There are two hosting sites of the dataset: <https://github.com/kreimanlab/HumorDB/> and <https://huggingface.co/datasets/kreimanlab/HumorDB>
3. URL to Croissant metadata record documenting the dataset/benchmark available for viewing and downloading by the reviewers: Available through huggingface API: <https://huggingface.co/api/datasets/kreimanlab/HumorDB/croissant>
4. The dataset is licensed under CC BY 4.0 license.
5. Hosting, licensing, and maintenance plan: All information in the Github repository will be maintained.

## B Appendix figures

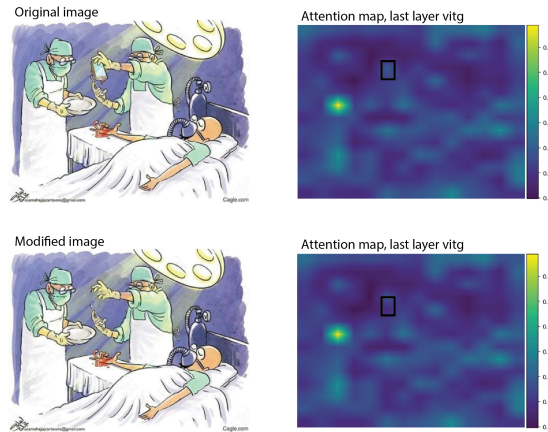
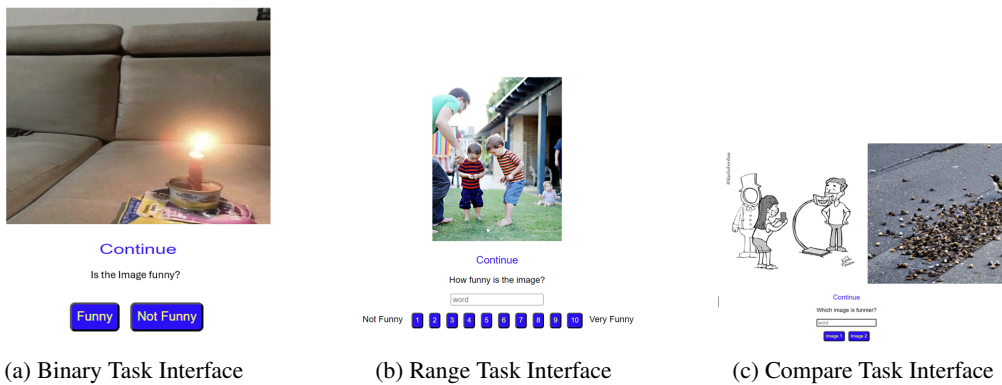


Figure 8: **Attention maps fail to capture elements critical to discern humor.** Attention maps based on the last layer of the vit huge model for the example images from **Fig. 1**. The black rectangle in the attention maps indicates the location of the phone. The maximum attention activation highlights the plate, which does not help distinguish between the original and modified images. Indeed, the model classified both images as funny.



(a) Binary Task Interface

(b) Range Task Interface

(c) Compare Task Interface

Figure 9: **Crowd sourcing interfaces for the three tasks.**

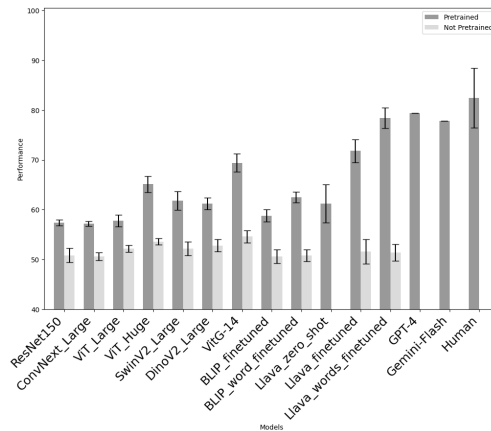
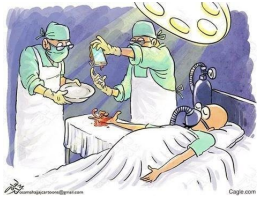




Figure 10: **Binary task results evaluated only with original images in the testAllSet** (i.e., excluding testing on modified images). Format and conventions as in Fig. 5. Several models improve their performance on the original images, indicating that the modified images serve as an important control. GPT-4o shows a particularly surprising increase in performance suggesting that it probably has seen the original images from the internet in its training set.




(a) LLM response: yes. the humor lies in the fact that the patient’s hand is being operated on with a smart-phone. Range rating: 7



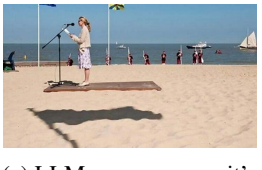
(b) LLM response: yes. the patient’s hand is detached from their body. Range rating: 8



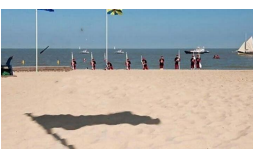
(c) LLM response: yes. the shadow of the women looks like two snakes. it is a dark and ironic humor. Range rating: 7




(d) LLM response: no. the image is not inherently funny. it depicts two people sitting on a couch, and there’s no clear comedic element. Range rating: 3




(e) LLM response: no. it’s a standard photo of a person speaking. there is nothing humorous. Range rating: 2



(f) LLM response: no. it is a picture of a beach with some people on it. Range rating: 1



(g) LLM response: yes. the dog is wearing sunglasses and enjoying a coconut drink. Range rating: 8



(h) LLM response: no. the image is of a child sleeping. there is no humor in the image. Range rating: 2

Figure 11: **Example Gemini-Flash binary classification, range ratings, and explanations.** The human binary and range ratings are as follows for the images: (a) 83%, 6.75, (b) 14%, 5.70, (c) 60%, 3.60, (d) 33%, 3.00, (e) 60%, 2.80, (f) 0%, 5.57, (g) 100%, 5.60, (h) 33%, 4.86.



(a) LLM response: Yes. The humor comes from the juxtaposition of a serious medical situation with the casual and nonchalant behavior of the doctors, who seem more interested in using their phones than focusing on the patient. Range rating: 6



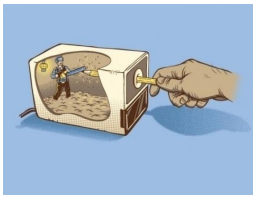
(b) LLM response: Yes, the image is intended to be humorous. It depicts a cartoon of a patient in a hospital bed with two doctors, one holding a tray of food and the other holding a tray of surgical instruments, both seemingly confused about which tray to use. Range rating: 6



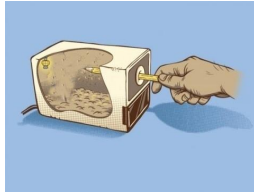
(c) LLM response: Yes, the image is funny. It depicts two people sitting on a couch, seemingly engaged in a conversation, but the shadow on the wall behind them forms a shape that looks like a snake eating its own tail. Range rating: 5



(d) LLM response: Yes, the image is funny. It depicts two people sitting on a couch, but their bodies are mirrored, with their heads on the wrong sides of their bodies, creating a humorous and surreal visual effect. Range rating: 6



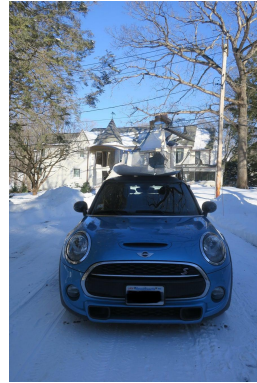
(e) LLM response: Yes, the image is funny. It depicts a cartoon of a person inside a box, seemingly being "pulled out" by a hand using a toothpick. The humor comes from the absurdity and unexpectedness of the situation. Range rating: 6



(f) LLM response: Yes, the image is funny because it shows a hand pushing a button on a box that is labeled "Emergency Button," but instead of a typical emergency response, the box is filled with a bunch of pencils, which is an unexpected and humorous outcome. Range rating: 5



(g) LLM response: Yes, the image is funny because it shows a car with a large inflatable pool toy, specifically a white inflatable animal, strapped to its roof. This is an unusual and humorous sight because it's not common to see such a large and whimsical item being transported in this manner. Range rating: 7



(h) LLM response: No. The image shows a car parked in a snowy environment, which is a common sight during winter. There is nothing inherently funny about the image; it appears to be a typical winter scene. Range rating: 2

**Figure 12: Example LLaVA binary classification, range ratings, and explanations.** The human binary and range ratings are as follows for the images: (a) 83%, 6.75, (b) 14%, 5.70, (c) 60%, 3.60, (d) 33%, 3.00, (e) 60%, 2.80, (f) 0%, 3.00, (g) 50%, 5.25, (h) 0%, 1.00.