

# Not all Fake News is Written: A Dataset and Analysis of Misleading Video Headlines

**Yoo Yeon Sung**  
University of Maryland  
yysung53@umd.edu

**Jordan Boyd-Graber**  
University of Maryland  
jbg@umiacs.umd.edu

**Naeemul Hassan**  
University of Maryland  
nhassan@umd.edu

## Abstract

Polarization and the marketplace for impressions have conspired to make navigating information online difficult for users, and while there has been a significant effort to detect false or misleading text, multimodal datasets have received considerably less attention. To complement existing resources, we present multimodal Video Misleading Headline (VMH), a dataset that consists of videos and whether annotators believe the headline is representative of the video’s contents. After collecting and annotating this dataset, we analyze multimodal baselines for detecting misleading headlines. Our annotation process also focuses on *why* annotators view a video as misleading, allowing us to better understand the interplay of annotators’ background and the content of the videos.

## 1 Introduction

Social media platforms are used by half of US adults for everyday news consumption (Walker and Matsa, 2021). They have supplanted television as the most common purveyor of news (Wakefield, 2016). However, content created on these online platforms are often lower quality than traditional sources and more prone to false stories. Vosoughi et al. (2018) contend that false news spreads six times faster online than offline.

This work focuses on one part of this problem: does a video headline match its content. We call this **misleading video headline** detection. In text, this is called incongruent headline detection (Chesney et al., 2017) and is an important problem because the headline is the first step to a reader accessing content (dos Rieis et al., 2015). While there has been work to automatically detect misleading headlines from text (Section 6), users are more likely to believe fake news when it is accompanied by videos (Wang et al., 2021)—and there are no datasets to train models for misleading video headline detection.

| VMH Dataset              |   |
|--------------------------|---|
| <b>Headline</b>          | Clinton Says Trump “Making Up Lies” About <b>New FBI Review</b>   |
| <b>Video</b>             | <a href="https://www.facebook.com/watch/?v=10154955844338812">https://www.facebook.com/watch/?v=10154955844338812</a> |
| <b>Label</b>             | <b>Misleading</b>   |
| <b>Rationale</b>         | The headline <b>implies more than what is introduced in the video.</b>  |
| <b>Subrationale</b>      | The headline <b>exaggerates</b> the video content.  |
| -----                    |   |
| <b>Annotator ID</b>      | A2P8V5SKYLL5I4  |
| <b>Annotator Profile</b> | Ages 30-49, Black, Democratic, Men, Post college  |
| <b>Venue</b>             | ABC News  |
| <b>Venue Kind</b>        | Broadcast   |
| <b>Venue Credibility</b> | High  |
| <b>News Topic</b>        | Politics  |
| <b>Headline Property</b> | Factual Statement   |
| <b>Transcript</b>        | ...is already making up lies about this he is doing his best to confuse misleading and discourage the American people |

Table 1: VMH includes video headline, video, annotator’s label, and rationales the label is grounded. In the video, the part about “New FBI Review” was not present, and thereby annotation is *misleading* because the headline was implying more than the video content.

Hence, it is necessary to investigate content outside the text (e.g., with videos) as it can help make a more informed decision by directly analyzing the relationship between the headline and the video.

To understand this new task, we create a new dataset<sup>1</sup>—Video Misleading Headline (VMH)—that includes 2,247 news articles labeled as *representative* or *misleading* (Section 2). A careful annotation process captures not just whether videos are misleading but *why*, with specific rationales. We further investigate videos, label rationales, and headline meta information (e.g., venues, news topics, and headline properties) to analyze the features that may contribute towards an instance being

<sup>1</sup><https://github.com/yysung/VMH/tree/master>

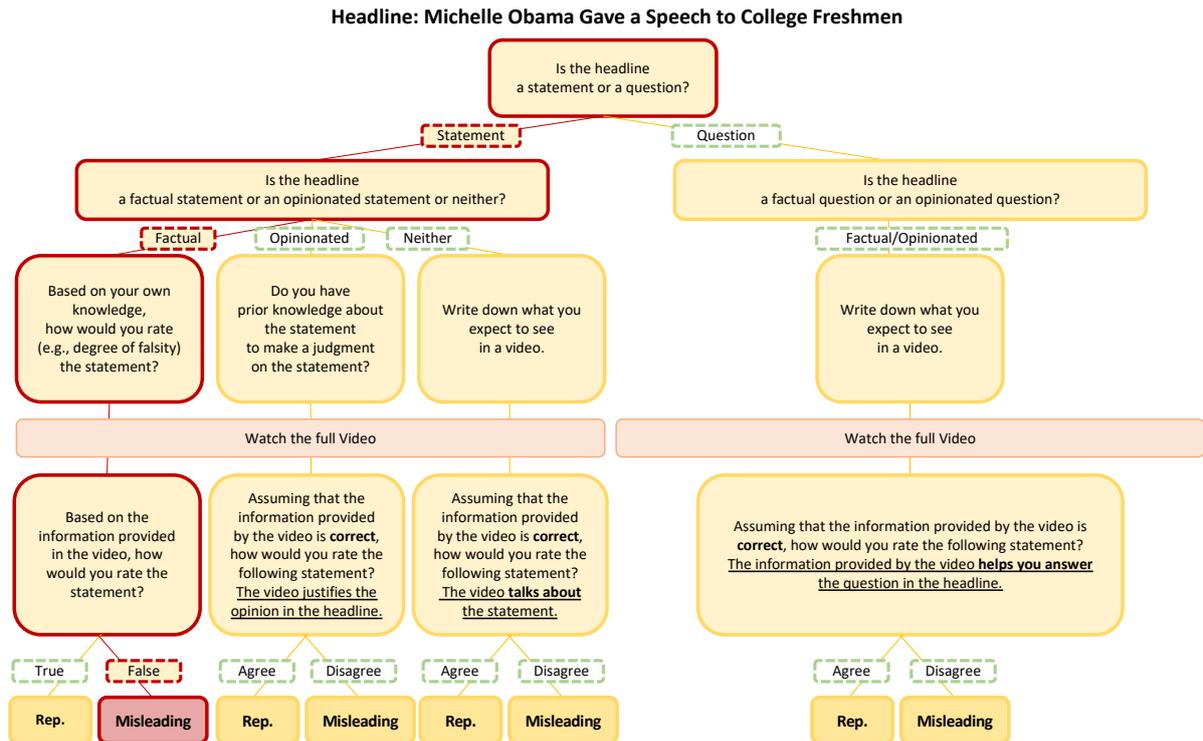


Figure 1: In the annotation tree, the annotators first consider if the headline “Michelle Obama Gave a Speech to College Freshmen” is a factual statement. Next, they answer the question, “Based on the information provided in the video, how would you rate the statement?” Because the answer was *False*, the implied label is *misleading*. The headline is indeed *misleading* because whether “College Freshman” were present in the video is unclear, making it impossible to assess the veracity. Rep. refers to *representative* label.

identified as misleading (Section 3). Section 4 shows that existing models fail to identify misleading video headlines, showing that this important but difficult task requires further research in both the text and visual domains.

## 2 Video Misleading Heading Dataset VMH

A *misleading headline* is when the headline distorts the underlying content (Wei and Wan, 2017) and facts in the news body, leading the audience to infer more or less than what was actually presented in the content. For example, in our task, the headline “Obama: I’m proud to be leaving *without scandal*” exaggerates the view of the content; the video plays Obama’s speech that he left the administration without a *significant* scandal. This distortion makes detecting misleading video headlines even more arduous because the audience has to watch the video to know if the headline is representative or—as in this case—has a subtle exaggeration or misrepresentation.

VMH consists of 2,247 video posts from 2014 to 2016. We focus on this period because it coincided

with the 2016 US presidential election, which was rife with disinformation, and is distant enough from current events that we believe annotators can be more confident about determining whether claims are true; as even news organizations are not immune to false news (Starbird et al., 2019).

Our Facebook video posts come from Rony et al. (2017), where we manually filtered any video that exceeded five minutes or had low-quality video or sound. The videos in VMH (Table 1) average two minutes long and come from fifty-two media venues, including the most circulated print and broadcast media and unreliable media in the US (Edelson et al., 2021; Samory et al., 2020, listed in Appendix A from a trustworthy journalism perspective).

We further collect venue-related information such as venue credibility<sup>2</sup> (e.g., High) and venue kind<sup>3</sup> (e.g., Broadcast). Also, we manually assigned news topics (e.g., Politics) inspired by News

<sup>2</sup><https://mediabiasfactcheck.com/>

<sup>3</sup><https://www.pewresearch.org/journalism/fact-sheet/newspapers/audience>

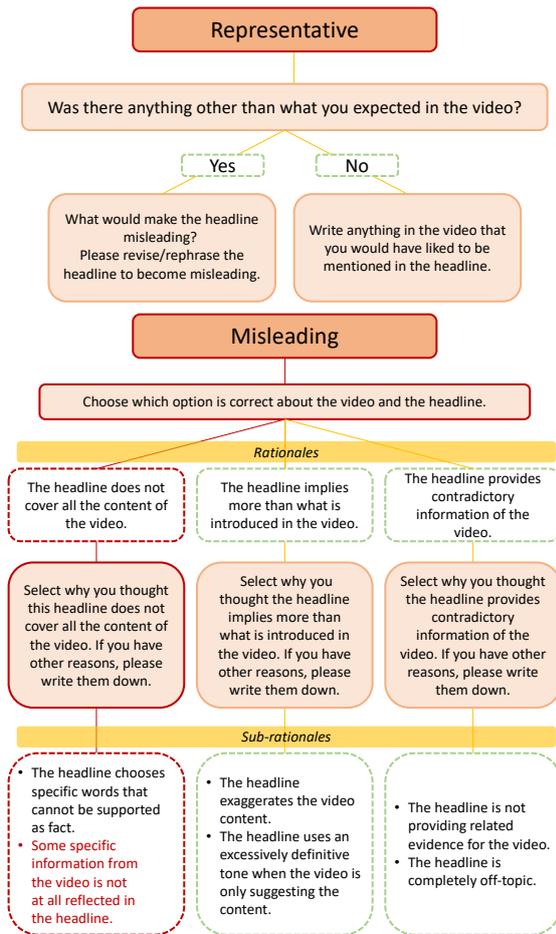


Figure 2: After label annotation, annotators provide grounding for the *misleading* labels by selecting rationales and subrationales hierarchically.

Areas<sup>4</sup> to each headline. We create audio transcripts (also released in our dataset) using automatic speech recognition software<sup>5</sup> whenever the video is accompanied by intelligible audio (Appendix I). Other features in the dataset include the number of tokens per headline (average 7.75 tokens) and annotator profile (e.g., gender).

## 2.1 Annotation

We ask Mechanical Turkers to identify misleading video headlines (Snow et al., 2008). We intentionally use non-experts to reflect the world knowledge of typical web users. For each task, the annotator goes through two phases, labeling and rationale annotation. We recruit three annotators per example (Chandler et al., 2014).

**Label Annotation** We structure the label annotation task as a series of questions to help annotators

<sup>4</sup><https://en.wikipedia.org/wiki/News>

<sup>5</sup><https://deepram.com>

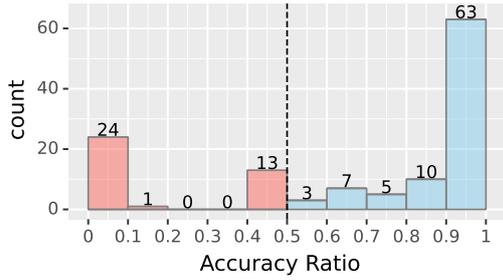
engage with the content of the headline and video (Figure 1). Because headlines can take different forms (statements of facts or opinions, questions, etc.), we first ask the user to determine the form of the headline. We refer to these forms as *headline property* in the rest of the paper. Annotators get different questions depending on the headline property: if the headline is an opinion, we ask if they agree; if the headline is a fact, we ask if they think it’s true (headline properties and associated questions in Appendix C). This helps them build a mental model of the content of the hypothetical video *before* viewing it. We adopt this format after initial pilots indicated that directly asking if a video was misleading is too ambiguous (pilot examples in Appendix B).

After the annotator has built a mental model, we ask the annotators to watch the video and answer whether the information provided in the video is consistent with the annotator’s mental model of the video. If it is, then it suggests the video is *representative*: it answered the question asked by the headline, justified an opinion, or gave evidence of a new event.

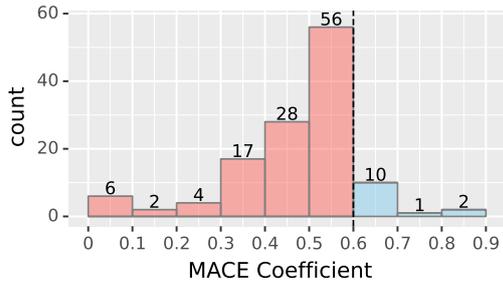
In contrast, if the video fails this check, we conclude that the headline is *misleading*. To reflect the subtle difference in participants’ opinions, we provide answer options that represent the levels of veracity or agreement with the headline (e.g., True, Mostly True, Mostly False, False, I don’t know). For the translation to binary labels, we regard the last three answers as *misleading*.

**Rationale Annotation** If their label is *misleading*, we ask the annotators to provide a *rationale*—justification—for their decision (Figure 2). For example, when an annotator labels a headline as *misleading* and chooses *The headline does not cover all the content of the video* as their rationale for the label, they then offer a subrationale to explain specifically what the headline omitted.

We offer pre-populated rationales to force objectivity in the annotator’s decision and exploit the rationales more systematically. Providing such annotations can improve not just data quality (Briakou and Carpuat, 2020)—by forcing the annotator to think about their reasoning—but also model accuracy (Zaidan et al., 2007). After the annotation is complete, final annotations are determined using a majority vote from the three annotators (Yang et al., 2015). Because subrationales can be free-form text, we do not apply majority voting for them.



(a) Qualified Workers by Accuracy Score Threshold



(b) Qualified Workers by MACE Score Threshold

Figure 3: The thresholds of accuracy ratio and MACE Coefficient are manually assigned to ensure *competent* workers are recruited after each annotation session.

## 2.2 Quality Control and Assessment

**Quality Control** We control the quality of VMH to select good crowdworkers using their accuracy score on synthetically created accuracy check questions. These questions are synthetically created to be always misleading. For each annotator, we calculate the ratio between the number of correct answers and the number of accuracy check questions they answered (examples in Appendix D).

To determine which users are reliable and to infer the labels annotators disagree on, we use a latent variable model, MACE (Paun et al., 2018), that explicitly estimates an annotator’s accuracy. This model, can correct for annotator-level biases (Martín-Morató et al., 2021, an annotator might overly favor a particular label, could have low overall accuracy, etc.). We use the point estimates—mean—from the posterior distributions of latent variables that stand for the trustworthiness of each worker (details about applying MACE to worker accuracy estimation in Appendix D).

As annotators enter the pool, we first vet them by asking for label annotations. After this “tryout” session, annotators are reinvited only if their accuracy (0.5) or MACE score (0.6) is high enough, yielding 88 and 13 qualified workers from each metric (Figure 3).

**Quality Assessment** Krippendorff’s  $\alpha$  reveals the difficulty of the task and the quality of the annotators: for the three annotators who passed the accuracy score threshold, it was 0.57 for labels and 0.33 for rationales. The Krippendorff’s  $\alpha$  values of the workers who qualified with the MACE cut-off are 0.68 (labels) and 0.21 (rationales). While the values have moderate-to-low agreement (Briakou and Carpuat, 2020), this is expected due to the inherent subjectivity of the annotation (Sandri et al., 2023; Kenyon-Dean et al., 2018; Akhtar et al., 2019; Daume III and Marcu, 2005). These inevitable disagreements are important as they can help capture the task’s nuance (Davani et al., 2022): the *source* of the disagreements can be revealing, as we discuss more in the next section.

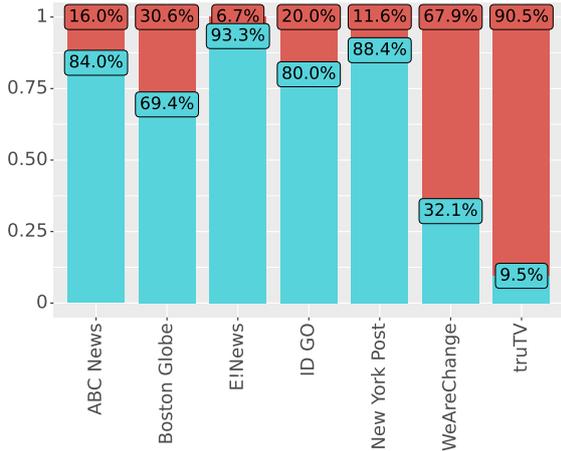
## 3 Dataset Analysis

Out of 2,247 video headlines, 1,906 headlines are annotated as *representative*, while 341 headlines are annotated as *misleading*, suggesting a high class imbalance. This section investigates VMH to understand what features contribute to (or correlate with) a headline being classified as misleading.

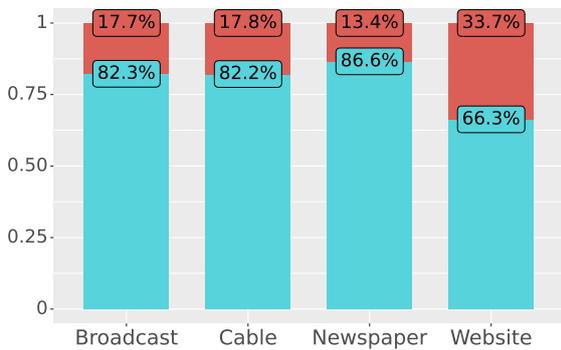
**Misleading Features** Figure 4 suggests that the venues *TruTV* and *WeAreChange.org* are strong indicators for misleading headlines. More generally, videos from the venue kind *Website* (as opposed to traditional media) are likely to be misleading (29%). The specific venue and the kind of venue may help detect misleading headlines (Appendix E).

**Clickbait** Misleading videos and clickbait both have the same goal: to entice more people to click on the underlying content. A reasonable hypothesis is that they would use similar tricks to lure in users. Thus, we reproduce the features found by (Dhoju et al., 2019) to be associated with clickbait headlines such as the number of demonstrative adjectives, numbers, and WH-words (e.g., what, who, how) for the headlines in VMH. Demonstrative adjectives do appear more frequently in misleading headlines, while numbers and superlative word features are less frequent (Table 2). Numbers and modal words appear in similar frequencies. Thus, misleading video headlines are not the same as clickbait.

**Investigation of Bias in Annotation** Because our dataset has many politically relevant videos, we also ask annotators’ political leaning to see if it biases annotations. A  $\chi^2$  test does not suggest that



(a) Venue Distribution



(b) Venue Kind Distribution

Figure 4: What proportion of headlines were misleading (red) or representative (blue) based on specific venues (top) and venue types (bottom). The venues *TruTV*, *WeAreChange.org* and venue kind *Website* were the strongest indicators of misleading headlines. The red and blue bars are proportions of *misleading* and *representative* labels. Not all venues are shown.

| Clickbait Patterns | Presence Ratio      |            |
|--------------------|---------------------|------------|
|                    | Dhoju et al. (2019) | VMH (Ours) |
| Demonstrative Adj  | 0.80                | 0.61       |
| WH-Words           | 0.70                | 0.40       |
| Numbers            | 0.72                | 0.60       |
| Modal              | 0.27                | 0.20       |
| Superlative        | 0.30                | 0.06       |

Table 2: Clickbait patterns in misleading headlines in VMH to demonstrate the difference between clickbait detection and misleading video headline task.

annotations and political leanings are dependent ( $p$ -value 0.36); the marginal proportion of misleading videos are comparable (Democratic: 22.9%, Republican: 22.6%, and Independent: 33%).

We also manually check fifty video headlines to see if their ideologies affected a headline’s assigned label, finding no substantial consequences. For example, the headline “Charles Blow: Donald

Trump is a bigot”, presumably “anti-Trump”, was annotated *Representative* by an annotator with a “Republican” leaning.

**Task Subjectivity** Motivated by Section 2.2, we examine the annotations that fail to have consensus among annotator decisions: there were 1436 *representative* and 159 *misleading* instances with the perfect agreement, leaving 30% to annotations that had disagreement. In addition to disagreeing on labels, annotators disagree about why the headline is misleading (Table 3).

## 4 Experiments

The misleading headline detection task is challenging because of the inherent subjectivity of the task. It also requires multimodal approaches that can consider both the headline and the video to make inferences about whether the headline is *representative* or not. Thus, this section benchmarks both text-only and multimodal approaches typically used for detecting video-text similarity and video-text entailment tasks.

**Experiment Settings** We compare the performance of models when trained with various combinations of input features in our dataset. The features that we consider are headlines ( $H$ ) and their associated video clips ( $V$ ), transcripts ( $T$ ), rationales, and sub-rationales ( $R$ ).

For textual features, we concatenate features as:<sup>6</sup> [SEP] {Headline [SEP] Transcript [SEP] rationale [SEP] sub-rationale}. We also extract embeddings corresponding to two multimodal models. We use VideoCLIP (Xu et al., 2021b) and VLM models (Xu et al., 2021a) that adopt zero-shot transfer learning to video-text understanding tasks.<sup>7</sup> VideoCLIP trains a transformer model using a contrastive objective on paired examples of video-text clips that maximize association between temporarily overlapping text and video segments (Xu et al., 2021b). In contrast, VLM is a task-agnostic multimodal learning model that uses novel masking schemes to improve the learning of multimodal

<sup>6</sup>While gold rationales might not be available during inference, our objective to study them as features are to highlight and understand if and how rationales can help improve detection accuracy in this task. We leave automatic prediction of the rationales to future work.

<sup>7</sup>The benchmark results in our study are to suggest baseline features and models that could be used in solving the detection task, rather than demonstrating them as a sole approach to validate the dataset or improve the detection performance.

| Headlines   | ID  | Ann. | Rationales   | Subrationales  |
|---|-----|------|--|--|
| Lester Holt Interrupted Trump Repeatedly                          | 81  | M    | The headline does not cover all the content of the video                       | The headline is not providing related evidence for the video                                   |
|   | 111 | M    | Neither of above: The headline provides contradictory information of the video | The headline chooses specific words that cannot be supported as fact                           |
|   | 97  | R    | -  | -  |
| Emily Blunt Weighs In On John Kransinskis Obsession With The D... | 42  | M    | The headline does not cover all the content of the video                       | The headline chooses specific words that cannot be supported as fact                           |
|   | 45  | M    | The headline does not cover all the content of the video                       | Some specific information from the video is not at all reflected in the headline               |
|   | 97  | R    | -  | -  |
| Did This Man Murder A Beautiful Country Music Producer            | 77  | M    | Neither of above: The headline provides contradictory information of the video | The headline is not providing related evidence for the video                                   |
|   | 12  | M    | The headline implies more than what what is introduced in the video            | The headline uses an excessively definitive tone when the video is only suggesting the content |
|   | 10  | M    | Neither of above: The headline provides contradictory information of the video | (Free Form Input) No mention of her being a country music producer                             |

Table 3: Examples of Samples with Subjectivity. The second headline shows that each annotator’s rationales are different even when the annotations are the same. The third headline shows an example where annotated subrationales all vary in their content (e.g., free-form text). ID is Annotator’s ID and Ann. is the annotation result from each annotator (M: Misleading, R: Representative)

fusion between the text and the video. We fine-tune a classification layer that takes input features extracted from video and text-based encoders as described to predict the label associated with a given video-headline pair (details in Appendix G).

**Data and Evaluation Metrics** We divide VMH into three sets: 70% for the training set, 15% for the validation set, and 15% for the test set. We evaluate using the following metrics: F1, precision, recall, AUPRC score, and accuracy. We report the precision and recall scores of the positive class, *misleading*. Each metric is estimated by averaging five replicates of stratified random splits.

## 5 Experimental Results and Model Analyses

**Experiment Results** Table 4 reports the main results: the multimodal models that use all the features, {Video Frame + Headline + Transcript + Rationale (V+H+T+R)} result in the best performance across the board, outperforming text-only based model. Adding rationales obviously helps, as these were the foundation of the annotator labels, and *subrationales* help even more (Appendix F).

Next, we validate the utility of the multimodal features in a partial-input setting. We explore how the subjectivity can affect the detection.

**Partial Input Analysis** Validating a dataset with a partial-input baseline is common in multimodal datasets (Thomason et al., 2019). Artifacts in the

dataset can lead the models to *cheat* using shortcut features that can result in poor generalizability (Feng et al., 2019). Thus, in our case, we also experiment with unimodal settings (partial input)—{Video} and {Headline}—to ensure VMH does not contain such artifacts. Using only video or text-based features result in poor F1 (0.16 – 0.18) relative to multimodal features (F1-score: > 0.22).

**Model Subjectivity Analysis** To understand the subjectivity of the task (Section 3), we also report F1-scores on the subset of the dataset, *subjective* samples (30%), that had low consensus in the annotation process. Training on this subset, even the best model with all features: {Video from VideoCLIP + Headline + Transcript + Rationale} only obtains 0.12 F1; and it drops to 0.10 with VLM compared to 0.53 (VideoCLIP) and 0.56 VLM using the entire training set. Difficult instances for humans might not include any reliable features for the model.

**Video-Text Entailment Analysis** A sceptical reader might content that this task problem is just entailment: if the headline is entailed from the video, it is representative. However, this is not a complete solution: to investigate the relationship we use transcripts to stand in for the video and then ask the RoBERTa NLI model<sup>8</sup> whether the headline is entailed from the transcript. We average the entailment score between chunked sentences from transcripts and the headlines to compensate

<sup>8</sup>fine-tuned on SNLI, MNLI, FEVER-NLI, and ANLI

| Model     | Input         | Evaluation Metrics |                    |                    |                    |                    |
|-----------|---------------|--------------------|--------------------|--------------------|--------------------|--------------------|
|           |               | F1-Score           | Precision          | Recall             | AUPRC              | Accuracy           |
| BERT      | H             | <b>0.16 (0.07)</b> | <b>0.29 (0.14)</b> | 0.11 (0.05)        | <b>0.17 (0.02)</b> | <b>0.82 (0.01)</b> |
|           | H + T         | 0.16 (0.08)        | 0.26 (0.11)        | <b>0.12 (0.06)</b> | 0.15 (0.01)        | <b>0.82 (0.01)</b> |
| VideoCLIP | H             | 0.16 (0.06)        | 0.22 (0.05)        | 0.13 (0.06)        | 0.17 (0.01)        | 0.80 (0.01)        |
|           | V             | 0.17 (0.03)        | 0.25 (0.06)        | 0.14 (0.04)        | 0.16 (0.00)        | 0.79 (0.02)        |
|           | V + H         | 0.26 (0.09)        | 0.32 (0.13)        | 0.24 (0.09)        | 0.20 (0.04)        | 0.79 (0.05)        |
|           | V + H + T     | 0.21 (0.04)        | 0.29 (0.06)        | 0.17 (0.03)        | 0.17 (0.01)        | 0.80 (0.01)        |
|           | V + H + T + R | <b>0.53 (0.06)</b> | <b>0.65 (0.08)</b> | <b>0.44 (0.06)</b> | <b>0.41 (0.05)</b> | <b>0.88 (0.01)</b> |
| VLM       | H             | 0.18 (0.05)        | 0.20 (0.06)        | 0.19 (0.09)        | 0.16 (0.01)        | 0.76 (0.04)        |
|           | V             | 0.00 (0.00)        | 0.00 (0.00)        | 0.00 (0.00)        | 0.15 (0.00)        | 0.83 (0.00)        |
|           | V + H         | 0.22 (0.06)        | 0.23 (0.05)        | 0.22 (0.06)        | 0.18 (0.02)        | 0.77 (0.02)        |
|           | V + H + T     | 0.23 (0.04)        | 0.23 (0.04)        | 0.56 (0.01)        | 0.17 (0.01)        | 0.76 (0.01)        |
|           | V + H + T + R | <b>0.56 (0.03)</b> | <b>0.63 (0.02)</b> | <b>0.52 (0.05)</b> | <b>0.40 (0.03)</b> | <b>0.88 (0.00)</b> |

Table 4: Benchmark Evaluation Results. Rows for each model shows performance with different input features: headlines (H), videos (V), transcripts (T), and rationales (R). The reported metrics are the average F1-score, average Precision score, average Recall score, average AUPRC score, and average accuracy score of 5 replicates of stratified random splits of the train, valid, and test sets. The brackets indicate standard deviation for each metric.

for different lengths. To calculate if there is correlation between entailment predictions and the labels, we conduct a  $t$ -test (Gerald, 2018). The  $p$ -value is 0.01, which indicates that the difference between the two is statistically significant: this is a signal.

However, it is not a stand-alone solution; Table 5 shows examples when entailment decisions contradict the annotator’s judgments. For example, the first headline shows a high entailment score with the transcript while annotated as misleading with the rationale of “The headline does not cover all the video content”. The second and third headlines are predicted with low entailment scores or “not entail” while being annotated as *representative* by majority annotators.

## 6 Related Work

One of the major factors of misinformation is inaccurate headlines, which pervade social media platforms (Wei and Wan, 2017). Clickbait is characterized by misleading headlines, depending on the degree of deception the audience experiences (Bourgonje et al., 2017). However, clickbait detection problems are distinguished from misleading headlines as they may exaggerate the content but are not particularly misleading (Chen et al., 2015).

As the spread of fake news appears in many forms of multimedia (Aïmeur et al., 2023), several works are on constructing datasets to enable research on multimodal misleading headline detection (Bu et al., 2023). Ha et al. (2018) introduces an image-based dataset and focuses on misrepresented headlines on Instagram. Also, Shang et al. (2019)

introduces a dataset of Youtube videos with manual annotations generated by misleading seed videos from the Youtube recommendation system. Zannettou et al. (2018) proposes a misleading-labeling mechanism with both manual and automatic. In this case, annotated videos could be biased as manual and automatic annotation may not be in consensus; they can lead to erroneous annotations of misleading headlines.

Apart from dataset research, previous works focus on detecting multimodal fake news by including multimedia features such as false videos, images, audio, and caption (Qi et al., 2023; Masciari et al., 2020; Demuyakor and Opata, 2022; McCrae et al., 2022). However, these works feature general forms of fake news (i.e., deep-fake videos), not misleading headlines.

For multimodal models built for misleading headline detection tasks, Song et al. (2016) identified the video thumbnails, Li et al. (2022) uses uploader and environment features (e.g., number of likes received, the date of most recent upload), Choi and Ko (2022) uses comments and domain knowledge, and Zannettou et al. (2018) uses video’s meta statistics (e.g., number of shares) to develop a deep variational autoencoder with semi-supervised learning. Shang et al. (2019) uses a convolutional neural network approach to find the correlation between the neural net features and the headline. You et al. (2023) uses model-selected video frames as input features to the classifier to detect dissimilarity between the video and the text.

| Headlines                  | Transcripts   | Entail | Score | Answer |
|----------------------------|---|--------|-------|--------|
| The sounds of emotions     | ... We use the principles of music to work with rhythm and melody to regain the functional use of language. Phrase is if we... Nice job. Let's all. Well You wanna skip this up? Okay. Do you wanna skip it or singing it? You wanna try to sing it? Let's jump to the chorus. Okay? So darling then. Music is what emotions sound like...  | ✓      | 0.71  | M      |
| There is a double standard | ... Is there a double standard when it comes to transparency between Trump and Clinton? Well, of course, there's a double standard. ... He's doing over a hundred foreign deals and he wants to be both the commander chief and the representative in the world for the United States... I mean, the difference between telling somebody you had pneumonia on Sunday instead of Friday is not even in the same league really... | ✗      | 0.20  | R      |
| Honor a Vet I Warfighters  | ... Having worked with veterans throughout my career, I know firsthand the importance of honoring our troops. This veterans day our series the war fighters and history are partnering with Team Rub con to create honor event. ... Honor the vets and more fighters in your life, and share a photo and a story today. Learn more history dot com honor that...  | ✓      | 0.53  | R      |

Table 5: Examples that show entailment is not enough to discover misleading headlines. The first headline shows high entailment score with the transcript while annotated as *misleading* with the rationale of “The headline does not cover all the content of the video”. The second and third headline are predicted with low entailment score or “not entail” while being annotated as “representative” by majority annotators.

## 7 Conclusion and Future Work

We present VMH, a dataset of misleading headlines from social media videos. Our annotation scheme reduces the task’s subjectivity, and we verify the reliability of the annotations. We believe incorporating the crowd workers’ distinct opinions (e.g., headline types and rationales) on misleading headlines allows crude reflection of the current social media misinformation phenomenon. Through their lenses, we anticipate a better understanding of how people perceive misinformation in misleading video headlines and for future work, use it to generalize the detection models that are soon to be deployed.

To obtain even more realistic examples for this task, we encourage applying a dynamic adversarial generation pipeline. Motivated by Eisenschlos et al. (2021), misleading headlines could be authored by humans guided to break the existing video headline detection models. For example, while they are writing a *misleading* headline, if the model falsely predicts the headline as *representative*, it would become an adversarial, *realistic* example (Ma et al., 2021). These examples can prevent the model from learning superficial patterns (Kiela et al., 2021) and further be developed to become a *robust* tool for journalists to prevent them from making “honest” mistakes when writing video headlines (Dhiman, 2023).

## 8 Limitations

Although the rationales advance the model’s knowledge in detecting misleading headlines, the limitation of this paper is that gold rationales are not realistic. Thus, the current rationale setting can be set as an upper bound for the generic model evaluation. Also, when building the model, we suggest including features that are alike with “subrationale” features in VMH, which informs *how* a headline is misleading.

Moreover, we acknowledge that the visual grounding of the headline may help the model to learn how the headline is (partially) relevant to the video’s visual content. It would be interesting to see what other multimodal models with visual grounding ability could be applied to our task; a multimodal model could be designed so that it addresses the questions of whether the headline represents the message the video conveys or identifying the gap between the video message and the headline.

## 9 Ethical Considerations

We address ethical considerations for dataset papers, given that our work proposes a new dataset VMH. We reply to the relevant questions posed in the ACL 2022 Ethics FAQ.<sup>9</sup>

<sup>9</sup><https://www.acm.org/code-of-ethics>

To collect VMH videos, we follow the community guidelines by Meta by using publicly available videos that are accessible with *public-view only* accounts. Our study was pre-monitored by an official IRB review board to protect the participants' privacy rights. Moreover, the identity characteristics of the participants were self-identified by the workers by answering the survey questions.

Before distributing the survey, we collected consent forms for the workers to agree that their answers would be used for academic purposes. All workers who make good faith annotations are paid regardless of their accuracy. The MTurkers were compensated over 10 USD an hour (a rate higher than the US national minimum wage of 7.50 USD).

Although we understand that VMH may be exploited to make misleading content in the future, we emphasize the impact of its social goods; it provides the resource to combat multimodal misinformation online today. As VMH is the first dataset that introduces video for misleading headline detection, we believe it will serve as a starting point in the research community to overcome the task.

**Acknowledgements** We thank CLIP and CJ Lab members and the anonymous reviewers for their insightful feedback. We thank the user study participants for supporting this work through annotating data. Yoo Yeon Sung, Naeemul Hassan, and Jordan Boyd-Graber are supported in part by NSF Grant "BaitBuster 2.0: Keeping Users Away From Clickbait" and DARPA Grant "SHADE" projects.

## References

- Esma Aïmeur, Sabine Amri, and Gilles Brassard. 2023. [Fake news, disinformation and misinformation in social media: a review](#). *Social Network Analysis and Mining*, 13(1):30.
- Sohail Akhtar, Valerio Basile, and Viviana Patti. 2019. A new measure of polarization in the annotation of hate speech. In *AI\* IA 2019—Advances in Artificial Intelligence: XVIIIth International Conference of the Italian Association for Artificial Intelligence, Rende, Italy, November 19–22, 2019, Proceedings 18*, pages 588–603. Springer.
- Hunt Allcott, Matthew Gentzkow, and Chuan Yu. 2019. [Trends in the diffusion of misinformation on social media](#). *Research & Politics*, 6(2):2053168019848554.
- Peter Bourgonje, Julian Moreno Schneider, and Georg Rehm. 2017. [From clickbait to fake news detection: an approach based on detecting the stance of headlines to articles](#). In *Proceedings of the 2017 EMNLP workshop: natural language processing meets journalism*, pages 84–89.
- Eleftheria Briakou and Marine Carpuat. 2020. Detecting fine-grained cross-lingual semantic divergences without supervision by learning to rank. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1563–1580.
- Yuyan Bu, Qiang Sheng, Juan Cao, Peng Qi, Danding Wang, and Jintao Li. 2023. [Online misinformation video detection: A survey](#). *arXiv e-prints*, pages arXiv–2302.
- Jesse Chandler, Pam Mueller, and Gabriele Paolacci. 2014. [Nonnaïveté among amazon mechanical turk workers: Consequences and solutions for behavioral researchers](#). *Behavior research methods*, 46:112–130.
- Yimin Chen, Niall J Conroy, and Victoria L Rubin. 2015. [Misleading online content: recognizing clickbait as "false news"](#). In *Proceedings of the 2015 ACM on workshop on multimodal deception detection*, pages 15–19.
- Sophie Chesney, Maria Liakata, Massimo Poesio, and Matthew Purver. 2017. [Incongruent headlines: Yet another way to mislead your readers](#). In *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing meets Journalism*, pages 56–61.
- Hyewon Choi and Youngjoong Ko. 2022. [Effective fake news video detection using domain knowledge and multimodal data fusion on youtube](#). *Pattern Recognition Letters*, 154:44–52.
- Hal Daume III and Daniel Marcu. 2005. [Bayesian summarization at duc and a suggestion for extrinsic evaluation](#). In *Proceedings of the Document Understanding Conference, DUC-2005, Vancouver, USA*.
- Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022. Dealing with disagreements: Looking beyond the majority vote in subjective annotations. *Transactions of the Association for Computational Linguistics*, 10:92–110.
- John Demuyakor and Edward Martey Opatá. 2022. [Fake news on social media: Predicting which media format influences fake news most on facebook](#). *Journal of Intelligent Communication*, 2(1).
- Bharat Dhiman. 2023. [Does artificial intelligence help journalists: A boon or bane?](#)
- Sameer Dhoju, Md Main Uddin Rony, Muhammad Ashad Kabir, and Naeemul Hassan. 2019. [Differences in health news from reliable and unreliable media](#). In *Companion Proceedings of The 2019 World Wide Web Conference*, pages 981–987.
- Julio Cesar Soares dos Rieis, Fabrício Benvenuto de Souza, Pedro Olmo S Vaz de Melo, Raquel Oliveira Prates, Haewoon Kwak, and Jisun

- An. 2015. [Breaking the news: First impressions matter on online news](#). In *Ninth International AAAI Conference on Web and Social Media*.
- Laura Edelson, Minh-Kha Nguyen, Ian Goldstein, Oana Goga, Damon McCoy, and Tobias Lauinger. 2021. [Understanding engagement with us \(mis\) information news sources on facebook](#). In *Proceedings of the 21st ACM Internet Measurement Conference*, pages 444–463.
- Julian Eisenschlos, Bhuwan Dhingra, Jannis Bulian, Benjamin Börschinger, and Jordan Boyd-Graber. 2021. [Fool me twice: Entailment from wikipedia gamification](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 352–365.
- Shi Feng, Eric Wallace, and Jordan Boyd-Graber. 2019. [Misleading failures of partial-input baselines](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5533–5538.
- Banda Gerald. 2018. [A brief review of independent, dependent and one sample t-test](#). *International journal of applied mathematics and theoretical physics*, 4(2):50–54.
- Yui Ha, Jeongmin Kim, Donghyeon Won, Meeyoung Cha, and Jungseock Joo. 2018. [Characterizing clickbaits on instagram](#). In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12.
- Dirk Hovy, Taylor Berg-Kirkpatrick, Ashish Vaswani, and Eduard Hovy. 2013. [Learning whom to trust with mace](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1120–1130.
- Kian Kenyon-Dean, Eisha Ahmed, Scott Fujimoto, Jeremy Georges-Filteau, Christopher Glasz, Barleen Kaur, Auguste Lalande, Shruti Bhandari, Robert Belfer, Nirmal Kanagasabai, Roman Sarrazingendron, Rohit Verma, and Derek Ruths. 2018. [Sentiment analysis: It’s complicated!](#) In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1886–1895, New Orleans, Louisiana. Association for Computational Linguistics.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. [Dynabench: Rethinking benchmarking in NLP](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4110–4124, Online. Association for Computational Linguistics.
- Xiaojun Li, Xvhao Xiao, Jia Li, Changhua Hu, Junping Yao, and Shaochen Li. 2022. [A cnn-based misleading video detection model](#). *Scientific Reports*, 12(1):6092.
- Zhiyi Ma, Kawin Ethayarajh, Tristan Thrush, Somya Jain, Ledell Wu, Robin Jia, Christopher Potts, Adina Williams, and Douwe Kiela. 2021. [Dynaboard: An evaluation-as-a-service platform for holistic next-generation benchmarking](#). *Advances in Neural Information Processing Systems*, 34:10351–10367.
- Irene Martín-Morató, Manu Harju, and Annamaria Mesaros. 2021. [Crowdsourcing strong labels for sound event detection](#). In *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 246–250. IEEE.
- Elio Masciari, Vincenzo Moscato, Antonio Picariello, and Giancarlo Sperlí. 2020. [Detecting fake news by image analysis](#). In *Proceedings of the 24th symposium on international database engineering & Applications*, pages 1–5.
- Scott McCrae, Kehan Wang, and Avidah Zakhor. 2022. [Multi-modal semantic inconsistency detection in social media news posts](#). In *MultiMedia Modeling: 28th International Conference, MMM 2022, Phu Quoc, Vietnam, June 6–10, 2022, Proceedings, Part II*, pages 331–343. Springer.
- Silviu Paun, Bob Carpenter, Jon Chamberlain, Dirk Hovy, Udo Kruschwitz, and Massimo Poesio. 2018. [Comparing bayesian models of annotation](#). *Transactions of the Association for Computational Linguistics*, 6:571–585.
- Peng Qi, Yuyan Bu, Juan Cao, Wei Ji, Ruihao Shui, Junbin Xiao, Danding Wang, and Tat-Seng Chua. 2023. [Fakesv: A multimodal benchmark with rich social context for fake news detection on short video platforms](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI.
- Md Main Uddin Rony, Naeemul Hassan, and Mohammad Yousuf. 2017. [Diving deep into clickbaits: Who use them to what extents in which topics with what effects?](#) In *Proceedings of the 2017 IEEE/ACM international conference on advances in social networks analysis and mining 2017*, pages 232–239.
- Mattia Samory, Vartan Kesiz Abnoui, and Tanushree Mitra. 2020. [Characterizing the social media news sphere through user co-sharing practices](#). In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, pages 602–613.
- Marta Sandri, Elisa Leonardelli, Sara Tonelli, and Elisabetta Jezek. 2023. [Why don’t you do it right? analysing annotators’ disagreement in subjective tasks](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2428–2441, Dubrovnik, Croatia. Association for Computational Linguistics.

- Lanyu Shang, Daniel Yue Zhang, Michael Wang, Shuyue Lai, and Dong Wang. 2019. [Towards reliable online clickbait video detection: A content-agnostic approach](#). *Knowledge-Based Systems*, 182:104851.
- Rion Snow, Brendan O’Connor, Daniel Jurafsky, and Andrew Ng. 2008. [Cheap and fast – but is it good? evaluating non-expert annotations for natural language tasks](#). In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 254–263, Honolulu, Hawaii. Association for Computational Linguistics.
- Yale Song, Miriam Redi, Jordi Vallmitjana, and Alejandro Jaimes. 2016. [To click or not to click: Automatic selection of beautiful thumbnails from videos](#). In *Proceedings of the 25th ACM international conference on information and knowledge management*, pages 659–668.
- Kate Starbird, Ahmer Arif, and Tom Wilson. 2019. [Disinformation as collaborative work: Surfacing the participatory nature of strategic information operations](#). *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–26.
- Jesse Thomason, Daniel Gordon, and Yonatan Bisk. 2019. [Shifting the baseline: Single modality performance on visual navigation & QA](#). 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)*, pages 1977–1983, Minneapolis, Minnesota. Association for Computational Linguistics.
- Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. [The spread of true and false news online](#). *science*, 359(6380):1146–1151.
- Jane Wakefield. 2016. [Social media ‘outstrips tv’ as news source for young people](#). *BBC News*.
- Mason Walker and Katerina Eva Matsa. 2021. [News consumption across social media in 2021](#). *Pew Research Center*.
- Shuting Ada Wang, Min-Seok Pang, and Paul A Pavlou. 2021. [Seeing is believing? how including a video in fake news influences users’ reporting the fake news to social media platforms](#). *How Including a Video in Fake News Influences Users’ Reporting the Fake News to Social Media Platforms (August 23, 2021)*.
- Yuxi Wang, Martin McKee, Aleksandra Torbica, and David Stuckler. 2019. [Systematic literature review on the spread of health-related misinformation on social media](#). *Social science & medicine*, 240:112552.
- Wei Wei and Xiaojun Wan. 2017. [Learning to identify ambiguous and misleading news headlines](#). In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4172–4178.
- Hu Xu, Gargi Ghosh, Po-Yao Huang, Prahal Arora, Masoumeh Aminzadeh, Christoph Feichtenhofer, Florian Metz, and Luke Zettlemoyer. 2021a. [Vlm: Task-agnostic video-language model pre-training for video understanding](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4227–4239.
- Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metz, Luke Zettlemoyer, and Christoph Feichtenhofer. 2021b. [Video-clip: Contrastive pre-training for zero-shot video-text understanding](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6787–6800.
- Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. [Wikiqa: A challenge dataset for open-domain question answering](#). In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 2013–2018.
- Jinpeng You, Yanghao Lin, Dazhen Lin, and Donglin Cao. 2023. [Video rumor classification based on multi-modal theme and keyframe fusion](#). In *Computer Supported Cooperative Work and Social Computing: 17th CCF Conference, Chinese CSCW 2022, Taiyuan, China, November 25–27, 2022, Revised Selected Papers, Part I*, pages 58–72. Springer.
- Omar Zaidan, Jason Eisner, and Christine Piatko. 2007. [Using “annotator rationales” to improve machine learning for text categorization](#). In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 260–267, Rochester, New York. Association for Computational Linguistics.
- Savvas Zannettou, Sotirios Chatzis, Kostantinos Papadamou, and Michael Sirivianos. 2018. [The good, the bad and the bait: Detecting and characterizing clickbait on youtube](#). In *2018 IEEE Security and Privacy Workshops (SPW)*, pages 63–69. IEEE.
- Melissa Zimdars. 2016. [My ‘fake news list’ went viral. but made-up stories are only part of the problem](#). *The Washington Post*.

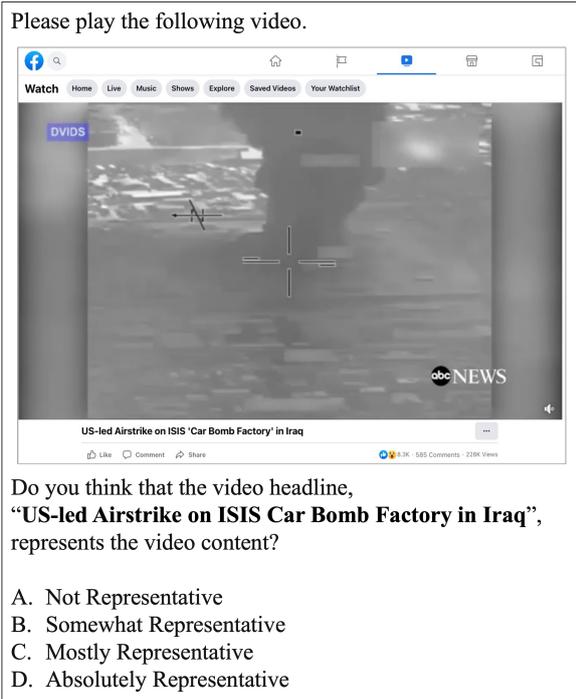


Figure 5: Example of Pilot Study. The word "represents" was too ambiguous for the audience, causing the annotators to interpret the task differently; thus it was difficult for them to consider the misleadingness of a headline.

## A Selection of Venues

We selected videos introduced by [Rony et al. \(2017\)](#) where the videos were created by mainstream media consisting of 25 most circulated print media and 43 most-watched broadcast media, and unreliable media cross-checked by two sources, information-beautiful<sup>10</sup> and [Zimdars \(2016\)](#) in the US. These were selected to include a broad range of media outlets that may include misinformation.

## B Annotation Task

**Example of Pilot Study** As demonstrated in Figure 5, our pilot study revealed that asking one question whether the video headline represented the video caused much confusion around the word *represents*, making it too ambiguous for the workers to answer the question properly. After a few interactions with workers, we found that this was due to the inherent subjectivity of the *Misleading Video Headline Detection Task*.

<sup>10</sup>Unreliable Fake News Sites

## C Questions for Headline Property

We found out from a preliminary survey that merely asking a question, *how well do you think the video headline represents the video content* causes confusion among workers due to the question’s inherent subjectivity. We assume that for different types of headlines, people follow different cognitive processes when assessing the headline’s misleadingness. Thus, we first assess the properties of the headline and ask the following questions. Examples are in Table 6 and Table 7.

**Opinionated Statement** If the worker chooses that a given headline is a *opinionated statement*, the consecutive question would be *Do you have prior knowledge about the statement in the headline to make a judgment on the statement?* to assess their original opinion stated in the headline. After watching the video, the workers are asked **Assuming that the information provided by the video is correct, how would you rate the following statement? The video justifies the opinion in the headline.** This question specifically asks to find the congruence between the video’s message and the opinion stated in the headline. If the worker finds the video content appropriate enough to match the headline, they are expected to select *Agree*. Then we conclude that the final label of the video headline is *representative*.

**Neither Statement** If the worker chooses that a given headline is a *neither statement*, the consecutive question would be *Write down what you expect to see in a video* to assess their background knowledge about the headline and what they expect to see in the video. After watching the video, the workers are asked **Assuming that the information provided by the video is correct, how would you rate the following statement? The video talks about the video.** This question specifically asks to find the congruence between the video’s message and the information in the headline. If the worker finds the video content appropriate enough to match the headline, they are expected to select *Agree*. Then we conclude that the final label of the video headline is *representative*.

**Factual/Opinionated Question** If the worker chooses that a given headline is in the form of *question*, we ask the same questions for both factual and opinionated questions. Before watching the video, the consecutive question would be *Write down what you expect to see in a video* to assess

| Factual Statement                    | Opinionated Statement  | Neither Statement          |
|--------------------------------------|--|----------------------------|
| Biden was not elected in 2020        | Best ways to make oatmeal<br>(The word 'best' is open to interpretation)                         | Great Depression           |
| Trump has 10 children                | The power of healthy food<br>(The word 'healthy' is open to interpretation)                      | Make your own coconut milk |
| She provided tips for making oatmeal | Vulgar language from Trump<br>(The word 'vulgar' is open to interpretation)                      | Tips for making oatmeal    |
| Trump to Biden: 'You're the Puppet'  | 5 minutes of truth<br>(The word 'truth' may imply different things depending on your experience) | Trump's wife               |

Table 6: Examples for Selecting Statement Headline Categories

| Factual Question                          | Opinionated Question  |
|---|---|
| Did Trump win the election?               | VP debate: Do you want a "you're hired" president?<br>(The question is asking for your personal preference)         |
| When were the first automobiles invented? | What started the French revolution?<br>(The question is asking something that is open to different interpretations) |
| Do you check the temperature every day?   | What if I made you eat worms?<br>(The question is asking for your personal preference)                              |

Table 7: Annotators are given five headline properties to choose what kind of sentence headline is.

| Original Headline  | Synthesized Headlines  | Groundings   |
|--|--|--|
| This woman takes some of the most dangerous selfies in the world | This man takes some of the most dangerous selfies in the world | False (because it is a "woman" not a man who is taking selfies in the video) |
| Baby Girl Gets Adorably Upset When Parents Kiss In Front Of Her  | Baby Boy Gets Adorably When Parents Kiss In Front Of Him       | False (because it is a "girl" not a boy who cries in the video)              |
| Trump to Clinton: 'You're the Puppet'                            | Trump to Biden: 'You're the Puppet'                            | False (because It is "Clinton" not Biden that counters Trump in the video)   |
| Toyota created a mini robot companion                            | Honda created a mini robot companion                           | False (because It is "Toyota" not Honda mentioned in the video)              |

Table 8: Examples of Synthesized Headlines for Accuracy-check Questions

their background knowledge about the headline and what they expect to see in the video. After watching the video, the workers are asked **Assuming that the information provided by the video is correct, how would you rate the following statement? The information provided by the video helps you answer the question in the headline.** This question specifically asks to find an answer to the question in the headline, assuming that video content is expected to contain the information that the headline is inquiring about. If the worker decides that the video content cannot answer or has insufficient information, they are expected to select *Disagree*. Then we conclude that the final label of

the video headline is *misleading*.

## D Quality Control and Assessment

**Pre-qualification Test** We restrict this task to the workers in the United States given that they have a higher possibility of being fluent in the verbal and literal understanding of English. Before the workers participate in the HIT, we prepare a preliminary qualification test that the workers must pass to start the HIT. All the participants must take this pre-qualification test, given multi-choice questions such as "How *representative* is the video?" and "How would you rewrite the headline." When they receive a score of 100, they are qualified to

participate in the following HITs. This process is included to ensure that the participants have the capacity to integratively comprehend the video content and video headline, and then draw out an accurate video label.

**Synthesized Headlines in Accuracy Check Questions** Table 8 shows examples of synthesized headlines in accuracy check questions. Accuracy check questions that are synthetically created to be always misleading (obviously false). For each annotator, we calculate the ratio between the number of correct answers and the number of accuracy check questions to select competent annotators.

**MACE** We compute MACE, a Bayesian approach-based metric that takes into account the credibility of the annotator and their spamming preference (Hovy et al., 2013).

$$\begin{aligned}
 &\text{for } i = 1, \dots, N : \\
 &\quad T_i \sim \text{Uniform} \\
 &\quad \text{for } j = 1, \dots, M : \\
 &\quad\quad S_{ij} \sim \text{Bernoulli}(1 - \theta_j) \\
 &\quad\quad \text{if } S_{ij} = 0 : \\
 &\quad\quad\quad A_{ij} = T_i \\
 &\quad\quad \text{else :} \\
 &\quad\quad\quad A_{ij} \sim \text{Multinomial}(\xi_j),
 \end{aligned}$$

where  $N$  denotes the number of headlines,  $T$  denotes the number of the true labels, and  $M$  denotes the number of workers.  $S_{ij}$  denotes the spam indicator of worker  $j$  for annotating headline  $i$ , while  $A_{ij}$  denotes the annotation of worker  $j$  for headline  $i$ .  $\theta$  and  $\xi$  each denotes the parameter of worker  $j$ 's trustworthiness and spam pattern. We add Beta and Dirichlet priors on  $\theta$  and  $\xi$  respectively. The assumption in the generative process is that an annotator always produces the correct label when he does not show a spam pattern which helps in excluding the labels that are not correlated with the correct label. Here, our parameter of interest is  $\theta$  which stands for the trustworthiness of each worker. We apply Paun et al. (2018)'s implementation to obtain posterior distributions (samples) of  $\theta$  and calculate point estimates.

## E Other Feature Distribution

The venue kind *Website* show higher percentage (29%) of creating misleading headlines (Table 9).

| Venue Kind | Annotated Labels |            |
|------------|------------------|------------|
|            | Representative   | Misleading |
| Broadcast  | 0.85             | 0.15       |
| Cable      | 0.85             | 0.15       |
| Newspaper  | 0.87             | 0.13       |
| Website    | 0.71             | 0.29       |

Table 9: *Website* shows more proportion of creating misleading headlines than other categories in the venue kind feature, which suggests that venue kind feature may be an indicator of representativeness of a headline.

| Headline Topics | Annotated Labels |            |
|-----------------|------------------|------------|
|                 | Representative   | Misleading |
| Entertainment   | 0.86             | 0.14       |
| Food            | 0.86             | 0.14       |
| Others          | 0.81             | 0.19       |
| Politics        | 0.85             | 0.15       |

Table 10: There was no significant difference in the proportions of topics, which suggests that topic feature is not strong indicator for misleadingness.

| Headline Properties   | Annotated Labels |            |
|-----------------------|------------------|------------|
|                       | Representative   | Misleading |
| Factual Statement     | 0.86             | 0.14       |
| Opinionated Statement | 0.84             | 0.16       |
| Neither Statement     | 0.83             | 0.17       |
| Factual Question      | 0.81             | 0.19       |
| Opinionated Question  | 0.72             | 0.28       |

Table 11: There was no significant difference in the proportions of properties, which suggests that property feature is not strong indicator for misleadingness.

On the other hand, because the proportions of misleading headlines are fairly uniform in the 1) proportions of news topics, 2) headline properties, and 3) venue credibility, it suggests that the three features are less prone to be an indicator for misleading headlines (The proportions of each label in the three features are reported in Table 10, 11 and 12).

## F What Makes for Misleadingness in Rationales?

To specifically understand how rationales help in predicting the correct *misleading* class, we trained Random Forest classifier using TF-IDF features of {Headline + Rationale + Subrationale}. Figure 6 shows the ratio of overlapping words between two types of rationales and top N important words. The top 10 words selected from the Random Forest Classifier to predict the correct label were mostly

| Venue Credibility | Annotated Labels |            |
|-------------------|------------------|------------|
|                   | Representative   | Misleading |
| High              | 0.86             | 0.14       |
| Mostly Factual    | 0.84             | 0.16       |
| Mixed             | 0.85             | 0.15       |
| Low               | 0.81             | 0.19       |
| Unknown           | 0.85             | 0.15       |

Table 12: There was no significant difference in the proportions of properties, which suggests that the headline property feature is not strong indicator for misleadingness.

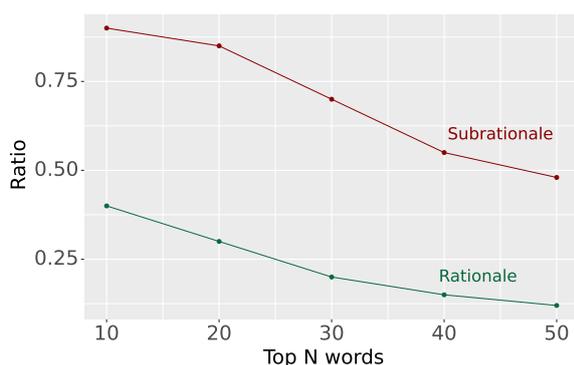


Figure 6: The top N words selected from the Random Forest Classifier to predict the correct label were mostly included in subrationales compared to rationales. As N increases, the ratio of overlapping words between the subrationale and top N important words stays higher than that of the rationale.

included in subrationales compared to rationales (Table 13).

## G Finetuning Details of Baseline Models

We finetune both VideoCLIP and VLM on a A6000 GPU using the Adam optimizer with a learning rate 0.00002, weight decay ratio of 0.001, and batch size 8 for 10 epochs. For text encoders and video encoders, we directly use the best checkpoints from the pretrained VideoCLIP and VLM models. We concatenate encoder outputs, the pooled video and text features, and learn fully connected layer optimized with Cross Entropy loss. For partial input experiments, we assign zeros to text or video encoder inputs.

## H Era of Fake News

People have been using social media platforms to converse, diffuse and broadcast their ideas in recent years. However, there has been widespread concern that misinformation is increasing on social media, which causes damage to societies (Allcott et al.,

2019). Some contemporary commentators even describe the current period as “an era of fake news” (Wang et al., 2019).

## I Censoring Audio Transcripts

We outsource transcript extractions from a software called Deepgram.<sup>11</sup> To validate its accuracy, we randomly sampled 55 videos that have transcripts and manually checked if the transcripts were accurate. These transcripts exactly matched the audio from the videos. VMH also includes transcript information on the timeframe that indicates when each word starts and ends in the video with its confidence score. We especially paid attention to this information when censoring the transcripts.

<sup>11</sup><https://deepgram.com/>

The question **MUST** be answered to proceed to the next question.

Is the headline a statement or a question?

**5-year-old Ukrainian boy rescued as he cries out for his mom**

Statement  Question

(a) Question 1

The question **MUST** be answered to proceed to the next question.

Is the headline stating a factual statement or an opinionated statement or neither? Choose from below. If you think that the headline is both factual and opinionated, select the option that best describes the statement. Click on the above Instructions button for reference.

**5-year-old Ukrainian boy rescued as he cries out for his mom**

Factual  Opinionated  Neither

(b) Question 2

The question **MUST** be answered to proceed to the next question.

Based on your own knowledge, how would you rate the statement? If you do not know, select I don't know.

Statement: **5-year-old Ukrainian boy rescued as he cries out for his mom**

False  Mostly False  Half True  Mostly True  True  I do not know

(c) Question 3



(d) Question 4

The question **MUST** be answered to proceed to the next question.

Based on the information provided in the video, how would you rate the statement? If you do not know, select I don't know. Please rate the following statement solely based on the knowledge from the video.

Statement: **5-year-old Ukrainian boy rescued as he cries out for his mom**

False  Mostly False  Half True  Mostly True  True  I do not know

(e) Question 5

The question **MUST** be answered to proceed to the next question.

Choose which option is correct about the video and the headline

**5-year-old Ukrainian boy rescued as he cries out for his mom**

The headline does not cover all the content of the video

The headline implies more than what is introduced in the video

Neither of above: The headline provides contradictory information of the video

(f) Question 6

Figure 7: Survey Example Distributed in Mturk

| Headline   | Rationale  | Subrationale   | Label      |
|--|--|--|------------|
| Tennessee Beats Georgia With Hail Mary                   | The headline does not cover all the content of the video       | Some specific information from the video is not at all reflected in the headline               | Misleading |
| President Obama Leaves For Final Overseas Trip           | The headline implies more than what is introduced in the video | The headline uses an excessively definitive tone when the video is only suggesting the content | Misleading |
| Protesters Gather Outside Chicagos Trump Tower           | The headline implies more than what is introduced in the video | Video shows a mob of people but does not provide location or reason for the protest.           | Misleading |
| Firefighters From Across US Battle Appalachian Wildfires | The headline implies more than what is introduced in the video | The headline exaggerates the video content   | Misleading |
| Tennessee Beats Georgia With Hail Mary                   | The headline does not cover all the content of the video       | The headline chooses specific words that cannot be supported as fact                           | Misleading |

Table 13: The top 10 words selected from Random Forest Classifier to predict the correct label were mostly included in subrationales compared to rationales. The word “implies” was included in the rationales, while “excessively” and “exaggerates” included in subrationales pointed the model to correctly predict *misleading*.