

Understanding Figurative Meaning through Explainable Visual Entailment

Anonymous ACL submission

Abstract

Large Vision-Language Models (VLMs) have demonstrated strong capabilities in tasks requiring a fine-grained understanding of literal meaning in images and text, such as visual question-answering or visual entailment. However, there has been little exploration of these models' capabilities when presented with images and captions containing *figurative* meaning, such as metaphors or humor. To close this gap, we propose a new task framing the *figurative meaning understanding* problem as an *explainable visual entailment* task, where the model has to predict whether the image (premise) entails a caption (hypothesis) and justify the predicted label with a textual explanation. The figurative phenomena can be present either in the image, the caption, or both. Utilizing a human-AI collaboration approach, we build the accompanying expert-verified dataset V-FLUTE, containing 6,027 {image, caption, label, explanation} instances spanning five diverse figurative phenomena: metaphors, similes, idioms, sarcasm, and humor. Through automatic evaluation, we find that VLMs struggle to generalize from literal to figurative meaning, particularly when it is present in images. We further conduct human evaluation to identify common errors in VLM reasoning about multimodal figurative meaning.

1 Introduction

Figurative language is integral to human communication, enabling a variety of communicative goals (Roberts and Kreuz, 1994), including affective communication (Fussell and Moss, 2014). Figurative language presents a significant challenge to computational approaches as it requires understanding of implicit meaning behind an expression (Stowe et al., 2022; Shutova, 2011; Veale et al., 2016; Zhou et al., 2021). Recently, Chakrabarty et al. (2022) proposed a task and dataset for Figurative Language Understanding through Textual Explanations (FLUTE) that frames the problem as an

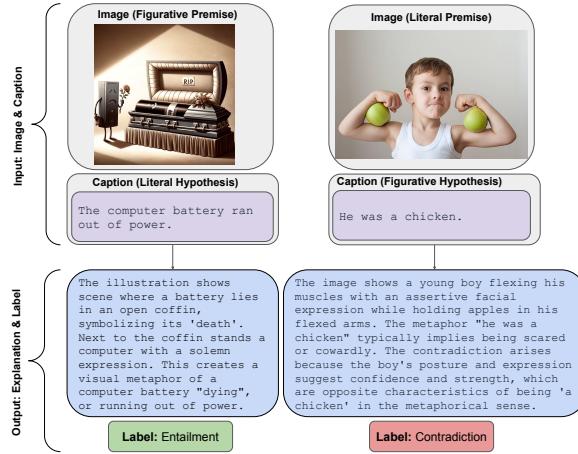


Figure 1: Explainable visual entailment for understanding figurative meaning: given an image and a caption output whether the image entails or contradicts the caption along with a textual explanation.

explainable textual entailment covering a variety of figurative language phenomena in text: metaphors, similes, idioms, and sarcasm. This dataset has been used successfully to advance and benchmark the capabilities of LLMs for understanding figurative language in text (Saakyan et al., 2022; Ziems et al., 2024; Sravanti et al., 2024; Dey et al., 2024).

However, figurative meaning is also prevalent in visual phenomena, such as visual metaphors (Akula et al., 2023; Chakrabarty et al., 2023), multimodal sarcasm (Desai et al., 2022), and humor (Hessel et al., 2023; Hwang and Shwartz, 2023). Yet so far most of the work on vision and language models (VLMs) has focused on understanding literal meaning in images and captions (e.g., ScienceQA (Lu et al., 2022), MMMU (Yue et al., 2024)) including work on explainable visual entailment (Kayser et al., 2021). Building on the idea of FLUTE (Chakrabarty et al., 2022) for text, we present a new dataset for understanding figurative meaning as explainable visual entailment, we call V-FLUTE. Our dataset contains 6,027 {image, caption, la-

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064

bel, explanation} instances spanning diverse figurative phenomena. Each instance contains an image (premise) and a caption (hypothesis) that is either entailed or contradicted by the image. Deciding the entailment relation requires the vision-language model to understand the implicit meaning in both the visual and textual modalities. Our dataset contains figurative phenomena present in the image, in the caption, or in both. In addition, to mitigate the dependence on spurious correlations, to more rigorously investigate reasoning capabilities, and to promote explainability, our task requires the model to generate a plausible explanation for the output label. See Figure 1 for two examples from our dataset.

We make the following contributions towards assessing VLMs ability to understand figurative meaning expressed multimodally:

- V-FLUTE, an expert-verified dataset of 6,027 {image, caption, label, explanation} instances built using a human-LLM collaboration framework covering several phenomena: metaphors, similes, idioms, sarcasm, and humor (Section 3). We will make the dataset available.
- A suite of evaluations to assess current VLMs’ capabilities on this new task of explainable visual figurative entailment (Section 4.2 and 4.3).
- A detailed human evaluation with error analysis yielding insights into the types of errors for different classes of models (Section 5).

2 Related Work

Textual entailment (MacCartney and Manning, 2008; Bowman et al., 2015) and visual entailment (Xie et al., 2019) tasks have been proposed to measure language and multimodal understanding. However, models trained to simply improve label accuracy on these data can be brittle and suffer from spurious correlations (Poliak et al., 2018; Gururangan et al., 2018; McCoy et al., 2019; Gardner et al., 2021). Datasets such as e-SNLI (Camburu et al., 2018) and e-SNLI-VE (Kayser et al., 2021) augment existing entailment datasets with natural language explanations and train models to not only predict the label, but also generate a textual explanation for the reason behind the prediction. However, they only focus on *literal meaning* in text and images. Recently, explainable entailment has been utilized to assess LLMs’ capabilities on understanding figurative language through the FLUTE dataset

(Chakrabarty et al., 2022). FLUTE frames figurative language understanding as an explainable textual entailment task. Recent progress in multimodal models (Li et al., 2022; Alayrac et al., 2022; OpenAI, 2023; Team, 2023; Liu et al., 2023b; Anthropic, 2024) prompts us to asses understanding of figurative meaning present in the multimodal setting, contained in both images and text beyond intent and sentiment (Zhang et al., 2021; Kruk et al., 2019). To this end, we present an equivalent of the FLUTE dataset for the visual modality: V-FLUTE.

3 V-FLUTE Task and Dataset

Following prior work on figurative language understanding in text defined as explainable textual entailment, FLUTE (Chakrabarty et al., 2022), we define *understanding figurative meaning* as an *explainable visual entailment task*: given an image (premise) p and a caption (hypothesis) h , output a textual explanation \hat{e} justifying whether the premise entails or contradicts the hypothesis and assign a label $\hat{y} \in \{\text{Entailment}, \text{Contradiction}\}$. We focus on the binary classification task, since for neutral labels, the explanations would be trivial (simply describing the image).

To build V-FLUTE, we start with existing multimodal figurative datasets which cover phenomena such as metaphors, similes, idioms, sarcasm or humor. We utilize human-AI collaboration frameworks with expert annotators (Chakrabarty et al., 2022; Wiegreffe et al., 2022; Liu et al., 2022) to augment them with expert-verified textual explanations and entailing/contradicting captions. Each instance then includes an image and a caption, and the figurative phenomenon can be either in the image, the caption or in both. An overview of V-FLUTE dataset and *our contributions w.r.t to the source datasets can be found in Table 1*. See examples corresponding to each source dataset in Table 2 as they appear in V-FLUTE. Below, we describe the construction of V-FLUTE by each phenomenon.

3.1 Metaphors, Similes and Idioms

To create visual entailment instances containing metaphors and similes in V-FLUTE, we rely on two existing resources: HAIMet (Chakrabarty et al., 2023) and IRFL (Yosef et al., 2023). Instances from HAIMet contain the metaphor/simile as a part of the premise (image), while those taken from IRFL have the metaphor/simile as a part of the hypothesis (text).

Phenomenon	Data Source	Figurative Part	Our Contribution	# instances
Metaphor/Simile	HAIIVMet (Chakrabarty et al., 2023)	Image	Image Selection Textual Explanations Expert Verification	857 (450 E, 407 C)
	IRFL (Yosef et al., 2023)	Caption	Image Selection Textual Explanations Expert Verification	1,149 (574 E, 575 C)
Idiom	IRFL (Yosef et al., 2023)	Caption	Image Selection Textual Explanations Expert Verification	370 (186 E, 184 C)
Sarcasm	MuSE (Desai et al., 2022)	Caption	Caption Generation Textual Explanations Expert Verification	1,042 (521 E, 521 C)
Humor	MemeCap (Hwang and Shwartz, 2023)	Image	Caption Generation Textual Explanations Expert Verification	1,958 (979 E, 979 C)
	NYCartoons (Hessel et al., 2023)	Image+Caption	Taken As Is	651 (651 E)

Table 1: V-FLUTE dataset composition: 5 figurative phenomena, source datasets, and our contributions. E denotes number of entailment instances, C - contradiction.

3.1.1 IRFL as Data Source

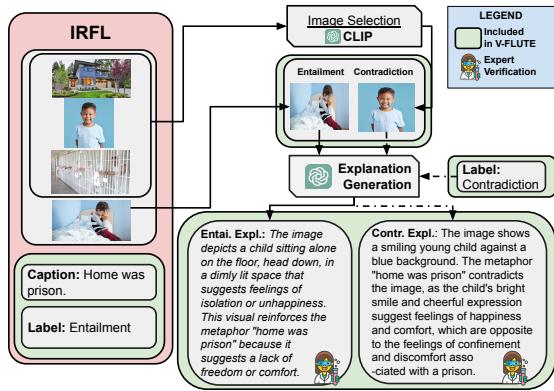


Figure 2: Creation of V-FLUTE instances for metaphors, similes, idioms from IRFL.

Yosef et al. (2023) proposed a benchmark (IRFL) where given a metaphor, a simile or an idiom the model has to distinguish which of the four associated images implies the figurative meaning of the expression. This dataset contains 1,440 figurative expressions, each associated with 4 distinct images. One of those images represents the figurative expression (see Figure 2), and the other 3 act as distractors.

Image Selection. We automatically select images using CLIP (Radford et al., 2021). We select one of the distractor images that have the highest CLIPScore (clip-vit-base-patch16) with the corresponding entailing image to create a challenging, contradictory instance (see where an unrelated

image of a house is discarded when selecting the contradiction instance in Figure 2).

Generating Textual Explanations. We prompt GPT-4 (gpt-4-vision-preview) with the ground truth label, caption, and the image to explain the relationship between the image and the caption.

Expert Verification. We recruit three expert annotators with significant experience in figurative language and visual metaphor understanding on Upwork and ask them to verify the explanation is correct, complete, and concise and if not, edit it (see details in Appendix B). We also ask the annotators to discard rare noisy instances where the caption, image, and label do not fit (due to automatic image selection). Due to relative simplicity of generating the explanation given a literal image, the experts only needed to edit $\approx 7\%$ of the explanations. They also removed $\approx 1\%$ the data, resulting in 1149 {image, caption, label, explanation} instances for metaphors and similes and 370 for idioms.

3.1.2 HAIIVMet as Data Source

Chakrabarty et al. (2023) use a human-AI collaboration framework to generate visual metaphors from linguistic metaphors (HAIIVMet dataset) and propose a visual entailment task as an extrinsic evaluation of dataset quality. The HAIIVMet data consists of 1,193 images of visual metaphors spanning over 958 distinct linguistic metaphors. Each image is associated with a caption that can be contradicting or entailing the image. In addition, each image is associated with a *visual elaboration* that

HAIIMet	IRFL	MuSE	MemeCap	NYCartoons
The faculty meeting was peaceful.	Their relationship is a house on fire.	Oh I just #love having to stare at this while I #work.	Even death won't exempt you from going to work.	Easy for you to say, you're cured!
Contradiction	Entailment	Contradiction	Entailment	Entailment
The image shows a faculty meeting transformed into a dramatic battlefield ... The visual metaphor suggests the faculty meeting was like a war, and not peaceful.	The photo suggests a conflict or an intense emotional situation ... which aligns with the symbolism of a house on fire representing a relationship filled with turmoil or heated arguments.	The image shows Disneyland Resort sign ... the person would like to experience it in person rather than just looking at the sign during work hours.	The image shows RoboCop ... it humorously illustrates a character who has been reanimated as a cyborg to continue working despite having died.	A play on the word "cured". People seek therapy to have their mental problems remedied or cured. But "cured" can also refer to a meat prep technique ...

Table 2: Sample dataset instances form V-FLUTE corresponding to the source datasets displaying images (premise), captions (hypothesis), labels, and explanations [Row 1-5].

210 presents a textual description of the image (See Ap-
 211 pendix A, Figure 7). This visual elaboration was
 212 used in the original paper to generate the visual
 213 metaphors (images).

214 **Generating Textual Explanations.** We aug-
 215 ment the dataset with candidate textual explana-
 216 tions. We prompt ChatGPT (gpt-3.5-0914) to
 217 generate an explanation for every tuple {visual elab-
 218 oration, caption, label} (See Appendix A, Figure
 219 7; and prompt in Appendix F.1.1).

220 **Expert Verification.** Each caption is paired with
 221 up to 5 images. However, since these images were
 222 automatically generated with DALLE-2 using the
 223 visual elaborations, not all are completely faith-
 224 ful. Moreover, some captions and labels were in-
 225 consistent. Finally, automatically generated LLM
 226 candidate explanations are not always correct and
 227 require refining. To tackle these issues, we employ
 228 an expert verification process recruiting the same
 229 three expert annotators as from the IRFL section
 230 above (see details in Appendix B). We ask the
 231 annotators to select the visual metaphor most faithful
 232 to the linguistic metaphor and the visual elaboration
 233 (see Image Selection in Appendix A, Figure
 234 7) or if none were. In addition, we ask them to verify
 235 and edit the explanation if necessary to ensure
 236 correctness, completeness, and conciseness. On
 237 average, experts edited $\approx 65\%$ of the explanations
 238 and 29% of captions, and rejected $\approx 30\%$ of visual
 239 metaphors, resulting in 857 {image, caption, label,

explanation} instances.

3.2 Sarcasm

To create visual entailment instances containing sarcasm, we rely on the MuSE data (Desai et al., 2022).

3.2.1 MuSE as Data Source

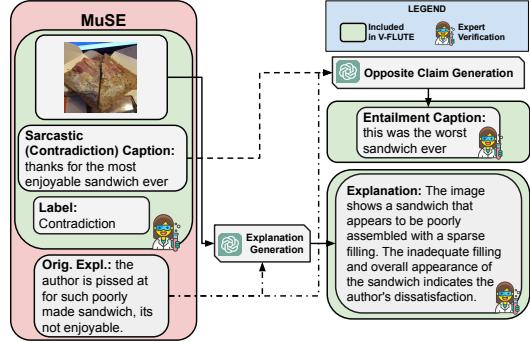


Figure 3: Creation of V-FLUTE instances for sarcasm from MuSE.

The MuSE dataset (Desai et al., 2022) consists of 3510 distinct images, the respective sarcastic captions that act as contradiction instances (see example in Figure 3), and crowd worker written explanations justifying the contradiction.

Generating Entailment Captions. Since the dataset only contains sarcastic instances, there are no captions with an entailment relationship. We

254 generate the entailing captions by prompting GPT-
 255 4 to generate a non-sarcastic version of the caption
 256 while maintaining the user-generated informal style
 257 of the text (see the generated entailment caption in
 258 Figure 3).

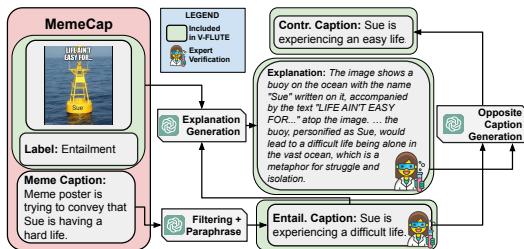
259 **Generating Textual Explanations.** While
 260 the dataset already contains crowdworker-written
 261 explanations, upon inspection, they were often
 262 deemed poor quality, lacking enough details, and
 263 formulaic (e.g., see the crowdworker explanation
 264 in Figure 3). To improve their quality, we use the
 265 dataset’s existing crowdworker explanations and
 266 prompt GPT-4 to rewrite and generate candidate
 267 textual explanations given the caption and the label
 268 (see the re-written explanation in Figure 3). See
 269 the prompt in Appendix F.3.

270 **Expert Verification.** Each image is now paired
 271 with a GPT-4-generated entailing caption, an origi-
 272 nal contradicting caption, and their respective la-
 273 bels and explanations. The same three expert an-
 274 notators checked if the generated explanations are
 275 adequate (i.e., complete, correct, and concise) and
 276 if not, asked to edit them. The experts were also
 277 instructed to discard noisy examples, e.g. when
 278 the image does not contradict the sarcastic cap-
 279 tion. On average, experts edited $\approx 13\%$ of the
 280 initial explanations and rejected $\approx 18\%$ of the ex-
 281 amples, resulting in 1,042 {image, caption, label,
 282 explanation} instances.

283 3.3 Humor

284 For multimodal humor, we rely on two datasets:
 285 MemeCap (Hwang and Shwartz, 2023) and New
 286 Yorker cartoons (Hessel et al., 2023).

287 3.3.1 MemeCap as Data Source



288 Figure 4: Creation of V-FLUTE instances for humor
 289 from MemeCap.

290 This dataset consists of memes along with their
 291 captions that describe the meme poster’s intent (see
 292 example in Figure 4). Memes frequently contain
 293 implicit, non-literal meaning (Lestari, 2019) and

294 rely on visual metaphors (Piata, 2016), posing a
 295 challenge to VLMs.

296 **Caption Generation.** Meme captions are not
 297 suited for an entailment task, so we prompt GPT-4
 298 with the original caption to generate an entailing
 299 caption in the form of a claim from it (see example
 300 in Figure 4). We filter these set of samples further
 301 with GPT-4 by asking whether the image entails
 302 the caption and only selecting positive instances.
 303 In addition to generating captions that entail the
 304 meme, we generate contradicting captions using
 305 GPT-4.

306 **Generating Textual Explanations.** We
 307 prompted GPT-4 with the ground truth label in
 308 the prompt to explain the relationship between the
 309 image and the caption. See prompts in Appendix
 310 F.4.

311 **Expert Verification.** We hire the same three
 312 expert annotators to ensure the correctness of the
 313 data. Each annotator is tasked with verifying that
 314 1) the generated caption fits the image and 2) the
 315 explanation is correct and complete, and if not,
 316 make the necessary changes. We also ask to discard
 317 samples with inappropriate content. Experts edited
 318 $\approx 35\%$ of the explanations and 15% of captions
 319 on average, and discarded $\approx 2\%$ of inappropriate
 320 instances, resulting in 1958 {image, caption, label,
 321 explanation} instances.

322 3.3.2 NYCartoons as Data Source

323 The NYCartoons dataset (Hessel et al., 2023) con-
 324 tains 651 high-quality instances from the New
 325 Yorker Cartoon Caption Contest. Each instance
 326 consists of an image paired with a humorous cap-
 327 tion and an explanation of why this combination
 328 of the caption and the image is funny. We utilize
 329 this data as is by treating the image as entailing
 330 the caption, so the explanation of the entailment
 331 relationship is the explanation of the joke.

332 3.4 Dataset Statistics

333 We split our data into 4,578 training, 726 valida-
 334 tion, and 723 testing instances. Detailed counts per
 335 phenomenon and dataset, as well as other statistics,
 336 are in Appendix C.

337 4 Experiments

338 We empirically study how several baseline mod-
 339 els perform on the task of explainable visual en-
 340 tailment. We investigate both off-the-shelf and
 341 fine-tuned model performance. We provide human

340 baseline performance in Appendix J. All hyperparameters are in Appendix E.
 341

342 4.1 Models

343 We select a variety of models for our study (see
 344 taxonomy in Appendix, Figure 10). For **off-the-**
 345 **shelf models**, we explore both *open* and *API-based*
 346 models. For *open* models, we select the (current)
 347 state-of-the-art LLaVA-1.6 models (Liu et al.,
 348 2024). LLaVA is one of the simplest, yet one of the
 349 most high-performing VLM architectures currently
 350 available. It utilizes a pretrained large language
 351 model (e.g., Mistral-7B (Jiang et al., 2023)) and
 352 a vision-language cross-modal connector (e.g., an
 353 MLP layer) to align the vision encoder (e.g., CLIP
 354 (Radford et al., 2021)) outputs to the language mod-
 355 els. We select LLaVA-1.6 models in their 7B and
 356 34B configurations (LLaVA-v1.6-7B and LLaVA-
 357 v1.6-34B respectively) and refer to them as *LLaVA-*
 358 *ZS-7B* and *LLaVA-ZS-34B*. Both models have been
 359 instruction-tuned on less than 1M visual instruc-
 360 tion tuning samples to act as general language and
 361 vision assistants. We also utilize *Compositional*
 362 *Chain-of-Thought Prompting* proposed by Mitra
 363 et al. (2023) denoted by LLaVA-ZS-7B-SG and
 364 LLaVA-ZS-34B-SG (see description and results
 365 discussion in Appendix H).

366 For *API-based* models, we select three widely
 367 available state-of-the-art VLMs: Claude-3
 368 Opus (claude-3-opus-20240229)(Anthropic,
 369 2024), GPT-4 (gpt-4-1106-vision-preview)
 370 (OpenAI, 2023) and GeminiPro
 371 (gemini-pro-vision)(Team, 2023).

372 For **fine-tuned** models, we focus on fine-tuning
 373 the LLaVA-1.5-7B model¹ (Liu et al., 2023a). To
 374 minimize bias for a single instruction, we fine-tune
 375 and evaluate the models on a set of 21 instruction
 376 paraphrases (see Appendix Table 7). Three model
 377 configurations are tested:

- 378 • *LLaVA-VF* is the same checkpoint fine-tuned on
 379 the training set of V-FLUTE. We also fine-tune
 380 the model with a white square instead of the V-
 381 FLUTE image (denoted by –Image).
- 382 • *LLaVA-eViL* and *LLaVA-eViL+VF* are check-
 383 points of LLaVA-v1.5-7B further fine-tuned on
 384 the eViL (e-SNLI-VE) dataset for explainable
 385 visual entailment (Kayser et al., 2021) converted
 386 to the instruction format or on both eViL and
 387 V-FLUTE. We removed neutral label instances,

388 ¹Fine-tuning code for 1.6 model was not published as of
 389 writing of this paper.

Model Name	F1@0	F1@53	F1@60
<i>Random Baseline</i>	49.82	-	-
<i>Fine-tuned</i>			
LLaVA-7B			
--> VF	72.78	60.66	47.12
--> – Image	64.77	53.28	39.37
--> eViL	54.34	4.11	0.55
--> + VF	74.91	62.34	48.80
<i>Off-the-shelf</i>			
<i>Open</i>			
LLaVA-ZS			
--> 7B	45.44	35.57	18.38
--> + SG	52.94	39.27	14.86
--> 34B	55.60	<u>48.32</u>	<u>31.83</u>
--> + SG	<u>58.08</u>	45.74	26.77
<i>API-based</i>			
--> Gemini	53.70	39.72	19.01
--> 5-shot	67.25	56.04	37.14
--> Claude	56.07	45.37	22.31
--> 5-shot	67.79	58.70	35.32
--> GPT-4	64.00	56.22	38.56
--> 5-shot	<u>69.36</u>	<u>61.95</u>	49.81

390 Table 3: F1 Score results for different models across
 391 thresholds 0.0, 0.53, and 0.6 for explanation score. Best
 392 result overall is in bold, best result in each category is
 393 underlined.

394 which resulted in 275,815 training instances and
 395 10,897 validation instances.

396 4.2 Automatic Metrics

397 Since our goal is to ensure models provide an an-
 398 swer for the right reasons, ideally, we would only
 399 count predictions as correct when the explanation is
 400 also correct. Similarly to prior work (Chakrabarty
 401 et al., 2022), we utilize both the standard F1 score
 402 and an adjusted score that accounts for explana-
 403 tion quality: F1@ExplanationScore. The Explan-
 404 ationScore computes the average of BERTScore
 405 (Zhang* et al., 2020) and BLEURT (Sellam et al.,
 406 2020) between model-generated and V-FLUTE ex-
 407 planations. We report F1@0 (simply F1 score),
 408 F1@53² (only predictions with ExplanationScore
 409 > 53 are considered correct), and F1@60.

410 4.3 Automatic Evaluation Results

411 We include results per phenomenon in Appendix I
 412 and discussion few-shot and scene graph prompting
 413 in Appendix H. Table 3 shows the results, inform-
 414 ing the following insights:

415 **A literal visual entailment dataset does not solve**
 416 **the figurative visual entailment task.** Fine-

417 ²Thresholds selected based on human evaluation of expla-
 418 nation quality in Appendix K.

tuning only on e-ViL barely improves over a random baseline (54.34 F1@0) and underperforms compared with the models fine-tuned on V-FLUTE (72.78 F1@0). Moreover, the explanations are of poor quality (0.55 F1@60). *This indicates that models trained on a literal visual entailment task struggle to generalize to figurative meaning, supporting the challenging nature of our dataset.*

The strongest model fine-tuned on V-FLUTE (LLaVA-7B-eViL+VF) outperforms the best off-the-shelf model (GPT-4-5shot) in terms of the F1@0 score ($p < 0.03^3$). It performs competitively when incorporating the reference-based ExplanationScore, with GPT-4 leading slightly as it is the model with which the candidate explanations were generated.

When figurative meaning is in the image rather than text, models perform worse. We plot the relative percentage decrease between F1@0 and F1@60 for LLaVA-eViL-VF, LLaVA-34B-SG, and GPT-4-5shot in Figure 5. Higher performance drop indicates higher difficulty of generating the correct explanation. For all models, we see a substantial decrease in performance, especially on challenging phenomena such as Humor (NYCartoons). The percentage drop is substantially higher for all models for the HAIVMet subset rather than the IRFL dataset, which contains metaphors in the image rather than in the text. *This suggests it is harder for models to generate correct explanations when the figurative meaning is contained in the image rather than in the text, indicating the need to expand the presence of figurative phenomena in existing visual datasets.*

VLMs benefit from visual information when dealing with figurative phenomena and do not just rely on the input text to make their prediction. We utilize a hypothesis-only baseline (Poliak et al., 2018) by including a model fine-tuned on the V-FLUTE dataset, but with a white square as the image input, denoted as –Image. Fine-tuning on the full V-FLUTE dataset shows an improvement of over 8 points in F1@0 (better with $p < 0.002$).

5 Human Evaluation and Error Analysis

We conduct human evaluation of generated explanations to assess their quality and identify key errors in reasoning about multimodal figurative meaning.

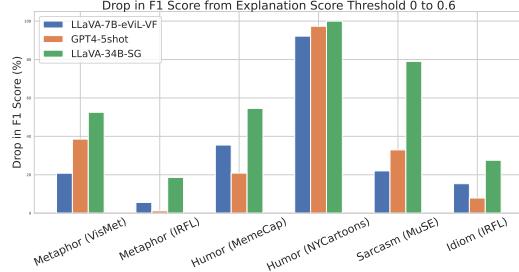


Figure 5: % Drop in F1 score for various models by source dataset between 0 to 0.6. Higher drop indicates higher proportion of wrongly generated explanations.

We recruit two expert annotators with background in linguistics for the task and sample 95 random instances from the test set. For each instance, we first provide the annotators with the image, caption and reference explanation and ask the annotators to choose the right label. If the annotator succeeds, they can view the rest of the task, which consists of 3 explanations from our top models by F1@0 in each category: LLaVA-eViL-VF, LLaVA-34B-SG, GPT-4-5shot. The explanations are taken for both correct and incorrect model predictions. For each explanation, we ask whether the explanation is adequate (accurate, correct, complete and concise). If not, we ask them to identify one of the errors based on the following taxonomy:

- **Hallucination:** explanation is not faithful to the image, indicating difficulties with visual comprehension (e.g., generates “blunt tip” when the pencil tip is actually sharp in row 1 of Table 4).
- **Unsound reasoning:** sentences do not adhere to natural logic or violate common sense (e.g., concluding that an upwards arrow and lots of money imply an economic crisis, see row 3).
- **Incomplete reasoning:** while overall the explanation makes sense, it does not address the key property reasons why the image entails or contradicts the caption (for example, does not address the figurative part in the image, see row 2).
- **Verbosity:** the explanation is too verbose.

5.1 How Do Models Perform According to Humans?

In Table 5, we show adequacy and preference rates for explanations from the 3 systems, where an explanation is deemed adequate or preferred if both annotators agreed it is, and inadequate if both agreed it is not. The average IAA using Cohen’s

³ p values reported via paired bootstrap test (Koehn, 2004)

Error Type	Image	Caption	Label and Explanation
Hallucination (<i>describes sharp pencil as having a blunt tip</i>)		The tip is as sharp as a knife.	Predicted Label: Contradiction <i>Predicted Explanation:</i> The image depicts a pencil with a blunt tip . [...]
Incomplete (<i>does not address metaphorical meaning of iceberg imagery</i>)		They looked at each other in an almost hostile way.	Predicted Label: Entailment <i>Predicted Explanation:</i> The image depicts two characters looking at each other with a tense and confrontational expression, which is consistent with the claim [...]
Unsound (<i>concludes that large amount of money + upward trending chart implies rough waters</i>)		The world recession of 2001 was one of the worst in history.	Predicted Label: Entailment <i>Predicted Explanation:</i> The image depicts a man in a boat, holding a large amount of money and a chart with an arrow pointing upwards. This suggests that the man is navigating through rough waters [...]

Table 4: Examples of error types generated explanations.

	LLaVA-7B eViL+VF	LLaVA-34B SG	GPT-4 (5 shot)
Adequate %	33.78	29.85	50.67
Preference %	23.08	7.69	44.23

Table 5: Adequacy and Preference rates for generated explanations.

κ is 0.47, indicating moderate agreement (Cohen, 1960). We observe that the teacher GPT-4 model is leading in terms of the adequacy of the explanations and preference rate, as expected from a larger system. Yet still only half of its explanations are considered adequate, confirming that despite good performance on the F1@0 scores, *the models are not yet capable of producing adequate textual explanations in many instances*.⁴

5.2 What Errors Do Models Make?

We perform an analysis of the types of errors from each model when the explanations are considered inadequate in the above evaluation. In Figure 6, we illustrate the normalized frequency of error types when both annotators agree that the explanation is not adequate (i.e., out of all errors for this model, what percentage is each type of error?). Overall, the annotators did not consider verbosity to be a major issue of the systems. For GPT-4, the leading error type is hallucination, indicating the need to improve faithful image recognition even in the

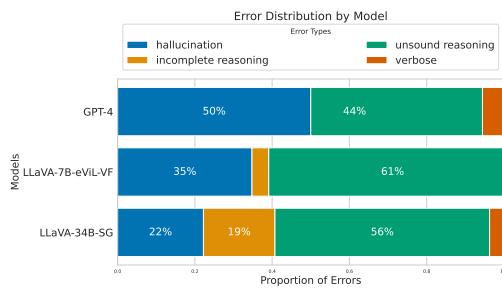


Figure 6: Normalized frequency of main error types in the explanation by model.

most advanced models. Comparing LLaVA-34B-SG and the fine-tuned model, we see that for the scene graph model a larger percentage of errors is due to incomplete reasoning (possibly due to focusing on the scene graph description rather than the underlying figurative phenomena). For both models, the main error type is unsound reasoning, indicating difficulty for the models to consistently reason about multimodal figurative inputs.

6 Conclusions

We introduce a novel dataset for understanding figurative meaning in multimodal input, V-FLUTE, via an explainable visual entailment task. Our dataset consists of 6,027 {image, caption, label, explanation} instances covering diverse phenomena. We find that VLMs struggle to generalize from literal to figurative meaning, particularly in images, and identify error types in VLM reasoning about multimodal figurative phenomena.

⁴Note this does not undermine our human-AI dataset creation framework, as 1) the LLM is conditioned on the correct label, 2) its explanation is edited by an expert annotator.

535 7 Ethics

536 Following prior work in human-AI collaboration
537 for complex text and image generation
538 (Chakrabarty et al., 2022, 2023; CH-Wang et al.,
539 2023; Saakyan and Muresan, 2023), we opt for an
540 expert-AI collaboration framework where experts
541 edit the initial generations by the language model.
542 Expert feedback is essential to improve the quality
543 of the data, as previous work has identified that
544 crowdworkers on platforms such as Amazon Me-
545 chanical Turk could be unreliable for open-ended
546 generation tasks (Karpinska et al., 2021), and might
547 even rely on ChatGPT to provide their answers
548 (Veselovsky et al., 2023). To mitigate these effects,
549 in this work, annotators were recruited through the
550 Upwork platform, allowing to select for relevant
551 level of expertise and verify, e.g., educational and
552 professional background of the annotators. All re-
553 cruited annotators have significant background in
554 figurative language understanding and have formal
555 educational background in linguistics or literature.
556 All of the annotators are fluent or native/bilingual
557 level in English. Workers on UpWork were in-
558 formed that that the work they were doing was going
559 to be used for research purposes. All are fairly
560 compensated with USD \$20 to \$25 per hour with
561 self-reported time needed to complete the tasks.
562 The total budget for the annotation and GPT-4 gen-
563 erations was \approx \$5,000 USD. We estimate that it
564 would take approximately 3 times longer to com-
565 plete the annotation task without the pre-generated
566 explanation, so we estimate that the cost would
567 have at least tripled if the human-AI collaboration
568 approach was not utilized. Workers were paid their
569 wages in full immediately upon the completion of
570 their work. All data collected by human respon-
571 dents were fully anonymized. We do not report
572 demographic or geographic information, given the
573 limited number of respondents, so as to maintain
574 full anonymity.

575 8 Limitations

576 We would like to acknowledge the following lim-
577 itations of our work. The textual explanations in
578 V-FLUTE dataset were generated with the help of
579 the strongest LLM available at the time of writ-
580 ing the paper, GPT-4. Despite our best efforts in
581 mitigating biases with expert human verification,
582 idiosyncrasies pertaining to GPT-4 outputs may
583 still be present in the text. This means that it is
584 potentially possible for the underlying biases of

585 source datasets of language model generations to
586 propagate into our resource, which we wish to mit-
587 iate by carefully examining each dataset instance
588 by one of the 3 expert annotators.

589 Reference-based evaluation has fundamental
590 flaws such as not considering all possible expla-
591 nations, which would be impossible to collect.
592 However, current reference-free metrics for free-
593 text rationales may still have flaws such as bias
594 towards length or the evaluator LLM (Stureborg
595 et al., 2024). When evaluating textual expla-
596 nations against these references, as is the case with
597 any reference-based evaluation, there may also be a
598 preference towards models which output text closer
599 in distribution to the GPT-4 model. Because of that,
600 it is important to utilize the dataset in order to com-
601 pare models other than the teacher model and pay
602 more attention to the F1@0 scores, which represent
603 simple classification scores and do not require out-
604 puts to be similar in distribution. While we showed
605 a relatively high predictive power of automatic ex-
606 planation scores to predict human judgements (see
607 Appendix K), future work may focus on increasing
608 reliability of reference-based and reference-free
609 textual explanation evaluation methods.

610 We also would like to note that images from the
611 HAIIVMet dataset (Chakrabarty et al., 2023) are
612 AI-generated. However, the majority of the rest
613 of images in V-FLUTE are not AI-generated but
614 naturally occurring or created by humans. Still,
615 to mitigate potential biases from the AI-generated
616 images, every instance of the data was examined
617 during the expert verification stage as described in
618 the paper.

619 Label predictions by language models can vary
620 significantly with slight differences in prompt word-
621 ing (Sclar et al., 2023), which is why during fine-
622 tuning and inference we utilize over 20+ different
623 templates of instructions (see Table 7). Neverthe-
624 less, it is important to consider the models’ ex-
625 planations to better assess their understanding of
626 the phenomena, which we hope to enable with our
627 explainable figurative visual entailment dataset.

628 References

- Arjun R Akula, Brendan Driscoll, Pradyumna Narayana,
Soravit Changpinyo, Zhiwei Jia, Suyash Damle,
Garima Pruthi, Sugato Basu, Leonidas Guibas,
William T Freeman, et al. 2023. Metacue: Towards
comprehensive visual metaphors research. In *Pro-
ceedings of the IEEE/CVF Conference on Computer
Vision and Pattern Recognition*, pages 23201–23211.

636	Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. <i>Advances in Neural Information Processing Systems</i> , 35:23716–23736.	693
637		694
638		695
639		696
640		697
641		698
642	Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www.anthropic.com/news/clause-3-family .	699
643		700
644		701
645		702
646		703
647		704
648	Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.	705
649		706
650		707
651		708
652	Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. In <i>Advances in Neural Information Processing Systems</i> , volume 31. Curran Associates, Inc.	709
653		710
654		711
655		712
656		713
657	Sky CH-Wang, Arkadiy Saakyan, Oliver Li, Zhou Yu, and Smaranda Muresan. 2023. Sociocultural norm similarities and differences via situational alignment and explainable textual entailment. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 3548–3564, Singapore. Association for Computational Linguistics.	714
658		715
659		716
660		717
661		718
662	Tuhin Chakrabarty, Arkadiy Saakyan, Debanjan Ghosh, and Smaranda Muresan. 2022. FLUTE: Figurative language understanding through textual explanations. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 7139–7159, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	719
663		720
664		721
665		722
666		723
667		724
668		725
669		726
670		727
671	Tuhin Chakrabarty, Arkadiy Saakyan, Olivia Winn, Artemis Panagopoulou, Yue Yang, Marianna Apidianaki, and Smaranda Muresan. 2023. I spy a metaphor: Large language models and diffusion models co-create visual metaphors. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 7370–7388, Toronto, Canada. Association for Computational Linguistics.	728
672		729
673		730
674		731
675		732
676		733
677		734
678		735
679	Jacob Cohen. 1960. A coefficient of agreement for nominal scales. <i>Educational and psychological measurement</i> , 20(1):37–46.	736
680		737
681		738
682	Poorav Desai, Tanmoy Chakraborty, and Md Shad Akhtar. 2022. Nice perfume. how long did you marinate in it? multimodal sarcasm explanation. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 36(10):10563–10571.	739
683		740
684		741
685		742
686		743
687	Gourab Dey, Adithya V Ganesan, Yash Kumar Lal, Manal Shah, Shreyashee Sinha, Matthew Matero, Salvatore Giorgi, Vivek Kulkarni, and H Andrew Schwartz. 2024. Socialite-llama: An instruction-tuned model for social scientific tasks. <i>arXiv preprint arXiv:2402.01980</i> .	744
688		745
689		746
690		747
691		748
692		749
693	Susan R Fussell and Mallie M Moss. 2014. Figurative language in emotional communication. In <i>Social and cognitive approaches to interpersonal communication</i> , pages 113–141. Psychology Press.	693
694		694
695		695
696		696
697	Matt Gardner, William Merrill, Jesse Dodge, Matthew Peters, Alexis Ross, Sameer Singh, and Noah A. Smith. 2021. Competency problems: On finding and removing artifacts in language data. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1801–1813, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	697
698		698
699		699
700		700
701		701
702		702
703		703
704		704
705	Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.	705
706		706
707		707
708		708
709		709
710		710
711		711
712		712
713		713
714	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. {DEBERTA}: {DECODING}-{enhanced} {bert} {with} {disentangled} {attention}. In <i>International Conference on Learning Representations</i> .	714
715		715
716		716
717		717
718		718
719	Jack Hessel, Ana Marasovic, Jena D. Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff, and Yejin Choi. 2023. 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.	719
720		720
721		721
722		722
723		723
724		724
725		725
726		726
727		727
728	Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8003–8017, Toronto, Canada. Association for Computational Linguistics.	728
729		729
730		730
731		731
732		732
733		733
734		734
735		735
736		736
737	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	737
738		738
739		739
740		740
741		741
742	EunJeong Hwang and Vered Shwartz. 2023. MemeCap: A dataset for captioning and interpreting memes. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1433–1445, Singapore. Association for Computational Linguistics.	742
743		743
744		744
745		745
746		746
747		747
748	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego	748
749		749

750	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	805
751	laume Lample, Lucile Saulnier, Lélio Renard Lavaud,	806
752	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	807
753	Thibaut Lavril, Thomas Wang, Timothée Lacroix,	808
754	and William El Sayed. 2023. <i>Mistral 7b</i> . Preprint,	809
755	arXiv:2310.06825.	
756	Marzena Karpinska, Nader Akoury, and Mohit Iyyer.	
757	2021. <i>The perils of using Mechanical Turk to eval-</i>	
758	<i>uate open-ended text generation</i> . In <i>Proceedings of the</i>	
759	<i>2021 Conference on Empirical Methods in Natural</i>	
760	<i>Language Processing</i> , pages 1265–1285, Online and	
761	Punta Cana, Dominican Republic. Association for	
762	Computational Linguistics.	
763	Maxime Kayser, Oana-Maria Camburu, Leonard	
764	Salewski, Cornelius Emde, Virginie Do, Zeynep	
765	Akata, and Thomas Lukasiewicz. 2021. e-vil: A	
766	dataset and benchmark for natural language expla-	
767	nations in vision-language tasks. In <i>Proceedings of</i>	
768	<i>the IEEE/CVF international conference on computer</i>	
769	<i>vision</i> , pages 1244–1254.	
770	Philipp Koehn. 2004. <i>Statistical significance tests for</i>	
771	<i>machine translation evaluation</i> . In <i>Proceedings of the</i>	
772	<i>2004 Conference on Empirical Methods in Natural</i>	
773	<i>Language Processing</i> , pages 388–395, Barcelona,	
774	Spain. Association for Computational Linguistics.	
775	Julia Kruk, Jonah Lubin, Karan Sikka, Xiao Lin, Dan	
776	Jurafsky, and Ajay Divakaran. 2019. <i>Integrating</i>	
777	<i>text and image: Determining multimodal document</i>	
778	<i>intent in Instagram posts</i> . In <i>Proceedings of the</i>	
779	<i>2019 Conference on Empirical Methods in Natu-</i>	
780	<i>ral Language Processing and the 9th International</i>	
781	<i>Joint Conference on Natural Language Processing</i>	
782	(EMNLP-IJCNLP), pages 4622–4632, Hong Kong,	
783	China. Association for Computational Linguistics.	
784	Widia Lestari. 2019. <i>Irony analysis of memes on insta-</i>	
785	<i>gram social media</i> . <i>Pioneer: Journal of Language</i>	
786	<i>and Literature</i> , 10(2):114–123.	
787	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	
788	Hoi. 2022. Blip: Bootstrapping language-image pre-	
789	training for unified vision-language understanding	
790	and generation. In <i>International Conference on Ma-</i>	
791	<i>chine Learning</i> , pages 12888–12900. PMLR.	
792	Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and	
793	Yejin Choi. 2022. <i>WANLI: Worker and AI collabora-</i>	
794	<i>tion for natural language inference dataset creation</i> .	
795	In <i>Findings of the Association for Computational</i>	
796	<i>Linguistics: EMNLP 2022</i> , pages 6826–6847, Abu	
797	Dhabi, United Arab Emirates. Association for Com-	
798	putational Linguistics.	
799	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	
800	Lee. 2023a. Improved baselines with visual instruc-	
801	tion tuning.	
802	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan	
803	Zhang, Sheng Shen, and Yong Jae Lee. 2024. <i>Llava-</i>	
804	<i>next: Improved reasoning, ocr, and world knowledge</i> .	
500	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	805
501	Lee. 2023b. <i>Visual instruction tuning</i> . In <i>Ad-</i>	806
502	<i>vances in Neural Information Processing Systems</i> ,	807
503	volume 36, pages 34892–34916. Curran Associates,	808
504	Inc.	809
505	Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-	810
506	Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter	811
507	Clark, and Ashwin Kalyan. 2022. Learn to explain:	812
508	Multimodal reasoning via thought chains for science	813
509	question answering. In <i>The 36th Conference on Neu-</i>	814
510	<i>ral Information Processing Systems (NeurIPS)</i> .	815
511	Bill MacCartney and Christopher D. Manning. 2008.	816
512	<i>Modeling semantic containment and exclusion in nat-</i>	817
513	<i>ural language inference</i> . In <i>Proceedings of the 22nd</i>	818
514	<i>International Conference on Computational Linguis-</i>	819
515	<i>tics (Coling 2008)</i> , pages 521–528, Manchester, UK.	820
516	Coling 2008 Organizing Committee.	821
517	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. <i>Right</i>	822
518	<i>for the wrong reasons: Diagnosing syntactic heuri-</i>	823
519	<i>tics in natural language inference</i> . In <i>Proceedings of</i>	824
520	<i>the 57th Annual Meeting of the Association for Com-</i>	825
521	<i>putational Linguistics</i> , pages 3428–3448, Florence,	826
522	Italy. Association for Computational Linguistics.	827
523	Chanchik Mitra, Brandon Huang, Trevor Darrell, and	828
524	Roei Herzig. 2023. <i>Compositional chain-of-thought</i>	829
525	<i>prompting for large multimodal models</i> . Preprint,	830
526	arXiv:2311.17076.	831
527	OpenAI. 2023. Gpt-4v(ision) system card.	832
528	https://cdn.openai.com/papers/GPTV_System_Card.pdf .	833
529		834
530	Anna Piata. 2016. When metaphor becomes a joke:	835
531	Metaphor journeys from political ads to internet	836
532	memes. <i>Journal of Pragmatics</i> , 106:39–56.	837
533	Adam Poliak, Jason Naradowsky, Aparajita Haldar,	838
534	Rachel Rudinger, and Benjamin Van Durme. 2018.	839
535	<i>Hypothesis only baselines in natural language infer-</i>	840
536	<i>ence</i> . In <i>Proceedings of the Seventh Joint Confer-</i>	841
537	<i>ence on Lexical and Computational Semantics</i> , pages	842
538	180–191, New Orleans, Louisiana. Association for	843
539	Computational Linguistics.	844
540	Amy Pu, Hyung Won Chung, Ankur P Parikh, Sebastian	845
541	Gehrman, and Thibault Sellam. 2021. Learning	846
542	compact metrics for mt. In <i>Proceedings of EMNLP</i> .	847
543	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	848
544	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	849
545	try, Amanda Askell, Pamela Mishkin, Jack Clark,	850
546	et al. 2021. Learning transferable visual models from	851
547	natural language supervision. In <i>International confer-</i>	852
548	<i>ence on machine learning</i> , pages 8748–8763. PMLR.	853
549	Richard M. Roberts and Roger J. Kreuz. 1994. <i>Why</i>	854
550	<i>do people use figurative language?</i> <i>Psychological</i>	855
551	<i>Science</i> , 5(3):159–163.	856
552	Arkadiy Saakyan, Tuhin Chakrabarty, Debanjan Ghosh,	857
553	and Smaranda Muresan. 2022. A report on the	858

859	figlang 2022 shared task on understanding figurative language. In <i>Proceedings of the 3rd Workshop on Figurative Language Processing (FLP)</i> , pages 178–183.	Adina Williams, Nikita Nangia, and Samuel Bowman.	915
860		2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	916
861			917
862			918
863	Arkadiy Saakyan and Smaranda Muresan. 2023. Iclef: In-context learning with expert feedback for explainable style transfer. <i>Preprint</i> , arXiv:2309.08583.		919
864			920
865			921
866	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In <i>The Twelfth International Conference on Learning Representations</i> .		922
867			923
868			
869			
870			
871			
872	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7881–7892, Online. Association for Computational Linguistics.		924
873			925
874			926
875			927
876			
877			
878	Ekaterina V Shutova. 2011. Computational approaches to figurative language. Technical report, University of Cambridge, Computer Laboratory.		
879			
880			
881	Settaluri Lakshmi Sravanthi, Meet Doshi, Tankala Pavani Kalyan, Rudra Murthy, Pushpak Bhattacharyya, and Raj Dabre. 2024. Pub: A pragmatics understanding benchmark for assessing llms’ pragmatics capabilities. <i>arXiv preprint arXiv:2401.07078</i> .		
882			
883			
884			
885			
886	Kevin Stowe, Prasetya Utama, and Iryna Gurevych. 2022. IMPLI: Investigating NLI models’ performance on figurative language. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5375–5388, Dublin, Ireland. Association for Computational Linguistics.		928
887			929
888			930
889			931
890			932
891			
892			
893	Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators. <i>arXiv preprint arXiv:2405.01724</i> .		933
894			934
895			935
896			936
897	Gemini Team. 2023. Gemini: A family of highly capable multimodal models. <i>Preprint</i> , arXiv:2312.11805.		937
898			938
899			939
900			940
901	Tony Veale, Ekaterina Shutova, and Beata Beigman Klebanov. 2016. <i>Metaphor: A computational perspective</i> . Morgan & Claypool Publishers.		941
902			
903			
904			
905			
906	Veniamin Veselovsky, Manoel Horta Ribeiro, and Robert West. 2023. Artificial artificial artificial intelligence: Crowd workers widely use large language models for text production tasks. <i>Preprint</i> , arXiv:2306.07899.		
907			
908			
909			
910			
911			
912			
913			
914	Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. 2022. Reframing human-AI collaboration for generating free-text explanations. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 632–658, Seattle, United States. Association for Computational Linguistics.		954
915			955
916			956
917			957
918			958
919			959
920			
921			
922			
923			
924	Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. Visual entailment: A novel task for fine-grained image understanding. <i>ArXiv</i> , abs/1901.06706.		924
925			925
926			926
927			927
928	Ron Yosef, Yonatan Bitton, and Dafna Shahaf. 2023. IRFL: Image recognition of figurative language. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 1044–1058, Singapore. Association for Computational Linguistics.		928
929			929
930			930
931			931
932			932
933	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> .		933
934			934
935			935
936			936
937			937
938			938
939			939
940			940
941			941
942	Dongyu Zhang, Minghao Zhang, Heting Zhang, Liang Yang, and Hongfei Lin. 2021. MultiMET: A multimodal dataset for metaphor understanding. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3214–3225, Online. Association for Computational Linguistics.		942
943			943
944			944
945			945
946			946
947			947
948			948
949			949
950	Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> .		950
951			951
952			952
953			953
954	Jianing Zhou, Hongyu Gong, and Suma Bhat. 2021. PIE: A parallel idiomatic expression corpus for idiomatic sentence generation and paraphrasing. In <i>Proceedings of the 17th Workshop on Multiword Expressions (MWE 2021)</i> , pages 33–48, Online. Association for Computational Linguistics.		954
955			955
956			956
957			957
958			958
959			959
960	Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? <i>Computational Linguistics</i> , pages 1–55.		960
961			961
962			962
963			963

A Dataset Pipeline Illustrations

We illustrate the pipelines for creating V-FLUTE from each underlying figurative phenomenon in Figures 7, 2, 3, 4.

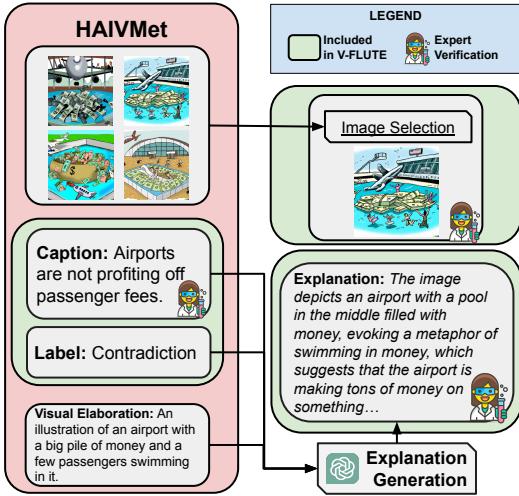


Figure 7: Creation of V-FLUTE instances for metaphors and similes from HAIVMet.

Type	Dataset	Train	Valid	Test
Metaphor /Similes	HAIVMET	649	107	101
	IRFL (metaphor /simile)	912	117	120
Idioms	IRFL (idiom)	170	100	100
Sarcasm	MuSE	830	106	106
Humor	MemeCap	1566	196	196
	NYCartoons	451	100	100
Total		4,578	726	723

Table 6: Data counts per phenomenon and dataset.

Length distribution Average length of a caption in V-FLUTE is ≈ 61 characters. Average length of an explanation is ≈ 367 characters. Figure 8 shows the distribution of caption lengths, and Figure 9 shows the distribution of explanation lengths by source dataset. We manually verified that the outlier instances are correct.

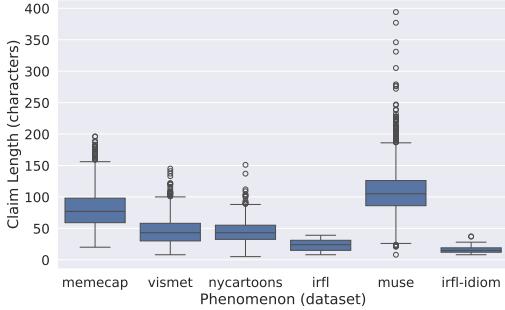


Figure 8: Distribution of lengths of captions by source dataset.

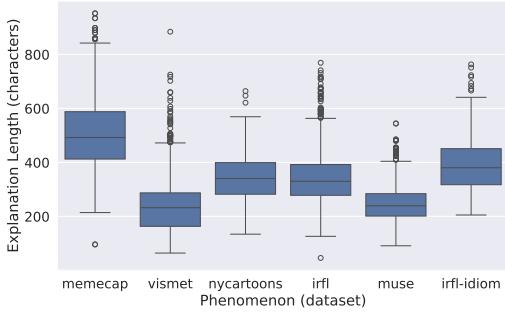


Figure 9: Distribution of lengths of explanations by source dataset.

1000	D API models Hyperparameters	1038
1001	D.1 Claude	1039
1002	• Model Name: claude-3-opus-20240229	1040
1003	• Max Tokens: 256	1041
1004	• Images greater than 5MB were resized maintaining aspect ratio	1042
1005		
1006	D.2 GPT-4	1043
1007	• Model Name: gpt-4-1106-vision-preview	1044
1008		
1009	• Max Tokens: 256	1045
1010	• Seed: 42	1046
1011	• Image URL detail: 'high'	1047
1012	D.3 Gemini	1048
1013	• Model Name: gemini-pro-vision	1049
1014		1050
1015	• Max Tokens: 256	1051
1016	• Safety Settings: 'BLOCK NONE'	1052
1017	• Images greater than 5MB were resized maintaining aspect ratio	1053
1018	E Fine-tuning Hyperparameters	1054
1019	LLava-v1.6-6B and 34B respectively utilize instruction-tuned LLMs as their backbone, Mistral-7BInstruct ⁵ and Yi-34B ⁶ .	1055
1020		1056
1021		
1022	We utilize LoRA (Hu et al., 2022) to fine-tune the models. We utilize the same hyperparameters for all fine-tunes outlined in Appendix E and use early stopping based on a V-FLUTE validation set to prevent overfitting. For evil and e-ViL+V-FLUTE we only fine-tuned for 2 epochs due to size of the e-ViL dataset and took the best checkpoint based on early stopping on V-FLUTE validation set. For eViL we only fine-tuned for 1 epoch to prevent overfitting. For VFLUTE, we trained for 3 epochs, and VFLUTE -Image for 10 epochs (to ensure performance does not increase even with larger number of epochs), for both we took the best checkpoint based on early stopping.	1057
1023		1058
1024		1059
1025		1060
1026		1061
1027		1062
1028		1063
1029		1064
1030		1065
1031		1066
1032		1067
1033		
1034		
1035		
1036		
1037		
	Fine-tuning	
	• Seed: 42	
	• Vision Tower: openai-clip-vit-large-patch14-336	
	• Number of Training Epochs: 3	
	• Train Batch Size (per device): 16	
	• Eval Batch Size (per device): 4	
	• Learning Rate: 2e-5	
	• Weight Decay: 0	
	• Warmup Ratio: 0.03	
	• Scheduler Type: cosine	
	• Number of epochs: 4 for eViL and eViL + vFLUTE, 10 for VFLUTE	
	• mm-projector-type: mlp2x gelu	
	• mm-vision-select-layer: -2	
	• mm-use-im-start-end: False	
	• mm-use-im-patch-token: False	
	• image-aspect-ratio: pad	
	• group-by-modality-length: False	
	LoRA	
	• lora r: 128	
	• lora alpha: 256	
	• mm-projector-lr: 2e-5	
	Deepspeed Configuration	
	• FP16 enabled: auto	
	• BF16 enabled: auto	
	• Micro Batch Size Per GPU: auto	
	• Train Batch Size: auto	
	• Gradient Accumulation Steps: auto	
	• Zero Optimization Stage: 3	

⁵huggingface.co/mistralai/Mistral-7B-Instruct-v0.1

⁶huggingface.co/NousResearch/Nous-Hermes-2-Yi-34B

1068	Training and Inference Instructions	1115
1069	All models are evaluated using beam search with	
1070	$n = 3$, temperature 0, max length 256. In the	
1071	case of generating scene graphs for the composi-	
1072	tional chain-of-thought method, we set the max	
1073	length to 256 for the graph generation step as rec-	
1074	ommended by Mitra et al. (2023). API models are	
1075	evaluated with default hyperparameters. We format	
1076	all fine-tuning data in the instruction format follow-	
1077	ing LLaVA (Liu et al., 2023a). To avoid overfitting	
1078	on a particular instruction for this task, we generate	
1079	20 similar instructions using an LLM (ChatGPT-4)	
1080	and randomly assign one of them to every instance	
1081	in the training, validation, and testing set. Same	
1082	instructions were sampled for the e-ViL dataset.	
1083	Table 7 shows the 20 instructions used.	
1084	The instructions were almost always followed. If	
1085	they were not followed during the data creation pro-	
1086	cess, we discarded those instances. For evaluation,	
1087	we looked at the sample outputs of each model and	
1088	designed rules to extract the label and the explana-	
1089	tion from the output, which was not too difficult	
1090	since mostly the instructions were followed well.	
1091	In the rare cases the model failed to follow instruc-	
1092	tions, that label would likely be incorrect.	
1093	Evaluation Hyperparameters	
1094	Following prior work, we utilize BERTScore	
1095	(Zhang* et al., 2020) based on the	
1096	microsoft-deberta-xlarge-mnli model (He	
1097	et al., 2021; Williams et al., 2018) and BLEURT	
1098	(Sellam et al., 2020) based on BLEURT-20 (Pu	
1099	et al., 2021) for the ExplanationScore.	
1100	F Prompts for LLMs	
1101	F.1 HAIVMET	
1102	F.1.1 One-shot Prompt for generating	
1103	explanations	
1104	We describe our one-shot prompts given to an LLM	
1105	(gpt-3.5-turbo-instruct-0914) for generating	
1106	explanations of entailment-contradiction relation-	
1107	ship. Refer to Table 8 for the detailed prompt.	
1108	F.2 IRFL	
1109	F.2.1 Zero-shot Prompt for generating	
1110	explanations	
1111	We provide our zero-shot prompt given to an LLM	
1112	(gpt-4-vision-preview) for generating the en-	
1113	tailment explanations given the claim and the im-	
1114	age. Refer Table 9 for the detailed prompt.	
1115	F.3 MuSE	
1116	F.3.1 Few-shot Prompt for generating	
1117	opposite claims	
1118	We provide our few-shot prompt given to an LLM	
1119	((gpt-4-0613)) for generating the opposite claims.	
1120	Refer Table 10 for the detailed prompt.	
1121	F.3.2 Zero-shot Prompt for Rephrasing	
1122	We provide our zero-shot prompt given to an LLM	
1123	(gpt-4-vision-preview) for rephrasing the ex-	
1124	planations given the claim and the crowd worker ex-	
1125	planation. Refer Table 11 for the detailed prompt.	
1126	F.4 MemeCap	
1127	F.4.1 Few-shot Prompt for generating	
1128	entailing claims	
1129	We describe our few-shot prompts given to an LLM	
1130	((gpt-4-0613)) for generating entailing captions as	
1131	part of the pipeline. Refer to Table 12 for the de-	
1132	tailed prompt.	
1133	F.4.2 Zero-shot Prompt for validating the	
1134	entailing captions	
1135	We describe our zero-shot prompt given to an	
1136	LLM (gpt-4-vision-preview) for validating the	
1137	claims generated in the previous step. Refer Table	
1138	13 for the detailed prompt.	
1139	F.4.3 Few-shot Prompt for generating	
1140	opposite claims	
1141	We provide our few-shot prompt given to an LLM	
1142	((gpt-4-0613)) for generating the opposite claims.	
1143	Refer Table 14 for the detailed prompt.	
1144	F.4.4 Zero-shot Prompt for generating	
1145	explanations	
1146	We provide our zero-shot prompt given to an LLM	
1147	(gpt-4-vision-preview) for generating the en-	
1148	tailment explanations given the claim and the im-	
1149	age. Refer Table 15 for the detailed prompt.	
1150	G Model Taxonomy	
1151	The taxonomy of all models used for automatic	
1152	evaluation is shown in Figure 10.	
1153	H Multimodal Strcutured	
1154	Chain-of-Thought Performance	
1155	In addition to zero-shot testing, we also test these	
1156	models using <i>Compositional Chain-of-Thought</i>	
1157	<i>Prompting</i> proposed by Mitra et al. (2023). The	
1158	method prompts the model <i>zero-shot</i> to generate a	
1159	scene graph in JSON format and then utilizes that	

No.	Instruction
1	Does the image's narrative confirm or disprove the claim REPLACE CLAIM? Discuss your reasoning and identify it as either entailment or contradiction.
2	Does this image confirm or deny the claim REPLACE CLAIM? Discuss your reasoning and determine a label: entailment or contradiction.
3	Is the image's message supporting or opposing the claim REPLACE CLAIM? Discuss your rationale and determine the appropriate label: entailment or contradiction.
4	Is there agreement or disagreement between the image and the claim REPLACE CLAIM? Provide your analysis and choose between entailment or contradiction.
5	Does the visual evidence support or counter the claim REPLACE CLAIM? Provide your explanation and assign it a label of entailment or contradiction.
6	Does the image agree with or dispute the claim REPLACE CLAIM? Explain your analysis and mark it as entailment or contradiction.
7	Does the illustration affirm or contest the claim REPLACE CLAIM? Provide your argument and choose a label: entailment or contradiction.
8	Is the visual content in agreement or disagreement with the claim REPLACE CLAIM? Offer your explanation and categorize it under entailment or contradiction.
9	Is the image in harmony with or in conflict with the statement REPLACE CLAIM? Explain your justification and label it as entailment or contradiction.
10	Is the portrayal in the image consistent with or contradictory to the claim REPLACE CLAIM? Offer your insights and select between entailment or contradiction.
11	Does the image's depiction validate or refute the claim REPLACE CLAIM? Explain your point of view and select a label: entailment or contradiction.
12	Is the content of the image endorsing or challenging the claim REPLACE CLAIM? Justify your position and label it as entailment or contradiction.
13	Is the image consistent with the statement REPLACE CLAIM? Justify your answer and classify it as either entailment or contradiction.
14	Does the illustration affirm or negate the claim REPLACE CLAIM? Articulate your reasoning and apply a label: entailment or contradiction.
15	Does the picture support or refute the assertion REPLACE CLAIM? Offer your rationale and select a label: entailment or contradiction.
16	Is the visual portrayal compatible with or adverse to the claim REPLACE CLAIM? Justify your viewpoint and label it as entailment or contradiction.
17	Does the image corroborate or dispute the claim REPLACE CLAIM? Outline your reasoning and categorize it under entailment or contradiction.
18	Is the depiction aligned with or against the claim REPLACE CLAIM? Share your evaluation and identify it as either entailment or contradiction.
19	Does the image entail or contradict the claim REPLACE CLAIM? Explain your reasoning and provide a label between entailment or contradiction.
20	Can the image be seen as validating or opposing the claim REPLACE CLAIM? Explain your thought process and assign a label of entailment or contradiction
21	Is the image's representation supportive of or contradictory to the claim REPLACE CLAIM? Articulate your analysis and assign the label: entailment or contradiction.

Table 7: Instruction variations for the figurative visual entailment task.

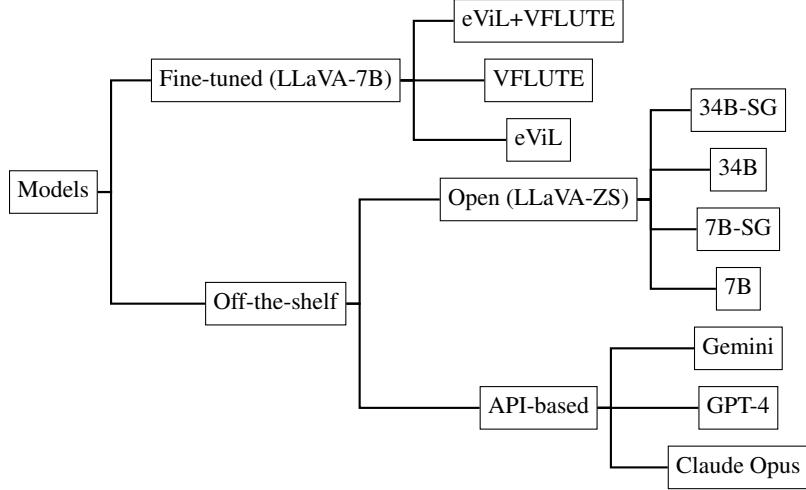


Figure 10: Taxonomy of models used for the study.

You will be provided a Caption describing what is in the image in detail. You will also be provided with a Claim that contradicts or is entailed by the image (as indicated by the Label). Your task is to explain why the claim contradicts or is entailed by the image. Be very brief in your explanation. Start your explanation by describing what the image depicts, displays or shows.
Caption: An illustration of a group of soldiers with red skin, horns, and pitchforks in hand with a fierce expression on their faces.
Claim: The soldiers were angels.
Label: Contradiction
Explanation: The image depicts soldiers with red skin, horns, and pitchforks, which are traditional characteristics associated with demons, not angels. Therefore, the claim that the soldiers were angels contradicts the image.
Caption:.....

Table 8: One shot prompt given to an LLM (gpt-3.5-turbo-instruct-0914) for generating explanations of entailment-contradiction relationship of the HAIVMET dataset.

You will be provided an image. You will also be provided with a **simile** that contradicts or is entailed by the image (as indicated by the Label). Your task is to explain why the **simile** contradicts or is entailed by the image. Be very brief in your explanation and remain consistent to the Label in your explanation. Start your explanation by describing what the image depicts, displays or shows.
Simile:
Label:
Explanation:

Table 9: Zero shot prompt given to an LLM (gpt-4-vision-preview) for generating explanations of entailment-contradiction relationship of the IRFL Dataset. The dataset contains similes, metaphors and idioms. For metaphors and idioms, the word simile in the prompt is replaced with the corresponding type.

u are an online redditor or flickr user and u always type in informal style. Convert the following sarcastic claim into a non-sarcastic claim. Preserve the informal style, including capitalization. Be super laid back and informal!!!
1. Sarcastic claim: stairs vs . escalator in airport . i wonder why we have an # obesity problem ? # publichealth # ncds # globalhealth # isometimesdothistoo
Explanation: no wonder we have an obesity problem since everyones using escalator instead of stairs in airport.
Non-sarcastic claim: **it s clear why we have an # obesity problem look at stairs vs. escalator in airport**
Claim:.....
Explanation:.....

Table 10: Few shot prompt given to an LLM (gpt-4-0613) for generating opposite claims utilizing the sarcastic claim and crowd worker explanation.

Paraphrase the draft explanation of why the image contradicts the literal interpretation of the claim. Be sure to first describe the image in one sentence. Keep your answer short. Do not refer to the claim or the draft explanation in your paraphrase. Stay close to the draft explanation.
Claim:
Draft Explanation:

Table 11: Zero shot prompt given to an LLM (gpt-4-vision-preview) for rephrasing the explanations given the claim and the.

You will be provided with a meme caption. Your task is to write the meme caption as a claim such that the meme poster is not mentioned in the claim.	1160
Caption: Meme poster is saying that searching Google plus the term you want to search on reddit is better than searching reddit itself.	1161
Claim: Searching on Google with the term you want to search plus 'reddit' is more effective than searching directly on Reddit.	1162
Caption: The person who wrote the post is saying people on Instagram are soft and reddit are funny.	1163
Claim: People on Instagram are soft, whereas those on Reddit are funny.	1164
Caption:.....	1165

Table 12: Two shot prompt given to an LLM (gpt-4-0613) for generating entailing claims utilizing the meme captions part of the MemeCap dataset.

You will be provided a meme image and a claim. Your task is to check whether the claim entails the image. Answer with a Yes or No.	1166
Claim:	1167

Table 13: Zero shot prompt given to an LLM (gpt-4-vision-preview) for validating the claims generated in F.4.1. The corresponding meme image is also attached with the prompt.

Claim: A useful feature has been removed on YouTube, causing disappointment.	1168
Explanation: The image shows a painting of a character with a distraught face and a speech bubble that reads "y tho," placed over text saying "When YouTube removed sort by oldest option." This implies that the removal of the sort by oldest option is a decision that users are questioning, hence indicating disappointment over the loss of a useful feature.	1169
Opposite claim: An unhelpful feature has been removed on YouTube, causing happiness.	1170
Claim:.....	1171
Explanation:.....	1172

Table 14: Few shot prompt given to an LLM (gpt-4-0613) for generating opposite claims utilizing the generated claim and explanation.

You will be provided a meme. You will also be provided with a claim that entails the image. Your task is to explain why the claim is entailed by the image. Be very brief in your explanation and start your explanation by describing what the image depicts, displays or shows.	1173
Claim:	1174
Explanation:	1175

Table 15: Zero shot prompt given to an LLM (gpt-4-vision-preview) for generating the entailment explanations. The corresponding meme image is also attached with the prompt.

scene graph in another prompt to answer the relevant question. We refer to these models as LLaVA-ZS-7B-SG and LLaVA-ZS-34B-SG for the 7B and 34B LLaVA configurations described above.	1176
Scene graph prompting and few-shot prompting improves performance on the figurative visual entailment task. By observing results in Table 3, we can see that the multimodal few-shot prompting and scene graph prompting, having demonstrated their effectiveness for literal inputs, also show improved performance on the figurative visual entailment task. However, the explanations generated by SG-models tend to overly focus on the contents of the scene graph rather than the underlying figurative phenomena, possibly causing a decrease in explanation score.	1177

I By-Phenomenon Performance

In Figure 11, we show the performance of the models by phenomenon and dataset across various thresholds.	1178
	1179

J Human Baseline

To find out how humans perform on the task, we hire two expert annotators with formal education in linguistics. We present them with 10 example instances and then ask them to complete 99 randomly sampled test set instances. We also evaluate our best model (see Table 3) on the same set. Results are shown in Table 16. Human performance is quite strong, almost reaching 90 F1@0 score overall. Human performance is better than our strongest fine-tuned model (LLaVA-7B-eVil+VF) performance with $p < 0.05$ for Annotator 1 and $p < 0.07$ for Annotator 2. Humans excel at interpreting memes, with both annotators reaching a 100% F1 score. Humans also perform noticeably better on the NYCartoons dataset and on the idiom subset of the task. The model has a slight	1180
---	------

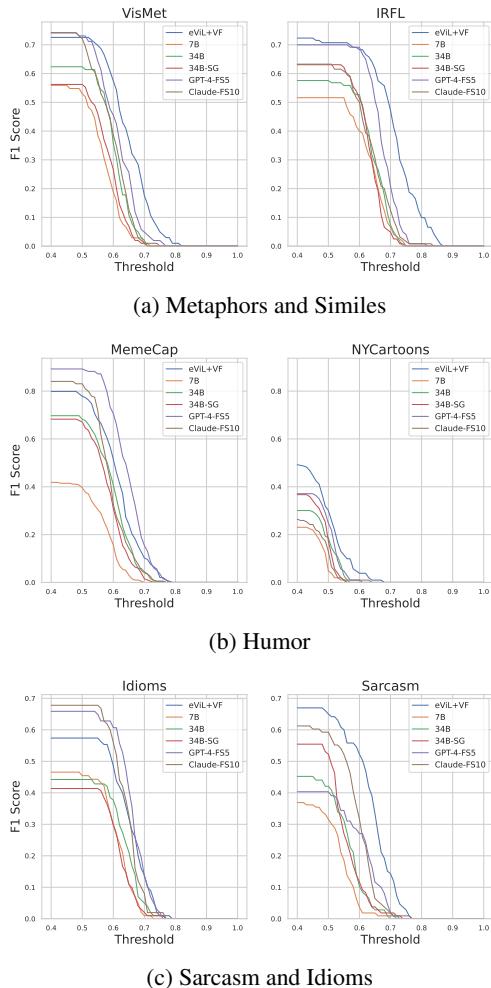


Figure 11: Performance of the models by phenomenon.

edge in performance on the sarcasm and visual metaphor subsets of the task, perhaps due to difficulty of these subsets and any potential spurious correlations during fine-tuning.

Phenomenon	Dataset	Human Avg	LLaVA-eViL+VF
Metaphor /Similes	HAIVMET	78.84	81.25
	IRFL (metaphor /simile)	94.36	77.78
Idioms	IRFL (idiom)	89.26	49.74
Sarcasm	MuSE	68.89	85.42
Humor	MemeCap	100.0	78.03
	NYCartoons	71.43	47.83
Overall		89.09	77.26

Table 16: Human baseline results (F1@0) by phenomenon and source dataset.

K How Well Does the Explanation Score Predict Human Judgment on Adequacy?

We explore whether the proposed explanation score can capture human judgement of explanation adequacy. We collect all instances from Section 5 where both annotators agreed on the adequacy judgement for the explanation. We evaluate if the explanation score described in Section 4.2 can act as a good predictor for the human adequacy judgment. We find that the area under the Precision-Recall curve is 0.79, and the maximum F1 score is 0.77, obtainable at the explanation score threshold of 0.53. Hence, we use this threshold to report the results in Table 3. We also use the threshold of 0.6 since it maximizes F1 such that both precision and recall are above 0.75.

L How Do Models Perform When Only Predicting the Label?

In our experiments, we found that predicting only the label improves accuracy compared to predicting label and explanation (this is expected and observed in other work on textual explanations such as e-SNLI (Camburu et al., 2018)). However, these predictions are less reliable since they could be due to spurious correlations (which is why we require the model to generate textual explanations). We also found when fine-tuning the model in a multi-task fashion with explanations (i.e., two tasks,

1230 one of generating explanations and one of predict-
1231 ing the label), the accuracy improves compared to
1232 when fine-tuning only for the prediction task (F1
1233 score of 80.85 vs. 83.26, $p < 0.1$), in line with
1234 previous findings by [Hsieh et al. \(2023\)](#).

1235 **M Annotation Interfaces**

1236 We provide the annotation interfaces below for
1237 HAIVMET (Figure 12), IRFL (Figure 13), Meme-
1238 Cap (Figure 14) and MuSE (Figure 15). In addition,
1239 instructions were explained in more detail to the
1240 annotators via chat on Upwork, and any of their
1241 doubts and questions were answered.

Welcome to Our Survey!

We are a team of researchers in computational linguistics exploring the visual entailment of image-text pairs. This survey involves three key tasks:

1. **Image Selection:** From a set of up to four images, select the one that best corresponds with the given caption. If none of them match with the given caption, select "None of the above" option. Note: In case of no match, select anything for the below two tasks.
2. **Label Verification:** Considering the caption, claim, and your chosen image, determine whether the label ('entailment' or 'contradiction') is accurate. Entailment refers to the claim and image being the same thing while Contradiction is them being different.
3. **Explanation Assessment:**
 - If the label is correct, evaluate whether the provided explanation is also correct. If not, suggest an appropriate explanation.
 - If the label is incorrect, provide a suitable explanation for a more fitting label.

Base your preference on correctness (Is the explanation correct? Is it consistent with the label?) and completeness (does the explanation mention everything I would want it to?). Please prioritize correctness.

Enter your username:

(Please use the same username throughout the study)

Caption: {{caption}}

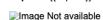
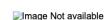
Image Not available

Image Not available

Image Not available

Image Not available

Select the output you prefer:

Output 1 Output 2 Output 3 Output 4 None of the above

Claim: {{claim}}

Label: {{label}}

Is the label correct?

Yes
 No

Explanation: {{explanation}}

Is the explanation correct and complete?

Yes
 No

If the label or the given explanation is incorrect, please provide the correct and complete explanation:

Sample No.: {{folder_name}}/{{og_max_len}}

<Previous Save Next>

Figure 12: Annotation interface for HAIVMET.

Welcome to Our Survey!

We are a team of researchers in computational linguistics exploring the visual entailment of image-text pairs. This survey involves the following task:

Explanation Assessment: You will be given the following:

- An image.
- A claim about the image which can be a Metaphor, Simile or an Idiom.
- A label ('entailment' or 'contradiction'). Entailment refers to the claim and image being the same thing while Contradiction is them being different.
- An explanation for the claim

Evaluate whether the provided explanation is correct. If not, suggest an appropriate explanation.

Base your preference on correctness (Is the explanation correct? Is it consistent with the label?) and completeness (does the explanation mention everything I would want it to?). Please prioritize correctness.

Enter your username:

(Please use the same username throughout the study)

Image Not available

Claim: {{claim}}

Label: {{label}}

Explanation: {{explanation}}

Is the explanation correct and complete?

Yes
 No

If the given explanation is incorrect, please provide the correct and complete explanation:

Sample No.: {{folder_name}}/{{og_max_len}}

<Previous Save Next>

Figure 13: Annotation interface for IRFL.

Welcome to Our Survey!

We are a team of researchers in computational linguistics exploring the visual entailment of image-text pairs. This survey involves the following task:

Explanation Assessment: You will be given the following:

- An image
- A claim about the meme in the image.
- An explanation for the claim entailing the image. Here entailment means the claim logically follows from the image, or they make the same statement.

First, evaluate whether the claim fits the image. This means that the claim is indeed entailed by the image and has correct grammar and lacks any other issues. If not, please provide a correction.

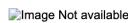
Next, evaluate whether the provided explanation is:

- **correct** (does the explanation accurately describe the image? Is it logical? Does it follow common sense?),
- **complete** (does the explanation mention everything I would want it to?), and
- **concise** (can I say the same thing using less words?)

If it is lacking, please suggest an appropriate explanation. Focus on brevity and correctness, no need to add missing detail if it is not relevant to the entailment between the claim and the image.

Enter your username:

(Please use the same username throughout the study)



Claim: {{claim}}

Explanation: {{explanation}}

Does the claim fit the image?

- Yes
 No

Is the explanation correct, complete, and concise?

- Yes
 No

Sample No.: {{index + 1}}/{{og_max_len}}

<Previous Save Next>

Figure 14: Annotation interface for MemeCap.

Welcome to Our Survey!

We are a team of researchers in computational linguistics exploring the visual entailment of image-text pairs. This survey involves the following task:

Explanation Assessment: You will be given the following:

- An image
- Two claims related to the image. One of them is sarcastic, the other one is literal.
- An explanation for the claim entailing or contradicting the image. Here entailment means the claim logically follows from the image or the combination of claim and image is true. So, for a sarcastic claim, the relationship is always a contradiction. For a literal claim, the relationship should always be entailment.

First, evaluate whether to discard the sample . You should discard the sample if the image does not relate to the claim or the claim does not make any sense.

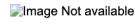
Next, evaluate whether the provided explanation is:

- **correct** (does the explanation accurately describe the image? Is it logical? Does it follow common sense?),
- **complete** (does the explanation mention everything I would want it to?), and
- **concise** (can I say the same thing using less words?)

The explanation should make sense for both of the claims. If it is lacking, please suggest an appropriate explanation. Focus on brevity and correctness, no need to add missing detail if it is not relevant to the entailment between the claim and the image.

Enter your username:

(Please use the same username throughout the study)



Sarcastic Claim (Contradiction): {{claim}}

Literal Claim (Entailment): {{contra_claim}}

Explanation (reference): {{expl_ref}}

Explanation (final): {{expl_fin}}

Discard Sample Options:

- Do not discard
 Image does not contradict the sarcastic claim
 Image does not entail the literal claim
 Other (please explain below)
Other:

Is the explanation correct, complete, and concise?

- Yes
 No

Sample No.: {{index + 1}}/{{og_max_len}}

<Previous Save Next>

Figure 15: Annotation interface for MuSE.