

# MultiCaption: Detecting disinformation using multilingual visual claims

Rafael Martins Frade<sup>1,2\*</sup>, Rrubaa Panchendrarajan<sup>3\*</sup>, Arkaitz Zubiaga<sup>3</sup>

<sup>1</sup>University of Santiago de Compostela, Spain

<sup>2</sup>Newtral Media Audiovisual, Spain

<sup>3</sup>Queen Mary University of London, United Kingdom

mfrade.rafael@gmail.com, {r.panchendrarajan, a.zubiaga}@qmul.ac.uk

## Abstract

Online disinformation poses an escalating threat to society, driven increasingly by the rapid spread of misleading content across both multimedia and multilingual platforms. While automated fact-checking methods have advanced in recent years, their effectiveness remains constrained by the scarcity of datasets that reflect these real-world complexities. To address this gap, we first present *MultiCaption*, a new dataset specifically designed for detecting contradictions in visual claims. Pairs of claims referring to the same image or video were labeled through multiple strategies to determine whether they contradict each other. The resulting dataset comprises 11,088 visual claims in 64 languages, offering a unique resource for building and evaluating misinformation-detection systems in truly multimodal and multilingual environments. We then provide comprehensive experiments using transformer-based architectures, natural language inference models, and large language models, establishing strong baselines for future research. The results show that *MultiCaption* is more challenging than standard NLI tasks, requiring task-specific finetuning for strong performance. Moreover, the gains from multilingual training and testing highlight the dataset’s potential for building effective multilingual fact-checking pipelines without relying on machine translation.

## Introduction

The rapid spread of misinformation online has its impacts identified in many aspects of society. Although automated fact-checking pipelines aim to mitigate this problem, the task remains highly challenging due to the complexity and diversity of misinformation. In particular, misinformation disseminated across multimedia content and multiple languages requires more generalizable and robust solutions. Figure 1 illustrates this challenge: it shows an image with its original caption describing a 2016 incident, alongside false captions in English and Urdu that repurposed the same image in connection with an unrelated event in 2023.

While early research in automated fact-checking focused primarily on textual claims (Thorne et al. 2018; Elsayed et al. 2019), recent work has expanded to incorporate images, videos, and other multimodal and multilingual content



This picture released on May 25, 2016 by the Italian Navy shows the shipwreck of an overcrowded boat of migrants off the Libyan coast today

GREECE BOAT DISASTER AND THE HUMAN TRAFFICKERS IN PAKISTAN...

یونان کشتی حادثہ کیسے ہوا  
(Translation: How did the Greek boat accident happen?)

Over 300 Pakistanis have lost their lives in #Greece but our media is still suffering...

Figure 1: An example of image circulating in social media, with true (green) and false captions (red) (afp 2025)

(Akhtar et al. 2023; Chakraborty et al. 2023; Aneja, Bregler, and Nießner 2023). However, datasets specifically designed to address misinformation spread through images or videos remain scarce, and most existing resources focus on textual claim verification (Thorne et al. 2018; Aly et al. 2021).

A closely related task aimed at identifying multimodal misinformation is out-of-context detection. Given an image and two associated claims, the objective is to determine whether one of the claims presents the image out of context. The dataset commonly used for this task is *COSMOS* (Aneja, Bregler, and Nießner 2023), which was primarily constructed from news articles and fact-checking websites. However, this dataset has certain limitations: it contains only claims in English and includes a relatively small number of instances that reflect claims used to spread disinformation in real-world scenarios.

To fill this gap, we first introduce *MultiCaption*, a multilingual dataset designed to identify disinformation through contradictory visual claims. We define two visual claims as contradictory if they cannot both be true simultaneously about the same image or video (Sepúlveda-Torres, Bonet-Jover, and Saquete 2023). We build our dataset on *Multi-*

\*These authors contributed equally.

*Claim* (Pikuliak et al. 2023) and *MultiClaimNet* (Panchendrarajan, Míguez, and Zubiaga 2025) as primary data sources, which consists of claims written by professional fact-checkers and based on social media content used to disseminate real-world misinformation across multiple countries and languages. We employ multiple labeling strategies, including manual validation, annotation using large language models (LLMs) and methods leveraging claim links from the original sources to label visual claim pairs as contradictory or non-contradictory. We apply a rigorous filtration process to extract high-quality visual claim pairs from the sources, constructing a dataset that reflects real-world challenges. This process yields 11,088 claim pairs across 64 languages, encompassing both monolingual and cross-lingual pairs, making *MultiCaption* a valuable resource for advancing research in multilingual disinformation detection.

Using the *MultiCaption* dataset, we then perform extensive experiments, carefully constructing a disjoint test set to prevent data leakage and applying data expansion strategies on the training set to enable robust fine-tuning. We provide benchmark results from fine-tuned transformers, natural language inference (NLI) models, and LLMs on both the *MultiCaption* test set and *COSMOS*, in monolingual and multilingual settings. Our results show that *MultiCaption* introduces a more challenging benchmark than *COSMOS*. Moreover, most models benefit significantly from the fine-tuning process, particularly in the multilingual setting, highlighting the value of our dataset as a resource for advancing multilingual misinformation detection.

We make the following contributions:

- We introduce *MultiCaption*<sup>1</sup>, the first multilingual dataset for detecting contradictory visual claims, containing 11,088 claim pairs across 64 languages.
- We carefully curate multiple labeling strategies, including both manual and LLM-based methods, to label visual claim pairs as contradictory or non-contradictory.
- We conduct extensive experiments to establish robust baselines in a multilingual setting and explore strategies for further expanding training data.

We believe our dataset and experiments will serve as a valuable resource and benchmarking for developing effective contradiction detection systems to combat disinformation in multilingual settings.

## Related Work

We provide an overview of existing research and datasets relevant to contradiction detection, as well as closely related tasks such as natural language inference (NLI), out-of-context detection, and claim matching.

**Contradiction Detection & NLI** Identifying contradictions between two pieces of text is typically framed as a NLI task (Li, Qin, and Liu 2017). More fine-grained task definitions address contradiction-type classification (Senouci,

Meziane, and Benbernou 2025), detecting contradictions between a query and a document (Xu et al. 2024), identifying contradictions within a single document (Li, Raheja, and Kumar 2024), and cross-modal contradictions between images and text (Popordanoska, Li, and Blaschko 2025).

The NLI task involves identifying entailment, neutrality or contradiction between two sentences (Gubelmann et al. 2024). NLI techniques are widely used in disinformation research, especially for claim verification. In this context, labels are often adapted from the standard NLI schema (CONTRADICTION, ENTAIL, NEUTRAL) to task-specific ones such as SUPPORTED, REFUTED, and NOT ENOUGH INFO (Thorne et al. 2018). As in many other natural language processing tasks, claim verification has recently shifted toward leveraging LLMs (Ahmad, Usmanova, and Rehm 2025; Dmonte et al. 2024). Two common datasets used in this task are FEVER (Thorne et al. 2018) and FEVEROUS (Aly et al. 2021). Both contain claims accompanied by textual and tabular evidence that supports or refutes each claim. A limitation of those datasets is that the data is only in English and was generated based on Wikipedia, lacking semantic characteristics commonly present on social media language.

**Out-of-context (OOC) Detection** Despite growing focus on deepfakes, a common form of disinformation involves real images or videos misrepresented through false context or claims. Closely related to the automated fact-checking pipeline, OOC detection evaluates the truthfulness of a claim given an image or video. The task typically assumes access to the original context as a triplet (*Candidate Claim*, *Original Claim*, *Image*) and aims to determine entailment or contradiction between the claims, optionally incorporating visual information. Most work in this setting relies on the *COSMOS* dataset and explores approaches such as context retrieval, text similarity, contradiction detection, entity linking, LLMs, and cross-modal consistency (Aneja, Bregler, and Nießner 2023; Abdelnabi, Hasan, and Fritz 2022; Nguyen and Tran 2023; La et al. 2022b; Tran et al. 2022; Nguyen, Suganuma, and Okatani 2022; La et al. 2022a).

Despite introducing the task and shaping research on the topic, the *COSMOS* dataset has two main limitations: It contains only claims in English and includes a very small labeled dataset (1,700 samples), containing only a limited number of real disinformation cases. Another limitation is that its OOC class contains a substantial amount of direct negations or linguistic features in candidate claims suggesting that they are fact-checks, making it less useful to develop solutions for real-life disinformation detection.

Another line of work assumes that the original context is unavailable and attempts to determine whether an image or video is presented out of context using only the claim-image pair. This approach relies on datasets such as NewsClipping (Luo, Darrell, and Rohrbach 2021) and VERITE (Papadopoulos et al. 2024), and typically applies multimodal feature extraction and classification (Papadopoulos et al. 2025). However, in many real-world social media scenarios, the veracity of accompanying claims cannot be reliably assessed without access to the original context; Figure 1 shows an example illustrating this limitation.

<sup>1</sup>Dataset is available at <https://doi.org/10.5281/zenodo.18230659>, and source code is available at <https://github.com/rfrade/multicaption>

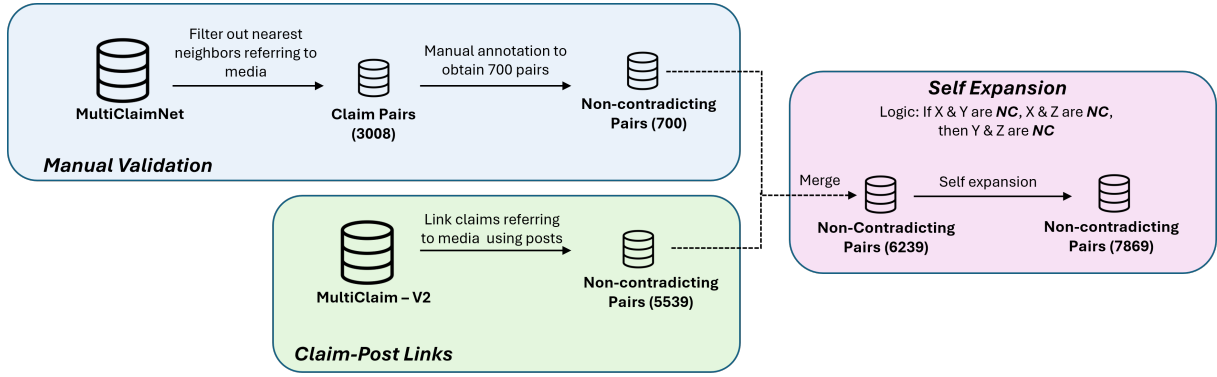


Figure 2: Work flow of generating non-contradicting (NC) pairs

**Claim Matching** Claim matching is a fundamental component of automated fact-checking pipelines. The task involves determining whether two claims convey similar information and can be extended to matching an unverified claim against a database of verified claims, a process often referred to as fact-checked claim retrieval (Panchendrarajan and Zubiaga 2024). While claim matching is framed as a binary task—deciding whether two claims are similar or not (Larraz, Míguez, and Sallicati 2023)—non-similarity does not necessarily imply contradiction. Claim matching typically focuses on fine-tuning multilingual transformers such as XLM-R (Larraz, Míguez, and Sallicati 2023) or LLMs (Pisarevskaya and Zubiaga 2025). While several claim matching datasets exist (Nielsen and McConville 2022; Singh et al. 2023; Pikuliak et al. 2023), they are designed for similarity detection and cannot be used to detect contradictions without significant adaptation.

## MultiCaption Dataset Construction

**Definition of Contradicting Visual Claims:** We define a pair of claims referring to the same image or video (referred to as visual claims) as *contradictory* if both cannot be true at the same time (Sepúlveda-Torres, Bonet-Jover, and Saquete 2023). A claim being less specific or only partially correct does not constitute a contradiction.

**Source of Visual Claims:** We use the *MultiClaim* (Pikuliak et al. 2023) dataset as the main source of multilingual claims referring images or videos for constructing the contradicting and non-contradicting pairs. The dataset comprises fact-checked claims along with references to the fact-checking articles and social media posts discussing these claims. The dataset was released in two versions (*MultiClaim* v1 and v2), with the final version containing 435K fact-checked claims and 89K linked social media posts. Each social media post is linked to at least one claim, resulting in a total of 105K claim–post links. The initial stages of our dataset construction were based on *MultiClaim* v1, while later stages used *MultiClaim* v2. We used the terms [“picture”, “image”, “photo”, “photograph”, “footage”, “document”, “video”, “clip”, “post”] combined with the verb “show” (e.g., photo shows) or preposition “of” (e.g., video

of) present in the English translation of the text to obtain the visual claims.

We employ a combination of manual and automatic methods to label candidate claim pairs, with the workflow for generating contradicting and non-contradicting pairs detailed below.

## Generating Non-Contradicting Pairs

We employ three labeling strategies to generate non-contradicting pairs—manual validation, claim–post links, and self-expansion—as shown in Figure 2.

**Generation using Manual Validation** Manual validation requires curating a set of visual claim pairs that are most likely to be non-contradicting to each other. We use *MultiClaimNet* (Panchendrarajan, Míguez, and Zubiaga 2025), derived from *MultiClaim* v1, which contains claim pairs constructed using a nearest-neighbor approach and annotated for similarity with large language models (LLMs). In the original work, the authors further grouped similar claims into clusters. From *MultiClaimNet*, we filtered out pairs in which at least one claim refers to an image or video, resulting in 3,008 similar claim pairs. However, as the authors relied solely on textual information for similarity annotation, these pairs do not necessarily refer to the same image or video. Therefore, our manual verification involved checking whether both claims refer to the same image or video, and ensuring that the claims are not contradictory. The original media was obtained from the fact-checking articles linked to each claim. Figure 3 shows an example of manually validated claim pairs ( $C_1, C_2$ ).

Since the claim pairs in *MultiClaimNet* were generated using a nearest-neighbor approach, we observed a similarity bias in the sample, where most pairs exhibit high similarity (ranging from 0.52 to 1) with an average of 0.86. Claim similarity was measured using cosine similarity between embeddings of their English translations, generated with Sentence Transformer (Reimers and Gurevych 2019). To mitigate this bias, we verified claim pairs in order of increasing similarity—starting from the least similar—until we obtained 700 valid non-contradicting pairs. The number of valid pairs was selected to closely match the number of manually validated

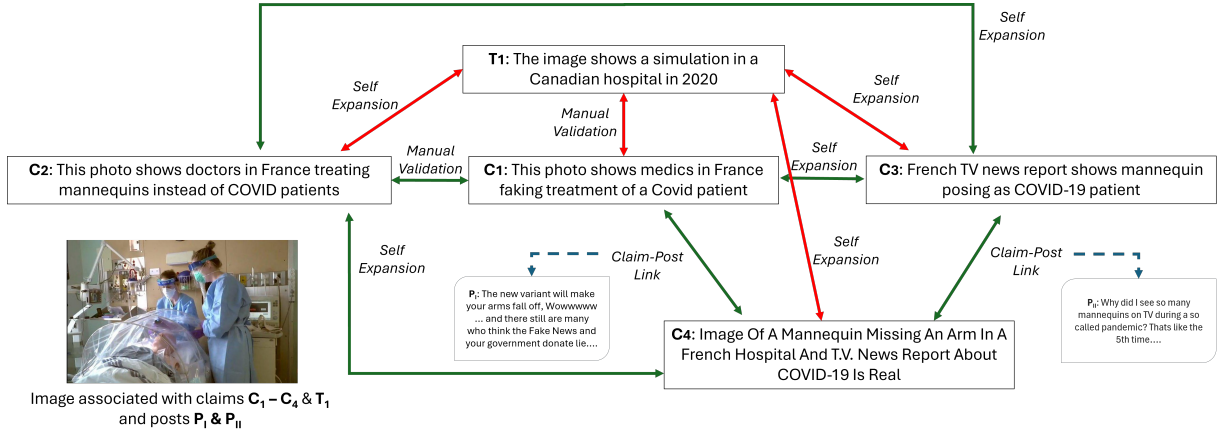


Figure 3: Example illustrating different annotation strategies in MultiCaption. Arrow colors encode relationship types: green for non-contradicting, red for contradicting, and blue for associated posts

contradicting pairs. To reach this target, we manually validated 1,144 claim pairs from the original sample. A substantial portion (approximately 33%) was discarded because the claims referred to different media, while an additional 5% were excluded due to missing fact-checking articles. Only a small fraction of pairs (less than 1%) were found to be contradictory within the *MultiClaimNet* sample.

**Generation using Claim-Post Links** *MultiClaim* dataset comprises claims and social media posts linked if the post discusses the claim. This relationship between claims and posts can be used to link similar claims referring to the same post. Therefore, we create further non-contradicting pairs when two visual claims share a post. Figure 3 shows claim pairs  $(C_1, C_4)$  and  $(C_3, C_4)$  automatically labeled non-contradicting via posts  $P_I$  and  $P_{II}$  respectively.

The latest version of *MultiClaim* comprises 105K claim-post links, created using multiple strategies, including back-linking, claim review schemas, and identical claims (Pikuliak et al. 2023). We excluded links generated through identical claims, as they inherently produce obvious non-contradicting pairs. Further, linking two claims via a post requires the post to be associated with at least two claims. Applying this constraint yielded 10.6K potential links, corresponding to 16.2K claim pairs. We retained only the pairs in which at least one claim referred to an image or video, resulting in 6,182 non-contradicting claim pairs.

Since these pairs were generated automatically, we manually inspected samples from both the lowest- and highest-similarity ranges at 0.05 similarity intervals. During this analysis, we observed that pairs with very low similarity scores ( $< 0.4$ ) were predominantly noisy, while those with very high similarity scores ( $> 0.95$ ) often contained identical claims. To eliminate noise and trivial examples, we discarded claim pairs from both extremes of the similarity distribution. Highly similar pairs ( $> 0.95$ ) were removed only when both claims were in the same language. Multilingual pairs were retained, as they contribute valuable cross-lingual context to the dataset despite their high similarity. Refer to the Appendix for the similarity distribution of the claim pairs

Label	Labeling Strategy	# Pairs
Contradiction	Manual Validation	722
	LLM Annotation	2072
	Self-Expansion	425
	<b>Total</b>	<b>3219</b>
Non-Contradiction	Manual Validation	700
	Claim-Pair Link	5508
	Self-Expansion	1661
	<b>Total</b>	<b>7869</b>

Table 1: Labeling strategies and No. of pairs in *MultiCaption*.

created via claim-post links and the discarded regions.

**Generation using Self-Expansion** Similar to the claim-post links discussed previously, there also exist direct links between claims. We identified that if claim  $X$  and claim  $Y$  are labeled as non-contradicting, and claim  $X$  and claim  $Z$  are also labeled as non-contradicting, then claim  $Y$  and  $Z$  can be automatically inferred to be non-contradicting. We refer to this process as self-expansion, through which additional non-contradicting pairs are automatically generated. This process yielded 1,768 new claim pairs labeled as non-contradicting via self-expansion. Figure 3 illustrates examples, including  $(C_1, C_3)$  (via  $C_4$ ),  $(C_2, C_3)$ , and  $(C_2, C_4)$  (via  $C_1$ ), all automatically labeled non-contradicting.

Together, the three labeling techniques—manual validation, claim-post linking, and self-expansion—produced a total of 7,869 non-contradicting claim pairs, as summarized in Table 1.

### Generating Contradicting Pairs

As mentioned earlier, *MultiClaim* contains claims and fact-checking articles written by professional fact-checkers. Since every claim is associated with at least one fact-checking article, we used the title of the fact-checking article as a source of contradicting claims. Similar to non-contradicting pairs, we generate contradicting pairs using



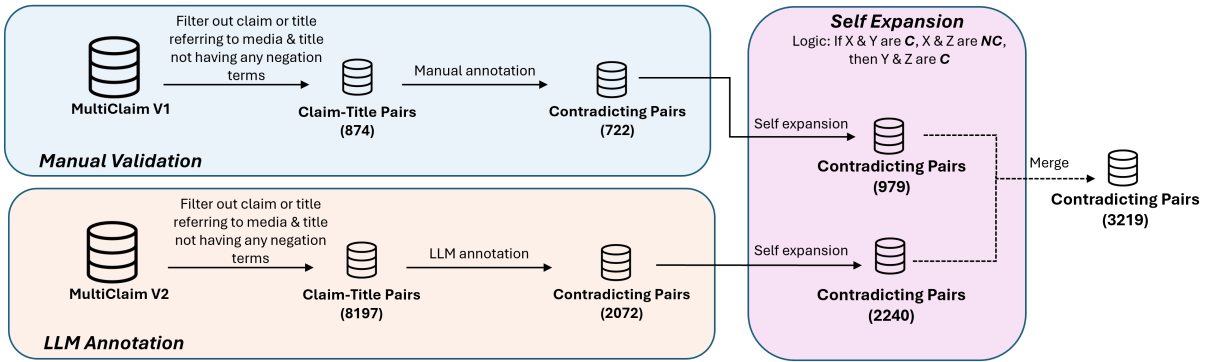


Figure 4: Work flow of generating contradicting (C) pairs

various labeling strategies—including manual validation, LLM annotation, and self-expansion as shown in Figure 4.

Fact-checking article titles typically follow three common patterns: (1) rewriting the claim as a question, (2) explicitly negating the claim, or (3) providing the original context of an image or video. The first pattern is not directly useful for our purposes, and the second—an explicit negation—would be trivially detectable in an automated fact-checking pipeline. We therefore focus on the third pattern, in which the title reveals the true context of the referenced image or video. The filtration aimed to create a dataset simulating real-world fact-checking, given an original context of an image or video, an automated agent determines whether a claim presents it out of context.

**Generation using Manual Validation** Based on claim-title pairs from *MultiClaim* v1, we first retained only the pairs in which either the claim or the title referred to a media item. We then removed titles that contained explicit negations or direct references to the claim. For example, for a claim such as "Image of the protest in France 2025", we excluded titles like "This image is not from France" or "Image of protest actually dates back to 2020". The exclusion of such cases represents a key advantage of our dataset relative to COSMOS. Refer to Table 5 in Appendix for negation terms used to filter out title with direct reference or denial of the claim. We then manually validated this sample to retain only the titles that did not contain any additional form of negation and did not make a direct reference to the claim, resulting in 722 contradicting claim-title pairs. The objective was to end up with pairs in which both claim and title were independent claims about the image/video. Figure 3 shows an example of manually validated claim-title pair  $(C_1, T_1)$ .

**Generation using Self-Expansion** Similar to the self-expansion technique used to generate non-contradicting claim pairs, the inherent relationships between contradicting and non-contradicting claims can be leveraged to automatically generate additional contradicting samples. Specifically, we also identified that if claim X contradicts claim Y, and claim X is non-contradicting with claim Z, then claim Y and claim Z are also contradicting. Figure 3 illustrates an example of this expansion: self-expanded contradicting pairs  $(T_1, C_2)$ ,  $(T_1, C_3)$ , and  $(T_1, C_4)$  are generated

as  $(T_1, C_1)$  is manually labeled as contradicting and  $C_1$  is non-contradicting with  $C_2 - C_4$ . Using this method, 722 manually validated contradicting pairs were expanded to a total of 979 pairs.

**Generation using LLM-Annotation** The number of contradictory pairs obtained through manual validation and self-expansion was relatively small compared with the non-contradictory pairs. Nevertheless, the *MultiClaim* v2 contains many potential contradicting claim-title pairs. Since manually validating all candidate pairs is infeasible, we leverage the LLM GPT-5 to perform validation as a proxy for manual review.

We applied the same filtering process described in Section to extract potential contradicting claim-title pairs from the *MultiClaim* v2 dataset. This procedure resulted in 8,197 claim-title pairs. For LLM annotation, we constructed the prompt to closely mimic the manual validation process. Although a claim and a title may appear contradictory, some titles directly reference the claim while explicitly denying it. To handle this, we instructed the LLM to label each claim-title pair as: *contradict* if both cannot be true for the same image/video, and *denial* if the title explicitly refutes or debunks, even with added context. Only claim-title pairs that are labeled as *contradict* but not *denial* are considered valid contradictory pairs.

To reduce potential noise from LLM annotations, we evaluated the approach using 874 claim-title pairs that were manually validated and labeled as either valid or invalid contradictory pairs. We iteratively refined the annotation prompt with clearer instructions and illustrative examples to guide accurate predictions. Once the prompt is finalized, we re-labeled the same 874 manually validated pairs using LLM to quantitatively assess its precision in identifying valid contradicting claim-title pairs. We measured the precision in deciding a pair as valid (*contradict*=True, *denial*=False) to limit potential noise introduced through LLM-based annotation. The final prompt yielded a precision of 0.887. This prompt was then applied to annotate the 8,197 candidate claim-title pairs. Refer to Appendix for the prompt used for LLM annotation and the distribution of *contradict* and *denial* labels. Consequently, 2,072 pairs (around 25%) were retained as valid contradicting claim-title pairs. These pairs

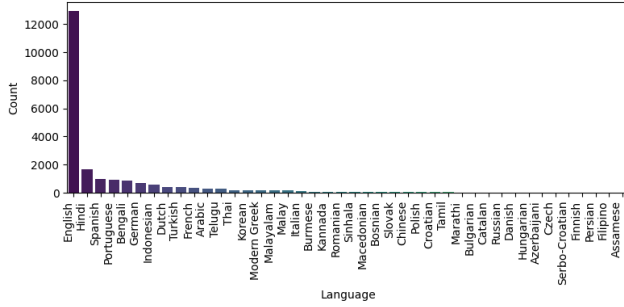


Figure 5: Language distribution of *MultiCaption*

were further expanded using the self-expansion technique and merged with manually validated pairs, resulting in a total of 3,219 contradicting pairs.

### Multilingual Claims

The original source, *MultiClaim*, is predominantly multilingual; consequently, *MultiCaption* also comprises claims written in 64 languages. Figure 5 illustrates the distribution of languages appearing at least 5 times in the dataset. Among these, English is the most frequent, followed by Hindi, Spanish, Portuguese, and Bengali. Among the 11K claim pairs, 8,131 are monolingual and 2,957 are multilingual, hence enabling multi- and cross-lingual research. Within the multilingual subset, Hindi–English combination constitutes the largest portion, with 973 pairs.

### Claim Topics

Table 2 presents the top 20 topics identified with BERTopic (Grootendorst 2022) clustering. As expected, many reflect major events from recent years. Topics with global relevance tend to appear across a larger number of languages. Other topics are more geographically specific, but strong diffusion within particular regions may explain their high volume of associated claims.

## Experiment Setup

### Datasets

In addition to our *MultiCaption* dataset that we introduce, we evaluate several baselines on the widely used out-of-context benchmark *COSMOS* (Aneja, Bregler, and Nießner 2023) to enable broader comparison. *COSMOS* contains 1,700 English image captions labeled as either out-of-context or not. We partition *MultiCaption* into training and test splits, fine-tune baseline models on the training split, and report performance on both the *MultiCaption* test split and *COSMOS*. To assess the multilingual capabilities of the baselines, we experiment with two language configurations in *MultiCaption*:

- **Monolingual** - English translations of the claims (obtained from *MultiClaim*) are used for both training and testing. Performance on the *COSMOS* dataset is evaluated using these monolingual models.

Topic Description	# Claims	# Languages
Russian invasion of Ukraine	348	28
Riots in France	177	17
COVID-19 vaccination	137	19
Accident in Kerala	133	18
Pandemic	123	21
UFO-related news	116	12
Gaza fake scenes	109	17
Military clashes in India	108	11
Prime Minister Modi	99	6
Wildfires	90	15
South Korean politics	90	11
Rescue operations in Turkey	88	23
Murder of Hindu girl	87	8
Donald Trump	85	11
Farmers’ protests	84	11
Ancient discoveries in Ayodhya	81	14
Joe Biden/Hunter Biden case	81	7
Conflict in Myanmar	77	8
Child kidnapping and trafficking	76	16
Israeli soldiers	74	14

Table 2: Top 20 topics in *MultiCaption*

- **Multilingual** - The original multilingual claims are used for both training and testing.

The following subsection details the procedure used to construct the *MultiCaption* train/test split.

**Test Partition Split** The goal of constructing a test partition was to create a challenging set of visual claim pairs for evaluation—one that is balanced across classes, reflects a diverse distribution of textual similarities, and remains fully disjoint from the training set. We initially considered the 1,422 manually validated pairs as the test set. However, this partition includes many highly similar non-contradicting pairs due to the source as discussed earlier. Further, due to the multiple annotation strategies, the same claim may appear across different samples. In the worst case, a random split could allow relational leakage: for example, if pairs  $(A, B)$  and  $(A, C)$  appear in training and  $(B, C)$  appears in testing, the test set would no longer be independent. It is therefore essential to construct two strictly disjoint sets such that no claim—and no pairwise relationship between claims—present in the test set appears in the training set. We therefore apply an iterative procedure to create a disjoint and balanced test partition as follows:

1. Initialize the test set with the 1,422 manually validated pairs.
2. Enforce disjointness: move any pair from the training set to the test set if it contains a claim already present in the test set. For example, if  $(B, C)$  appears in the test set, then  $(A, B)$  and  $(A, C)$  are moved from training to test.
3. Balance classes within similarity bins: for each similarity bin of width 0.05 between 0 and 1, randomly sample instances from training and move to the test set to equalize the number of contradictory and non-contradictory pairs whenever possible.

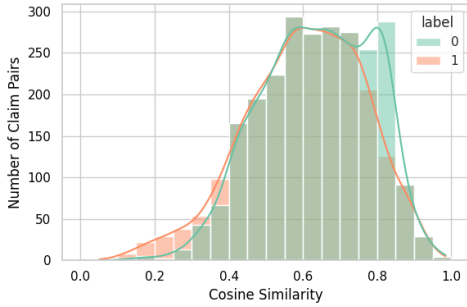


Figure 6: Textual similarity distribution in *MultiCaption* - test partition

4. Repeat step 2 & 3 until no additional samples can be moved from training to test.
5. Finalize class balance: within each similarity bin, randomly discard samples from the majority class so that both classes contain equal counts in test set. Manually annotated samples were not discarded at this step.

The iterative process produced a test partition of 4920 claim pairs spanning 52 languages, with an approximately equal number of contradicting and non-contradicting pairs and a balanced similarity distribution, as shown in 6. Minor imbalances within certain bins arise from retaining the manually validated pairs during Step 5.

**Training Data Expansion** Creating a highly balanced test partition resulted in highly imbalanced train set, containing 4747 non-contradicting pairs and only 804 contradicting pairs. Since our primary goal was to establish a robust benchmark for identifying contradictory visual claims, we prioritized the integrity of the test set.

To mitigate the imbalance of contradicting pairs in training, we adopted a training data expansion strategy using GPT-5. For the 804 contradicting pairs, we instructed the model to generate a paraphrase of the claim and a paraphrase of the title in English to avoid noise (refer to Appendix for more details and prompt used). Each contradicting pair (A,B) produced two additional pairs: (A, Paraphrase B) and (Paraphrase A, B), yielding 1,608 additional training pairs.

Our preliminary results revealed that some baselines trained in the *MultiCaption* training partition performed poorly on the *COSMOS* dataset. One contributing factor is that 56% of *COSMOS* contradicting examples (out-of-context pairs) contains one of the negation expressions we initially used to filter out claims from *MultiCaption*. To train models capable of detecting both types of contradiction—*independent claims* and *negations*—we therefore further expanded the training set with *MultiClaim* claim-title pairs in which the title is a direct negation of the claim. A similar set of negation expressions<sup>2</sup> used to filter out claim-

<sup>2</sup>We excluded certain expressions that frequently appear in titles in negation form but may not indicate true negation: “claim that”, “nothing”, “generated”, “fact check:”, “associated with”, “context”, “since”, “not”, “actually”, “as if”, “yes”, “in reports about”, “was taken in”, “attention”, “claiming”

Dataset	# C Pairs	# NC Pairs	Total Pairs	# Languages
MultiCaption - Train	4020	4767	8795	59
MultiCaption - Test	2415	2505	4920	52
COSMOS	850	850	1700	English

Table 3: Statistics of datasets used in the experiments

title pairs initially (Table 5) was used to select pairs to be included in the train set. Following the paraphrasing expansion strategy, we added 1,608 direct negation pairs. We also observed that non-contradicting pairs tend to be more similar than contradicting pairs. To decrease the imbalance, the 1,608 negation pairs were randomly selected from the high-similarity region (cosine similarity 0.5–0.85). The resulting expanded training set contains 4,020 contradicting pairs and 4,767 non-contradicting pairs across 59 languages.

Table 3 summarizes the statistics of the *MultiCaption* train and test partitions, as well as *COSMOS*.

## Baselines

We employ the following three types of models to establish stronger baselines for detecting contradictory visual claims. Refer to Appendix for the implementation of baselines.

**Transformer-based Classifiers.** We fine-tuned the following three multilingual transformer architectures for sentence classification: XLM-Roberta-large (XLM-R) (Conneau et al. 2020), Multilingual DeBERTa-V3 (mDeBERTa) (He, Gao, and Chen 2021), and Multilingual BERT (mBERT) (Devlin et al. 2019).

**Natural Language Inference (NLI) Models.** Since our task is closely related to NLI, we also evaluated multilingual NLI models. These models are transformer encoders equipped with a classification head fine-tuned on NLI datasets to predict one of three labels—entailment, neutral, or contradiction—for a given sentence pair. We fine-tuned the NLI models XLM-RoBERTa-large-xnli (Alotaibi, Nadeem, and Hamdy 2025) and mDeBERTa-v3-base-mnli-xnli (Ta et al. 2022) by replacing their three-way classification heads with a binary classifier. This approach allowed us to leverage the semantic reasoning abilities these models had already acquired while adapting them specifically to our task. We do not report the zero-shot performance for these models, as the neutral label cannot be directly mapped to contradiction or entailment. Moreover, we observed that the predicted label in the zero-shot setting was sensitive to the order of the input claims. For these reasons, we report only the results of the fine-tuned binary classifiers.

**Large Language Models.** We experiment with the following multilingual LLMs for identifying contradictory visual claims: Phi-4-mini-instruct (Phi-4) (Abouelenin et al. 2025), Mistral-7B-Instruct-v0.3 (Mistral) (Jiang et al. 2023), Llama-3.1-8B-Instruct (Llama3) (Dubey et al. 2024),

Gemma-7b (Gemma) (Team et al. 2024), and Qwen2.5-7B-Instruct (Qwen) (Yang et al. 2024). We evaluated the models in zero-shot and fine-tuned settings. In the zero-shot scenario, the pre-trained models were instructed to predict class labels without any task-specific training. In both settings, we used two different prompts for monolingual and multilingual language configurations. For the multilingual setup, original languages of claim pairs are provided in the multilingual prompt to guide LLMs.

## Metrics

Except for the zero-shot LLM setting, all other models were fine-tuned on the *MultiCaption* train set and evaluated on *MultiCaption* test set and *COSMOS*. Zero-shot LLMs and fine-tuned model with different random seeds was evaluated five times, and we report the mean precision, recall, F1-score, and accuracy across these runs.

## Results

### Baseline Performance

Table 4 reports the performance of the baseline and fine-tuned models on the *COSMOS* dataset and on the *MultiCaption* test partition. The results show that fine-tuned NLI models outperform the fine-tuned transformer baselines on *COSMOS*, suggesting that the NLI pretraining provides useful transferable knowledge for this dataset. However, on *MultiCaption*, these NLI-based models underperform compared to transformer models suggesting that *MultiCaption* poses challenges that extend beyond those captured in standard NLI tasks—even after fine-tuning. Moreover, none of the transformer-based classifiers exceed random-chance performance on *COSMOS*. On *MultiCaption*, however, the fine-tuned transformer model mDeBERTa achieves substantially higher F1-scores—both in monolingual and multilingual settings—surpassing many LLMs. This contrast highlights that these models adapt well to the data they are fine-tuned on, excelling when the task aligns well with their learned representations but struggling when it diverges.

A comparison of zero-shot LLM performance on *MultiCaption* and *COSMOS* reveals that, for some models, *MultiCaption* is the more challenging dataset. For instance, Mistral achieves an F1-score of 0.83 on *COSMOS* but drops to 0.72 on *MultiCaption*, with Llama-3 showing a similar pattern. However, once fine-tuned, all LLMs show substantial performance gains, underscoring the effectiveness of task-specific fine-tuning.

Overall, the fine-tuned LLMs achieve the strongest performance on both datasets regardless of language configuration. Mistral, in particular, performs consistently well, reaching an F1-score of 0.851 on *COSMOS* and 0.912 on *MultiCaption*. Notably, its performance remains stable even when trained and evaluated on claims in their original languages. Other fine-tuned LLMs, such as Gemma and Qwen, show substantial gains from multilingual training, with improvements of nearly 8% in F1-score. A similar trend is observed among multilingual fine-tuned transformer models, whose performance improves when trained and tested using the original, multilingual context. This underscores the value

		Baseline	Precision	Recall	F1-Score	Accuracy
<b>COSMOS</b>						
Finetuned Transformer	XLM-R		0.707	0.425	0.530	0.625
	mDeBERTa		0.712	0.421	0.528	0.625
	mBERT		0.653	0.405	0.500	0.595
Finetuned NLI	XLM-R		0.844	0.623	0.717	0.754
	mDeBERTa		<b>0.864</b>	0.578	0.693	0.744
Zero-shot LLM	Phi-4		0.502	0.988	0.665	0.503
	Mistral		0.789	0.882	0.833	0.823
	Llama3		0.673	0.805	0.733	0.706
	Gemma		0.500	<b>1.000</b>	0.667	0.500
	Qwen		0.706	0.626	0.663	0.683
Finetuned LLM	Phi-4		0.696	0.977	0.814	0.776
	Mistral		0.774	0.945	<b>0.851</b>	<b>0.834</b>
	Llama3		0.699	0.963	0.810	0.774
	Gemma		0.546	0.992	0.704	0.582
	Qwen		0.536	0.994	0.697	0.567
<b>MultiCaption — Test (Monolingual)</b>						
Finetuned Transformer	XLM-R		<b>0.894</b>	0.729	0.802	0.824
	mDeBERTa		0.879	0.767	0.819	0.833
	mBERT		0.884	0.675	0.766	0.797
Finetuned NLI	XLM-R		0.795	0.562	0.659	0.714
	mDeBERTa		0.808	0.654	0.723	0.754
Zero-shot LLM	Phi-4		0.490	0.990	0.656	0.489
	Mistral		0.584	0.954	0.725	0.644
	Llama3		0.515	0.973	0.673	0.536
	Gemma		0.491	<b>1.000</b>	0.658	0.491
	Qwen		0.558	0.782	0.651	0.589
Finetuned LLM	Phi-4		0.671	0.965	0.791	0.751
	Mistral		0.878	0.949	<b>0.912</b>	<b>0.910</b>
	Llama3		0.796	0.981	0.878	0.867
	Gemma		0.724	0.989	0.835	0.806
	Qwen		0.593	0.995	0.743	0.662
<b>MultiCaption — Test (Multilingual)</b>						
Finetuned Transformer	XLM-R		<b>0.956</b>	0.715	0.816	0.844
	mDeBERTa		0.951	0.767	0.849	0.866
	mBERT		0.952	0.747	0.837	0.857
Finetuned NLI	XLM-R		0.828	0.518	0.638	0.710
	mDeBERTa		0.837	0.619	0.712	0.754
Zero-shot LLM	Phi-4		0.489	0.990	0.655	0.487
	Mistral		0.528	0.975	0.685	0.560
	Llama3		0.508	0.984	0.669	0.524
	Gemma		0.491	<b>1.000</b>	0.658	0.491
	Qwen		0.522	0.723	0.606	0.539
Finetuned LLM	Phi-4		0.533	0.997	0.695	0.570
	Mistral		0.866	0.964	0.912	0.910
	Llama3		0.672	0.992	0.801	0.758
	Gemma		0.87	0.97	<b>0.917</b>	<b>0.913</b>
	Qwen		0.711	0.985	0.826	0.796

Table 4: Baseline performance on COSMOS & MultiCaption

of multilingual datasets such as *MultiCaption*, which enable models to better handle linguistic diversity while supporting more efficient multilingual fact-checking pipelines without costly content translation.

## Performance Analysis

We compare the accuracy of the best-performing baselines from each category—fine-tuned transformers, NLI models, and LLMs—under multilingual configuration across dif-



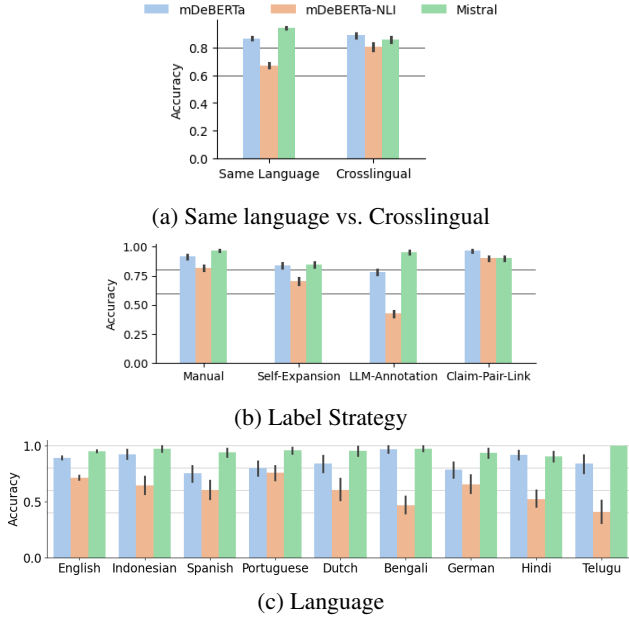


Figure 7: Accuracy of baselines across different settings: (a) same language vs. crosslingual, (b) labeling strategies, (c) language-specific accuracy.

ferent settings of the *MultiCaption* dataset. The results are shown in Figure 7. The comparison reveals that mDeBERTa-NLI achieves lower accuracy when predicting pairs from the same language (Figure 7a) as well on LLM annotated samples (Figure 7b). This is likely because most contradictory pairs fall into these two categories, and NLI models often struggle to correctly identify contradiction. In contrast, the fine-tuned Mistral model maintains more consistent performance across both settings. Figure 7c presents the accuracy of the models by language on monolingual pairs for languages with at least 150 claims in the test set. Again, Mistral is consistent, being the top performer in all languages. In contrast, mDeBERTa-NLI exhibits substantial variability, with performance dropping sharply for some languages. These results also suggest that fine-tuning a transformer model from scratch can yield better task-specific adaptation than trying to leverage previously learned knowledge from NLI models for robust multilingual performance.

## Discussion and Conclusion

This paper tackles the problem of identifying misinformation through contradictory visual claims, for which we first introduce *MultiCaption*, the first multilingual dataset for this purpose. We employ multiple annotation strategies to label claim pairs about an image or video as contradicting or non-contradicting. The dataset contains 11,088 visual claim pairs written in 64 languages, making *MultiCaption* a valuable resource for multilingual and crosslingual research.

Given the verified context of an image or video, *MultiCaption* enables the detection of misinformation by identifying claims that contradict the original context. The in-

clusion of timestamps for each claim further supports applications such as analyzing the temporal and geographical diffusion of misinformation around the same piece of media. In this work, we focus exclusively on visual claim pairs and leave multimodal extensions for future research. Nevertheless, our baseline experiments demonstrate that textual information alone is sufficient in most cases to detect contradictory visual claims. Adding vision models could potentially increase detection capabilities, but they also come with additional computational cost. Our work shows that smaller language models, like mDeBERTa, can be a strong option for environments with limited resources.

We conduct extensive experiments to establish strong baselines using transformer models, natural language inference (NLI) models, and large language models (LLMs). Overall, our results show that *MultiCaption* poses a greater challenge than standard NLI tasks: both zero-shot LLMs and fine-tuned NLI models struggle to achieve high performance. In contrast, fine-tuned transformer models and LLMs perform significantly better, underscoring the importance of task-specific training. Notably, multilingual models trained and evaluated on multilingual data achieve strong results, highlighting the dataset’s potential for building effective fact-checking pipelines without relying on machine translation. Our dataset and the fine-tuned models can easily be integrated in multilingual misinformation detection systems, allowing applications that involve claim verification and out-of-context detection. Experiments show that *MultiCaption* is more challenging than a comparable existing dataset, incorporating the multilingual dimension and a larger volume of data. Future work will focus on extending both the dataset and methodology to multimodal settings.

Beyond the empirical results, the dataset construction methodology itself constitutes a significant contribution to the field of automated fact-checking. In particular, the use of claim–post links and the self-expansion of contradicting and non-contradicting pairs provide effective mechanisms for increasing dataset scale while introducing a multilingual dimension. Widely used disinformation datasets such as MultiClaim (Pikuliak et al. 2023), MOCHEG (Yao et al. 2023) and Factify (Mishra et al. 2022; Suryavardan et al. 2023) could benefit from similar strategies, as they exploit latent links within already available data to expand coverage without requiring additional data collection, human annotation, or synthetic data generation.

## Limitations

The limitations of this work are as follows:

- A portion of the dataset is annotated using a Large Language Model (LLM). While possible annotation errors introduced by the LLM may propagate and affect the overall quality of the dataset, we tested the quality of the LLM annotations with a manually annotated dataset to mitigate this.
- Although the use of textual information is sufficient for the models to achieve high performance in contradiction detection, the approach can be extended to a multimodal setting. The current work is limited to text-only analysis.

- Detecting contradictions between two captions and determining which one constitutes misinformation requires access to the original or true context. Such context may not always be available, potentially necessitating manual intervention to identify the correct and incorrect captions.
- The construction of our dataset largely relies on the original data sources, inevitably resulting in an imbalanced distribution of samples across languages.

## Acknowledgments

Rafael Martins Frade and Rubaa Panchendrarajan are funded by the European Union and UK Research and Innovation under Grant No. 101073351 as part of Marie Skłodowska-Curie Actions (MSCA Hybrid Intelligence to monitor, promote, and analyze transformations in good democracy practices). We acknowledge Queen Mary's Apocrita HPC facility, supported by QMUL Research-IT, for enabling our experiments (King, Butcher, and Zalewski 2017).

## References

2025. <https://factcheck.afp.com/doc.afp.com.33LV9QC>. [Accessed 29-11-2025].
- Abdelnabi, S.; Hasan, R.; and Fritz, M. 2022. Open-domain, content-based, multi-modal fact-checking of out-of-context images via online resources. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 14940–14949.
- Abouelenin, A.; Ashfaq, A.; Atkinson, A.; Awadalla, H.; Bach, N.; Bao, J.; Benhaim, A.; Cai, M.; Chaudhary, V.; Chen, C.; et al. 2025. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*.
- Ahmad, R. A.; Usmanova, A.; and Rehm, G. 2025. The ClimateCheck shared task: Scientific fact-checking of social media claims about climate change. In *Proceedings of the Fifth Workshop on Scholarly Document Processing (SDP 2025)*, 263–275.
- Akhtar, M.; Schlichtkrull, M.; Guo, Z.; Cocarascu, O.; Simperl, E.; and Vlachos, A. 2023. Multimodal automated fact-checking: A survey. *arXiv preprint arXiv:2305.13507*.
- Alotaibi, A.; Nadeem, F.; and Hamdy, M. 2025. Weakly Supervised Deep Learning for Arabic Tweet Sentiment Analysis on Education Reforms: Leveraging Pre-trained Models and LLMs with Snorkel. *IEEE Access*.
- Aly, R.; Guo, Z.; Schlichtkrull, M.; Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; Cocarascu, O.; and Mittal, A. 2021. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*.
- Aneja, S.; Bregler, C.; and Nießner, M. 2023. COSMOS: catching out-of-context image misuse using self-supervised learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, 14084–14092.
- Chakraborty, T.; La Gatta, V.; Moscato, V.; and Sperli, G. 2023. Information retrieval algorithms and neural ranking models to detect previously fact-checked information. *Neurocomputing*, 557: 126680.
- Conneau, A.; Khandelwal, K.; Goyal, N.; Chaudhary, V.; Wenzek, G.; Guzmán, F.; Grave, E.; Ott, M.; Zettlemoyer, L.; and Stoyanov, V. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, 8440–8451.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 4171–4186.
- Dmonte, A.; Oruche, R.; Zampieri, M.; Calyam, P.; and Augenstein, I. 2024. Claim verification in the age of large language models: A survey. *arXiv preprint arXiv:2408.14317*.
- Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Yang, A.; Fan, A.; et al. 2024. The llama 3 herd of models. *arXiv e-prints*, arXiv-2407.
- Elsayed, T.; Nakov, P.; Barrón-Cedeno, A.; Hasanain, M.; Suwaileh, R.; Da San Martino, G.; and Atanasova, P. 2019. Overview of the CLEF-2019 CheckThat! Lab: automatic identification and verification of claims. In *International conference of the cross-language evaluation forum for European languages*, 301–321. Springer.
- FORCE11. 2020. The FAIR Data principles. <https://force11.org/info/the-fair-data-principles/>.
- Gebbru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.
- Grootendorst, M. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794*.
- Gubelmann, R.; Katis, I.; Niklaus, C.; and Handschuh, S. 2024. Capturing the varieties of natural language inference: A systematic survey of existing datasets and two novel benchmarks. *Journal of Logic, Language and Information*, 33(1): 21–48.
- He, P.; Gao, J.; and Chen, W. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Jiang, A. Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D. S.; de las Casas, D.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825*.
- King, T.; Butcher, S.; and Zalewski, L. 2017. *Apocrita - High Performance Computing Cluster for Queen Mary University of London*.
- La, T.-V.; Dao, M.-S.; Le, D.-D.; Thai, K.-P.; Nguyen, Q.-H.; and Phan-Thi, T.-K. 2022a. Leverage Boosting and Transformer on Text-Image Matching for Cheap Fakes Detection. *Algorithms*, 15(11): 423.

- La, T.-V.; Dao, M.-S.; Tran, Q.-T.; Tran, T.-P.; Tran, A.-D.; and Dang-Nguyen, D.-T. 2022b. A Combination of Visual-Semantic Reasoning and Text Entailment-based Boosting Algorithm for Cheapfake Detection. In *Proceedings of the 30th ACM International Conference on Multimedia*, 7140–7144.
- Larraz, I.; Míguez, R.; and Sallicati, F. 2023. Semantic similarity models for automated fact-checking: ClaimCheck as a claim matching tool. *Profesional de la información*, 32(3).
- Li, J.; Raheja, V.; and Kumar, D. 2024. Contradoc: understanding self-contradictions in documents with large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 6509–6523.
- Li, L.; Qin, B.; and Liu, T. 2017. Contradiction detection with contradiction-specific word embedding. *Algorithms*, 10(2): 59.
- Luo, G.; Darrell, T.; and Rohrbach, A. 2021. Newsclippings: Automatic generation of out-of-context multimodal media. *arXiv preprint arXiv:2104.05893*.
- Mishra, S.; Suryavardan, S.; Bhaskar, A.; Chopra, P.; Reganti, A. N.; Patwa, P.; Das, A.; Chakraborty, T.; Sheth, A. P.; Ekbal, A.; et al. 2022. FACTIFY: A Multi-Modal Fact Verification Dataset. In *DE-FACTIFY@ AAAI*.
- Nguyen, T.-S.; and Tran, M.-T. 2023. Multi-Models from Computer Vision to Natural Language Processing for Cheapfakes Detection. In *2023 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, 93–98. IEEE.
- Nguyen, V.-Q.; Suganuma, M.; and Okatani, T. 2022. Grit: Faster and better image captioning transformer using dual visual features. In *European Conference on Computer Vision*, 167–184. Springer.
- Nielsen, D. S.; and McConville, R. 2022. Mumin: A large-scale multilingual multimodal fact-checked misinformation social network dataset. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 3141–3153.
- Panchendrarajan, R.; Míguez, R.; and Zubiaga, A. 2025. MultiClaimNet: A Massively Multilingual Dataset of Fact-Checked Claim Clusters. *arXiv preprint arXiv:2503.22280*.
- Panchendrarajan, R.; and Zubiaga, A. 2024. Claim detection for automated fact-checking: A survey on monolingual, multilingual and cross-lingual research. *Natural Language Processing Journal*, 7: 100066.
- Papadopoulos, S.-I.; Koutlis, C.; Papadopoulos, S.; and Petrantonakis, P. C. 2024. Verite: a robust benchmark for multimodal misinformation detection accounting for unimodal bias. *International Journal of Multimedia Information Retrieval*, 13(1): 4.
- Papadopoulos, S.-I.; Koutlis, C.; Papadopoulos, S.; and Petrantonakis, P. C. 2025. Similarity over Factuality: Are we making progress on multimodal out-of-context misinformation detection? In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 5041–5050. IEEE.
- Pikuliak, M.; Srba, I.; Moro, R.; Hromadka, T.; Smoleň, T.; Melišek, M.; Vykopal, I.; Simko, J.; Podroužek, J.; and Bieliková, M. 2023. Multilingual Previously Fact-Checked Claim Retrieval. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 16477–16500.
- Pisarevskaya, D.; and Zubiaga, A. 2025. Zero-shot and Few-shot Learning with Instruction-following LLMs for Claim Matching in Automated Fact-checking. In *Proceedings of the 31st International Conference on Computational Linguistics*, 9721–9736.
- Popordanoska, T.; Li, J.; and Blaschko, M. B. 2025. CLASH: A Benchmark for Cross-Modal Contradiction Detection. *arXiv preprint arXiv:2511.19199*.
- Reimers, N.; and Gurevych, I. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Senouci, A.; Meziane, H.; and Benbernou, S. 2025. Claim polarity analysis from conflicting sources. *International Journal of Data Science and Analytics*, 19(3): 435–451.
- Sepúlveda-Torres, R.; Bonet-Jover, A.; and Saquete, E. 2023. Detecting misleading headlines through the automatic recognition of contradiction in spanish. *IEEE Access*, 11: 72007–72026.
- Singh, I.; Scarton, C.; Song, X.; and Bontcheva, K. 2023. Finding Already Debunked Narratives via Multistage Retrieval: Enabling Cross-Lingual, Cross-Dataset and Zero-Shot Learning. *arXiv e-prints*, arXiv–2308.
- Suryavardan, S.; Mishra, S.; Patwa, P.; Chakraborty, M.; Rani, A.; Reganti, A.; Chadha, A.; Das, A.; Sheth, A.; Chinnakotla, M.; et al. 2023. Factify 2: A multimodal fake news and satire news dataset. *arXiv preprint arXiv:2304.03897*.
- Ta, H. T.; Rahman, A. B. S.; Najjar, L.; and Gelbukh, A. F. 2022. Transfer Learning from Multilingual DeBERTa for Sexism Identification. In *IberLEF@ SEPLN*.
- Team, G.; Mesnard, T.; Hardin, C.; Dadashi, R.; Bhupatiraju, S.; Pathak, S.; Sifre, L.; Rivière, M.; Kale, M. S.; Love, J.; et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; and Mittal, A. 2018. FEVER: a large-scale dataset for fact extraction and VERification. *arXiv preprint arXiv:1803.05355*.
- Tran, Q.-T.; Tran, T.-P.; Dao, M.-S.; La, T.-V.; Tran, A.-D.; and Dang Nguyen, D. T. 2022. A Textual-Visual-Entailment-based Unsupervised Algorithm for Cheapfake Detection. In *Proceedings of the 30th ACM International Conference on Multimedia*, 7145–7149.
- Xu, H.; Lin, Z.; Sun, Y.; Chang, K.-W.; and Indyk, P. 2024. SparseCL: Sparse Contrastive Learning for Contradiction Retrieval. *arXiv preprint arXiv:2406.10746*.
- Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; Lin, H.; Yang, J.; Tu, J.; Zhang, J.; Yang, J.; Yang, J.; Zhou, J.; Lin, J.; Dang, K.; Lu, K.; Bao, K.; Yang, K.; Yu, L.; Li, M.; Xue, M.; Zhang, P;

Zhu, Q.; Men, R.; Lin, R.; Li, T.; Tang, T.; Xia, T.; Ren, X.; Ren, X.; Fan, Y.; Su, Y.; Zhang, Y.; Wan, Y.; Liu, Y.; Cui, Z.; Zhang, Z.; Qiu, Z.; et al. 2024. Qwen2.5 Technical Report. arXiv preprint arXiv:2412.15115.

Yao, B. M.; Shah, A.; Sun, L.; Cho, J.-H.; and Huang, L. 2023. End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2733–2743.

## Paper Checklist

### 1. For most authors...

- (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? **Yes. The research advances automated fact-checking and misinformation detection without violating privacy norms or social contracts. The dataset is constructed from publicly available fact-checking articles and social media claims, does not include personal or sensitive private data, and aims to mitigate societal harm caused by disinformation.**
- (b) Do your main claims in the abstract and introduction accurately reflect the paper’s contributions and scope? **Yes. The abstract and introduction accurately describe the creation of a multilingual dataset for contradictory visual claims, the labeling strategies employed, and the experimental evaluation of multiple model families, which align with the paper’s actual contributions.**
- (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? **Yes. The paper clearly motivates contradiction detection as a suitable formulation for identifying visual misinformation and justifies the use of data sources, labeling strategies, and multilingual modeling.**
- (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? **Yes. Refer to *MultiCaption Dataset Construction*.**
- (e) Did you describe the limitations of your work? **Yes. Refer to *Limitations*.**
- (f) Did you discuss any potential negative societal impacts of your work? **Yes. Refer to *Discussion and Conclusion* and *Limitations*.**
- (g) Did you discuss any potential misuse of your work? **Yes. Refer to *Limitations*.**
- (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? **Yes. Refer to *Baseline Implementation*. Source code will be released for reproducibility, and the dataset will be released with restricted access only for research purpose**

- (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? **Yes. The paper conforms to ICWSM ethics guidelines.**
- ### 2. Additionally, if your study involves hypotheses testing...
- (a) Did you clearly state the assumptions underlying all theoretical results? **NA. The study is empirical; no formal hypotheses are tested.**
  - (b) Have you provided justifications for all theoretical results? **NA.**
  - (c) Did you discuss competing hypotheses or theories that might challenge or complement your theoretical results? **NA.**
  - (d) Have you considered alternative mechanisms or explanations that might account for the same outcomes observed in your study? **Yes. Refer to *Results and Discussion* and *Conclusion*.**
  - (e) Did you address potential biases or limitations in your theoretical framework? **Yes. Refer to *Limitations*.**
  - (f) Have you related your theoretical results to the existing literature in social science? **Yes. Refer to *Related Work*.**
  - (g) Did you discuss the implications of your theoretical results for policy, practice, or further research in the social science domain? **Yes. Refer to *Discussion* and *Conclusion*.**
- ### 3. Additionally, if you are including theoretical proofs...
- (a) Did you state the full set of assumptions of all theoretical results? **NA.**
  - (b) Did you include complete proofs of all theoretical results? **NA.**
- ### 4. Additionally, if you ran machine learning experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **Yes. Source code and dataset are provided as supplementary materials and will be published upon acceptance.**
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **Yes. Refer to *Experiment Setup* and *Baseline Implementation* in *Appendix*.**
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **Yes. Refer to *Metrics*.**
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **Yes. Refer to *Baseline Implementation* in *Appendix*.**
  - (e) Do you justify how the proposed evaluation is sufficient and appropriate to the claims made? **Yes. Refer to *Experiment Setup* and *Results*.**
  - (f) Do you discuss what is “the cost” of misclassification and fault (in)tolerance? **Yes. Refer to *Discussion* and *Conclusion*.**

5. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets, **without compromising anonymity...**

- (a) If your work uses existing assets, did you cite the creators? **Yes.** Refer to [MultiCaption Dataset Construction](#).
- (b) Did you mention the license of the assets? **Yes.** Refer to [Artifact Availability in Appendix](#).
- (c) Did you include any new assets in the supplemental material or as a URL? **Yes.** Refer to [Artifact Availability in Appendix](#).
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **Yes.** Refer to [Artifact Availability in Appendix](#).
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or of-fensive content? **Yes.** Refer to [Artifact Availability in Appendix](#).
- (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR (see FORCE11 (2020))? **Yes.** Refer to [Artifact Availability in Appendix](#).
- (g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset (see Gebru et al. (2021))? **No.** A formal datasheet is not included; dataset documentation is provided instead. Refer to [Artifact Availability in Appendix](#).

6. Additionally, if you used crowdsourcing or conducted research with human subjects, **without compromising anonymity...**

- (a) Did you include the full text of instructions given to participants and screenshots? **NA.**
- (b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? **NA.**
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **NA.**
- (d) Did you discuss how data is stored, shared, and de-identified? **NA.**

### Similarity distribution of non-contradicting claims created via claim-post links

Figure 8 shows the similarity distribution of non-contradicting claim-pairs created via claim-post links and the discarded regions (noise and identical).

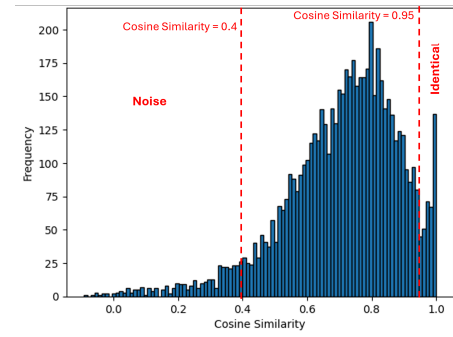


Figure 8: Similarity distribution of non-contradicting pairs generated via claim-post links

Table 5: List of negation terms used to filter out titles with direct reference or denial of the claim

Negation Terms
fake, ?, mislead, false, wrong, manipulat, edited, altered, not, no, photoshop, isn't, purport, fabricated, forged, mistaken, claim that, doesn't, doesn't, wasn't, watsn't, morphed, doctored, misuse, unrelated, rumor, nothing, generated, misrepresent, misinterpret, didn't, shared as, fact check:, fraudulent, deceive, misinform, misguide, invalid, inaccurate, untrue, erroneous, tamper, distort, modified, portrayed as, falsified, concoct, bogus, phony, hoax, associated with, aren't, isn't, neither, miscaption, context, since, fictional, montage, ,not, actually, incorrect, old, re-touched, reject, as if, doesnt, didnt, apocryphal, didnt, doesnt, yes, predate, before, previous, prior, another, none, far from, in reports about, satirical, attributed, was taken in, not, no, attention, never, claiming, staged, claim, joke, in fact, is from, are from, dates from, date from, but

### Negation Terms for Filtering Fact-Checked Article Titles

### LLM paraphrase generation of contradicting pairs

#### Prompt 1: Paraphrase generation

### Instruction:

You are given a list of image caption pairs. For each pair, generate a paraphrase of the {prefix\_param} caption only. The paraphrase should keep the same meaning, but be as different as possible.

Prompt 1 presents the instruction used to generate the paraphrases using gpt-5-mini-2025-08-07. "prefix\_param" was set to either the claim or title to generate two new pairs of contradicting samples.



## LLM Annotation of Contradicting Pairs

### Prompt 2: Contradiction and denial labeling prompt

### Instruction:

You are given two statements:

Claim { A statement describing the content of an image or video. Title { The title of a fact-checking article referring to that same image or video.

Your task is to assign two binary labels:

{"contradict": <True/False>, "denial": <True/False>}

Contradict:

Mark True if the claim and title cannot both be true for the same image or video | that is, they describe opposing or mutually exclusive facts (different events, people, locations, or circumstances). Mark False if they could both be true, or if the title simply adds compatible information. Contradict=True only when both cannot be true simultaneously, not when one is just less specific or a partial correction.

Examples:

True → \Photo shows protest in India" vs \Photo shows protest in Bangladesh."  
True → \Taken in 2021" vs \Taken in 2019."  
False → \Photo shows protest in Delhi" vs \Photo shows Indian protest in Delhi." (Compatible)

Denial:

Mark True if the title contains explicit refutation or debunking language, even if it also adds context. Typical denial cues include: - Direct negations: \does not," \did not," \is not," \no, this is not..."  
- Verdict terms: \false," \fake," \misleading," \fabricated," \debunked."  
- Authority denials: \officials denied," \government refuted," \police said it's fake."  
- Refutational phrasing: \other photos show...," \not what it claims to be."  
- Direct references: \original photo shows...," \real video is ...."

Mark False if none of these cues appear. Denial=True only if explicit refutation words are present, not merely implied by a differing fact.

Examples:

True → \Other photos from the event show them shaking hands."  
True → \Claim that this shows COVID-19 detention is false."  
False → \This video shows a 2017 airshow in Italy."

Claim:

Title:

Prompt 2 presents annotation instructions supplied to ChatGPT (snapshot *gpt-5-2025-08-07*) for assigning *contradict* and *denial* labels to candidate claim-title pairs. Table 6 presents the distribution of *contradict* and *denial* labels

### Prompt 3: Contradiction detection - Monolingual

You are an expert fact-checking assistant.

Your task is to decide if two claims about the same image or video contradict each other.

Definition of Contradict:

Answer 'Yes' if the two claims cannot both be true for the same image or video that is, they describe opposing or mutually exclusive facts (different events, people, locations, or circumstances). Answer 'No' if they could both be true, or if one simply adds compatible information. Two claims are contradicting only when both cannot be true simultaneously, not when one is just less specific or a partial correction.

Claim 1: {claim1}

Claim 2: {claim2}

Question: Do Claim 1 and Claim 2 contradict each other according to the definition above? Answer with only one word: Yes or No

Answer:

assigned by the LLM. Among the 8197 pairs, approximately 43% (*denial*=True) were discarded because the title explicitly denied the claim, and 31% were discarded as the title merely rephrased the claim (*contradict*=False, *denial*=False).

Table 6: Distribution of contradict and denial combinations

Contradict	Denial	# Pairs	Majority Cases
True	True	3227 (39%)	Title denies the fact-checked claim
True	False	2072 (25%)	Title contains only the true claim
False	True	355 (4%)	Title is a rephrase of fact-checked claim with denying phrases
False	False	2543 (31%)	Title is a rephrase of fact-checked claim

## Contradiction Detection

### LLM Prompts

Prompts 3 and 4 provide the instructions used for the LLMs to predict contradiction labels in monolingual and multilingual settings, respectively.

### Baseline Implementation

**Finetuning Transformers** Transformer models were trained with a learning rate of 2e-5, batch size 16, for 5 epochs, with 10% of training data used for validation. The model with best validation F1-Score among the 5 epochs was saved. This process was performed 5 times for each model and the average results were the ones reported. Our experiments utilized the following implementations: XLM-Roberta-large (XLM-R) <sup>3</sup>, Multilingual DeBERTa-V3 (mDe-

<sup>3</sup><https://huggingface.co/xlm-roberta-large>

---

**Prompt 4: Contradiction detection - Multilingual**

---

You are an expert multilingual fact-checking assistant. Your task is to decide if two claims about the same image or video contradict each other. The two claims below may be written in different languages, but always produce the final answer **\*\*in English only\*\***.

Definition of Contradict:

Answer 'Yes' if the two claims cannot both be true for the same image or video that is, they describe opposing or mutually exclusive facts (different events, people, locations, or circumstances). Answer 'No' if they could both be true, or if one simply adds compatible information. Two claims are contradicting only when both cannot be true simultaneously, not when one is just less specific or a partial correction.

Claim 1 written in {language1}: {claim1}

Claim 2 written in {language2}: {claim2}

Question: Do Claim 1 and Claim 2 contradict each other according to the definition above? Answer with only one word in English: Yes or No

Answer:

---

Table 7: Hyperparameters used for LLM finetuning.

Parameter	Value
<i>LoRA Configuration</i>	
LoRA rank ( $r$ )	64
LoRA alpha	16
LoRA dropout	0.1
Target modules	all-linear
<i>Training Arguments</i>	
Epochs	3
Batch size (per device)	1
Gradient accumulation steps	4
Optimizer	paged_adamw_8bit
Learning rate	2e-5
Learning rate scheduler	constant
Weight decay	0.001
Warmup ratio	0.03
Max gradient norm	0.3

BERTa)<sup>4</sup>, Multilingual BERT (mBERT)<sup>5</sup>.

**Finetuning NLI Models** To leverage the existing knowledge of the models previously trained for NLI, we set a small learning (1e-6). This allowed models to retain prior classification abilities, while still adapting to the binary contradiction/non-contradiction task. The training and evaluation process was the same for the other transformer models. We conducted our experiments using these pretrained NLI models: XLM-RoBERTa-large-xnli (XLM-RoBERTa-NLI)<sup>6</sup> and mDeBERTa-v3-base-mnli-xnli (mDeBERTa-

NLI)<sup>7</sup>

**Zero-shot LLMs** We conduct all Large language model (LLM)-based evaluations on a single GPU with 8 cores, each equipped with 11 GB of memory. Prompts 3 and 4 specify the templates used for the monolingual and multilingual evaluation settings, respectively. The LLMs are instructed to generate a single token ("Yes" or "No") with a temperature of 0.1. The following Hugging Face-hosted models were used in our experiments: microsoft/Phi-4-mini-instruct (Phi-4)<sup>8</sup>, mistralai/Mistral-7B-Instruct-v0.3 (Mistral)<sup>9</sup>, meta-llama/Llama-3.1-8B-Instruct (Llama3)<sup>10</sup>, google/gemma-7b (Gemma)<sup>11</sup>, Qwen/Qwen2.5-7B-Instruct (Qwen)<sup>12</sup>.

**Finetuning LLMs** We fine-tuned the LLMs described earlier using the same prompts applied during testing. To improve memory efficiency, we employed LoRA (Low-Rank Adaptation), a parameter-efficient fine-tuning (PEFT) method, in combination with 4-bit quantization. Table 7 summarizes the hyperparameters used during fine-tuning.

### Artifact Availability

Following research artifacts are released for reproducibility:

- The *MultiCaption* dataset<sup>13</sup>, including the finalized train and test splits used in all experiments.
- Source code<sup>14</sup> for data preprocessing, model training, and evaluation.

The source datasets *MultiClaim* and *MultiClaimNet* are available under restricted access for research purposes only, subject to their original licensing terms. The *MultiCaption* is released under the same restricted-access conditions, permitting use exclusively for non-commercial research.

No original social media content is redistributed. All claims and metadata included in *MultiCaption* are derived from publicly available fact-checking articles. The released artifacts will include documentation describing dataset structure, labeling methodology, and intended use to facilitate responsible and reproducible research.

**FAIR Principles.** The release of *MultiCaption* follows the FAIR data principles. The dataset will be *Findable* through a persistent repository link provided upon acceptance, *Accessible* under clearly defined research-only access conditions, *Interoperable* via standard, machine-readable formats, and *Reusable* through accompanying documentation describing dataset structure, labeling strategies, preprocessing steps, and intended use.

<sup>7</sup><https://huggingface.co/MoritzLaurer/mDeBERTa-v3-base-mnli-xnli>

<sup>8</sup><https://huggingface.co/microsoft/Phi-4-mini-instruct>

<sup>9</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

<sup>10</sup><https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

<sup>11</sup><https://huggingface.co/google/gemma-7b>

<sup>12</sup><https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

<sup>13</sup>Dataset is available at <https://doi.org/10.5281/zenodo.18230659>

<sup>14</sup>Source code is available at <https://github.com/rfrade/multicapTION>

<sup>4</sup><https://huggingface.co/microsoft/mdeberta-v3-large>

<sup>5</sup><https://huggingface.co/google-bert/bert-base-multilingual-cased>

<sup>6</sup><https://huggingface.co/joeddav/xlm-roberta-large-xnli>