

# Edited Media Understanding: Reasoning About Implications of Manipulated Images

Jeff Da<sup>♠</sup> Maxwell Forbes<sup>♠♥</sup> Rowan Zellers<sup>♠♥</sup> Anthony Zheng<sup>♠</sup>

Jena Hwang<sup>♠</sup> Antoine Bosselut<sup>♠♦</sup> Yejin Choi<sup>♠♥</sup>

<sup>♠</sup>Allen Institute for Artificial Intelligence <sup>♠</sup>University of Michigan

<sup>♥</sup>Paul G. Allen School of Computer Science & Engineering, University of Washington

<sup>♦</sup>Stanford University

## Abstract

Multimodal disinformation, from ‘deepfakes’ to simple edits that deceive, is an important societal problem. Yet at the same time, the vast majority of media edits are harmless – such as a filtered vacation photo. The difference between this example, and harmful edits that spread disinformation, is one of intent. Recognizing and describing this intent is a major challenge for today’s AI systems.

We present the task of Edited Media Understanding, requiring models to answer open-ended questions that capture the intent and implications of an image edit. We introduce a dataset for our task, EMU, with 48k question-answer pairs written in rich natural language. We evaluate a wide variety of vision-and-language models for our task, and introduce a new model PELICAN, which builds upon recent progress in pretrained multimodal representations. Our model obtains promising results on our dataset, with humans rating its answers as accurate 40.35% of the time. At the same time, there is still much work to be done – humans prefer human-annotated captions 93.56% of the time – and we provide analysis that highlights areas for further progress.

## 1. Introduction

The modern ubiquity of powerful image-editing software has led to a variety of new misinformation threats. From AI-enabled “deepfakes” to low-skilled “cheapfakes,” attackers edit media to engage in a variety of harmful behaviors, such as spreading disinformation, creating revenge porn, and committing fraud ([30, 8, 23], inter alia). Accordingly, we argue that it is important to develop systems to help spot harmful manipulated media. The rapid growth and virality of social media requires as such, especially as social media trends towards visual content [16].

Correspondence to Jeff Da (jzda@cs.washington.edu).

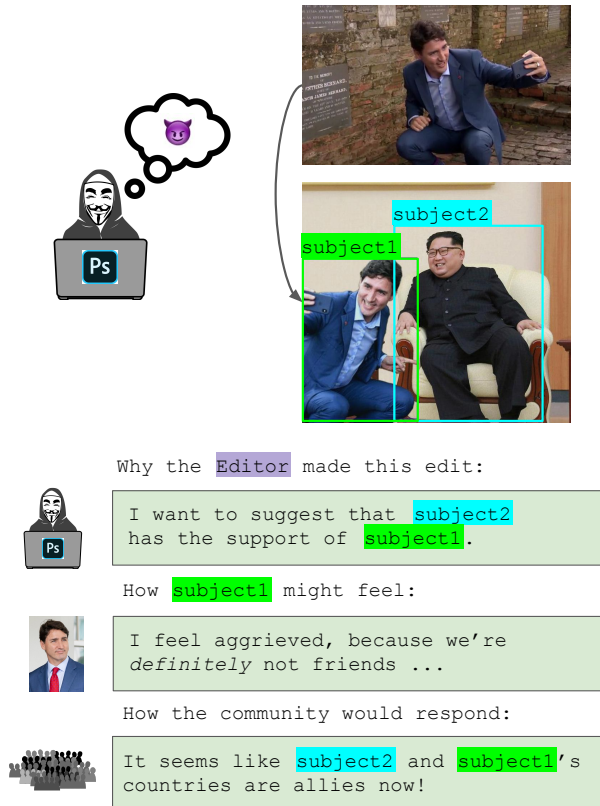


Figure 1. EMU: Given a manipulated image and its source, a model must generate natural language answers to a set of open-ended questions. Our questions test the understanding of the *what* and *why* behind important changes in the image – like that **subject1** appears to be on good terms with **subject2**.

To this end, identifying *whether* an image or video has been digitally altered (i.e., “digital forgery detection”) has been a long-standing problem in the computer vision and media forensics communities. This has enabled the development of a suite of detection approaches, such as analyzing pixel-level statistics and compression artifacts [13, 5, 4].

However, we argue that this framework is not a sufficient solution for defending against harmful manipulated media for two key reasons:

- a. **Intent.** Most manipulations are innocuous: a user might touch up a vacation photo in Adobe Photoshop, or use it to remove red-eye. These changes differ in *intent* from harmful manipulations that alter what can be reasonably inferred from the image.
- b. **Robustness.** Modern forgery detectors perform well on benchmarks by spotting lower-level patterns such as noise (e.g. [9, 17]). However, this runs the risk of overfitting to known manipulation approaches. A novel attack might fool today’s forgery detectors.

In this paper, we propose a new framework to respond to this threat: a machine must predict *why* media was edited, along with the likely implications of the edit. To do so is challenging in part because the output space, the intent and implications of the edit, is so open-ended and vast. Moreover, simply classifying an image as “harmful” is not completely helpful or explanatory. We argue that tackling this problem requires a joint approach between vision and language, with the open-ended nature of natural language being a good fit for the prediction space.

We make our framework concrete by introducing a new task, **Edited Media Understanding**, for language and vision systems (Figure 1). Given an edited image and its source, a machine must generate answers to a select set of open-ended questions regarding the edit. We introduce five such question types, which cover a diverse range of inferences necessary to fully understand the image edit. Each response involves both a classification and an explanation: for example, for question types involving disinformation, a model needs to classify the edit as misleading or not, and describe why this classification is the case.

We take two steps to require a high level of grounding through our task. First, we explicitly tie our questions and answers with the entities in the image: where applicable, a model must use tags such as `subject1`. This is important to allow models to refer to each subject individually without requiring models to identify subjects by name or use relative descriptions such as “person on the left”. Second, models must justify their predictions by providing a natural language *rationale* as part of each answer.

We then introduce a new dataset for our task, EMU, with 48k annotations over 8k image pairs. A central challenge in building this dataset at scale is finding an appropriate repository from which to source images. There is no large central resource of harmful image manipulations. Instead, we take a new approach of using images from photo-manipulation contests, known as ‘Photoshop battles,’ where different users post image edits of a (shared) source image. This has the additional benefit of letting us train and evaluate on different edits of the same image: if their semantic content differs,

then a good model should answer differently. In addition, while understanding of fake images is a difficult task, learning from EMU is an important first step, since many images are visually similar to images used in fake media.

To kickstart progress on our task, we introduce a new language and vision model, PELICAN. Our model leverages recent progress in building joint pretrained (Transformer) representations of images and text (e.g., [39, 28, 26]). A core contribution is prioritization. Since our task requires two images (and thus, twice the number of image regions), we teach models to prioritize image regions pertaining to important subjects and objects in the image – for example, the regions containing the designated main subjects.

We compare our model to a suite of strong baselines, including a standard VLP model, and show key improvement in terms of ability to reason about co-referent subjects in the edit. Nevertheless, our task is far from solved: a significant gap remains between the best machine and human accuracy.

Our contributions are thus as follows. First, we introduce a new task of Edited Media Understanding, which requires a deep understanding of *why* an image was edited, and a corresponding dataset, EMU, with 48k captions that cover diverse inferences. In addition, we introduce a new model, PELICAN, improving over recent advancements in building language-and-vision transformers. Our model makes progress, yet there is significant headroom left; we thus will release our dataset at to-be-provided upon publication.

## 2. Our Task: Edited Media Understanding

We introduce the new task of Edited Media Understanding, which focuses on holistically understanding an image edit through its context, intent, and likely implications. For example, to truly understand what changed in Figure 2, we need to understand the edit’s intent (by putting a gun in `subject1`’s hand, he is made out to be a criminal), and how the general public might react if it was distributed as a ‘real image’ (people might think of `subject1` as a totalitarian leader, threatening to kill his rivals).

Our task tests this rich image understanding through the format of open-ended question answering. A model is given the following:

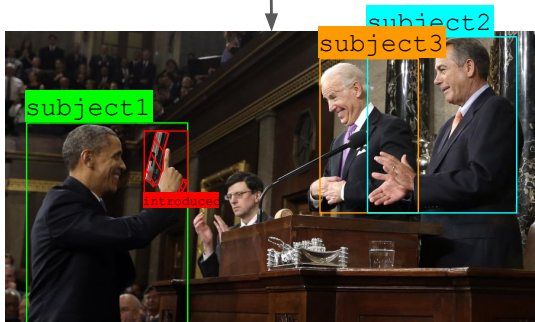
- Two images, a source image  $I_S$ , and an edited image  $I_E$ .
- A list of important entities: expressed as bounding boxes  $b_i$  for each entity (e.g. `subject3` for  $i=3$ ). These boxes ensures a high level of grounding, while avoiding the chunkiness of referring expressions (‘The man on with the red tie on the right’) or relying on explicit knowledge of the entity’s real name (‘John Boehner’).
- An open-ended question  $q$  (possibly referring to an entity); e.g., “How might `subject3` feel upon seeing this edit?”

Each question  $q$  has a binary classification label  $y$ , where the label-space is specific to the question. For example, for

source image



manual edit



fake image

subject1

**emotion:** How would subject1 feel? label: negative

subject1 would feel upset because he knows ...

**deception:** What is subject1's public image? label: negative

... subject1 is a because he's waving a gun in the air ...

subject2

**emotion:** How would subject2 feel? label: positive

subject2 would feel validated because he is part of the political opposition ...

subject3

**emotion:** How would subject3 feel? label: negative

subject3 be enraged because he knows that subject1 ...

**intent:** Why was this edit made? label: harmful

... make subject1 appear to be a man prone to violence ...

**implication:** What are potential implications? label: harmful

... subject1 seem like a criminal because of the way ...

**disinformation:** How could the public be misled? misleading

... believing that subject1, a former President, was an authoritarian leader who used force ...

Figure 2. An example from EMU. Given a source image and its edit, and a list of main subjects in the image, we collect natural language responses to applicable open-ended questions. In this particular image edit, we collect six such question-answer pairs: the first three cover the emotional reactions of subject1/2/3; the next three cover the intent and implications of the edit.

“How might subject1 feel upon seeing this edit?”, the valid options are ‘positive’ or ‘negative.’ Yet we also want to go beyond simple answering – we want models to answer *for the right reasons*, in an explainable way. Thus, given a model’s  $y$ , we require it to generate a *rationale*  $r$  explaining why its answer is true. For example, to justify why subject1 might feel “negative” in response to seeing Figure 2, a good rationale explains that the image of subject1 could be harmed because a gun was added to subject1’s hand. We evaluate by having human raters compare generated answers and rationales  $y/r$  to those written by annotators.

One last important point is how to structure the questions. In our task, we consider five open-ended question types – *intent*, *implication*, *emotion*, *deception*, and *disinformation*; descriptions of each are in Table 1. Each type focuses on a different aspect of the image edit, and is related one-to-one with an open-ended question  $q$ . Each question type may also reference a specific entity  $b$ . In these cases, the answer to the question would differ based on the main subject referred.

### 3. EMU: A Dataset of Edit Analysis

In this section, we describe how we collect data for EMU, which contains edited images, with open-ended questions and answers for each.

Our high-level goal is to construct a large dataset of semantic image manipulations: wherein an editor changes what can be reasonably inferred from an image. We argue that understanding and explaining the meaning of an image edit are two necessary subtasks for a reliable defense against deep- and cheap-faked media. We propose measuring this understanding using a question-answering format, with the open-endedness of natural language being a good fit for the open-endedness of possible image edits.

There are several challenges in building such a dataset, which we describe below. First, there is the question of *where to mine image edits*. As there is no (large) central database of harmful deepfakes that we know of, we instead use image edits from Reddit’s [r/photoshopbattles](#) community. Here,

editors tend to make complex and culturally implicative edits (e.g., reference to politics or pop culture). This makes them more similar to the deepfake detection problem than other types of edited images on the internet – such as enhanced vacation photos.

Second, there is the question of how to annotate a semantically complex image edit. We hire crowd workers on Amazon Mechanical Turk to then annotate the edits – highlighting the main subjects, and answering open-ended questions through natural language. We go into more detail in the subsections below.

### 3.1. Sourcing Image Edits

We source our image edits from the [r/photoshopbattles](#) community on Reddit. The community hosts regular Photoshop competitions, that work as follows. A competition starts with a *source photo* – then, members will comment with their own *edited photo*. One source photo may get a multitude of edited photos in response, and community members vote on which they think is the best edit. We collect image edits from this community by doing the following:

- 1. Download image edits.** To do this, we manually curated a list of more than 100 terms describing people, that also frequently appear in Photoshop battles posts (e.g. names like ‘Barack Obama’). Using our search terms, we screen over 100k posts for titles that contain one or more of the terms in question. This results in around 20k image pairs.
- 2. Filter non-people images.** To ensure that annotators do not see image pairs that do not contain any subjects, we additionally run an object detector [18] to determine if there is at least one person present in each image.

We annotated 8k of the image pairs identified through this process, some with multiple answers to the same questions, using the process described below.

### 3.2. Annotating Image Edits

In our next stage of annotation, we hire crowd workers to identify the main subjects in an image edit, and answer open-ended questions in natural language. Our annotation process is as follows:

- 1. Subject selection.** Annotators see a numbered set of people bounding boxes (produced by Mask R-CNN [18]) over the edited image. They will then select which people are main subjects, as opposed to people in the background for whom the edit is not about.
- 2. Classifying image edits.** The annotators are given a template containing all five possible question types, and are tasked with providing classification labels in regards to each question type. We gather three classifications per label and use the majority, with Cohen’s Kappa = 0.67.
- 3. Answering questions in natural language.** The annotators are tasked with filling out the template with answers

Question type		% of dataset
INTENT	Why would someone create this edit?	21.5%
IMPLICATION	What are the potential implications of this edit?	22.1%
Dis- INFORMATION	If the edit was portrayed as real news, how might it mislead the viewer?	8.9%
IMPLICATION for <b>subjectX</b>	How could this edit mislead public perception of <b>subjectX</b> ?	15.8%
EMOTION for <b>sub jectX</b>	How might this image edit make <b>sub jectX</b> feel?	31.6%

Table 1. Question types of EMU. We consider five question types, which in aggregate require a strong understanding of the image edit. The first three types are subject agnostic, though annotations refer explicitly to subjects through subject tags; the last two (with **subjectX**) are subject-specific.

for relevant questions. Some (emotion and deception, which require subjects) can be filled out several times, once for each main subject selected. Some, (deception and disinformation) do not need to be filled out for all image edits, since they may not apply. All question types can be filled out more than once (if needed for more complex edits).

- 4. Bounding edited regions.** In addition to classification labels and natural language responses, we also provide bounding boxes on each edited image denoting the edited regions in the image. We define a taxonomy {*introduced, altered, missing*} in which workers are tasked with labeling the important sections of the image that are modified.

We paid workers based on how many questions they answered per image, tracking completion time to ensure they were paid at least \$15/hr. We took several steps to ensure a high level of data quality. We used a qualification exam and checked the annotations of each worker, ensuring that they fully understood the task and that their answers were high-quality. Then, we manually reviewed annotations regularly, scoring each worker’s set of questions. We consistently gave feedback, and gave monetary benefits to high quality captions. On average, each image pair took 10-15 minutes for an annotator to label. Our answers are longer and more complex than answers for visual question answering and image captioning – in part due to the inherent open-endedness of image edits (e.g. Figure 1).

## 4. Modeling Edited Media Understanding

In this section, we present a new model for Edited Media Understanding, with a goal of kickstarting research on this challenging problem. As described in Section 2, our task differs from many standard vision-and-language tasks both in terms of format and required reasoning: a model must take as input two images (a source image and its edit), with



a significant change of implication added by the editor. A model must be able to answer questions, grounded in the main subjects of the image, describing these changes. The answers are either boolean labels, or open-ended natural language – including explainable rationales.

#### 4.1. Challenges with Multimodal Transformers

Recently, the dominant paradigm for language-and-vision tasks has shifted in favor of pretrained multi-modal Transformer models [39, 29, 26, 45, 7]. The idea is similar to how models in pure computer vision tasks often have a backbone built around Imagenet pretraining [10], and models for natural language processing are directly pretrained through language modeling [33, 11]. A Transformer model [41] is built that takes Faster-RCNN visual regions [36] as input, along with words; it is pretrained on a large dataset of paired vision-language data (like image captions on COCO [27] or Flickr30K [32]) and then finetuned for another task of interest.

These models transfer well to vision-and-language tasks with a single image, and relatively simpler closed-ended questions (such as Visual Question Answering [1]). Yet we argue that transfer to EMU is far more difficult, for a few reasons (that we confirm experimentally in Section 5):

- a. **Distinguishing subject.** Our dataset refers to important entities (like **subject1**) in an unambiguous and grounded way. Though this is trivial for humans, it differs from what models see while pretraining on image captioning data (e.g. noun phrases like “the woman” [25] and unlinked image regions).
- b. **Sparsity.** In contrast with image captioning, where the overall *gestalt* of the image is described in language, in EMU, many image regions are irrelevant. Adding an extra “source” image doubles the number of image regions, e.g. to 200. In reality, a significant number of these image regions are not needed for a human to understand the image edit, but a pretrained multimodal transformer will still attend to – and possibly be confused by – all source regions.
- c. **Open-endedness.** The questions in EMU are inherently difficult and open-ended, often requiring a significant amount of visual commonsense reasoning [44] between various image regions to answer. We hypothesize that this issue magnifies issues **a** and **b**: the challenge of the task might make it more likely for a model with suboptimal inductive biases to overfit to dataset patterns that are not representative of the true task.

For these reasons, in the next section we introduce a new model, PELICAN that is built on top of a transformer backbone, yet with added structure to handle linked image regions and the sparsity of edited image understanding.

#### 4.2. Our model: PELICAN

The challenges mentioned in Section 4.1 – linked subjects, sparsity, and inherent open-endedness – pose challenges to multimodal transformers trained on image captioning data. In other words, for Edited Media Understanding, *not all image regions are created equal*. Not only is the subject referred to in the question (e.g. **subject1**) likely important, so too are all of the regions in the image edit that are *introduced*, *altered*, or *missing*. We propose to use the annotations that collected for these regions as additional signal for the model to highlight where to attend.<sup>1</sup> Not only should a model likely attend to these important regions, it should prioritize attending to regions *nearby* (such as objects that an edited person is interacting with).

We propose to model the (likely) importance of an image region through graph propagation. We will build a directed graph with all regions of the image, rooted at a subject mentioned by the question (e.g. **subject1**). We will then *topologically sort* this graph; each region is then given an embedding corresponding to its sorted position – similar to the position embedding in a Transformer. This will allow the model to selectively attend to important image regions in the image edit. We use a different position embedding for the image source, and do not perform the graph propagation here (as we do not have *introduced/altered/missing* annotations); this separate embedding captures the inductive bias that the edited is more important than the source.

#### 4.3. Model details and Transformer integration

In this section, we describe integrating our *importance embeddings* with a multimodal transformer. In this paper, we adopt VLP [45] as our backbone, since our task is generative – for a given question expressed in natural language, we must either predict a binary token ‘yes/no’, or generate a natural-language response (consisting of several tokens).

Let the source image be  $I_S$  and  $I_E$ . We use the backbone feature extractor  $\phi$  (Faster-RCNN feature extractor [36, 3]) to extract  $N$  regions of interest for each region:

$$[\mathbf{s}_1, \dots, \mathbf{s}_N] = \phi(I_S) \quad [\mathbf{e}_1, \dots, \mathbf{e}_N] = \phi(I_E). \quad (1)$$

We note that some of these regions in  $\mathbf{e}_1, \dots, \mathbf{e}_N$  are provided to the model (as annotated regions in the image); the rest are detected by  $\phi$ . These, plus the language representation of the question, are passed to the Transformer backbone  $T$ :

$$[\mathbf{z}_1 \dots \mathbf{z}_{N+L}] = T([\mathbf{s}_1 \dots \mathbf{s}_N], [\mathbf{e}_1, \dots, \mathbf{e}_N], [\mathbf{x}_1 \dots \mathbf{x}_L]) \quad (2)$$

Important for EMU,  $\mathbf{z}_{2N+1}, \dots, \mathbf{z}_{2N+L}$  serve as language representations. Training under a left-to-right language mod-

<sup>1</sup>These annotations are collected from workers, but in theory, it would be possible to train a model to annotate regions as such. To make our task as accessible and easy-to-study as possible, however, we use the provided labels in place of a separate model however.

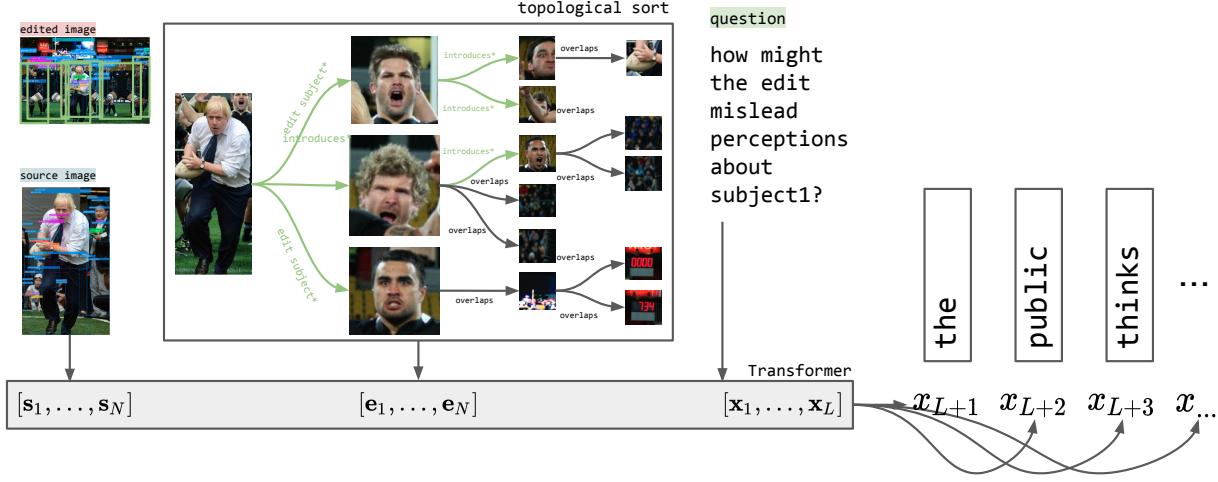


Figure 3. Overview of PELICAN. Our model takes as input all regions  $\mathbf{s}$  from the source image and  $\mathbf{e}$  from the edited image. We order the regions in  $\mathbf{e}$  using a topological sort of overlapping boxes, rooted at **subject1**. The green regions marked with an asterisk are additional regions that were introduced, and were labeled through annotators. This ordering allows the model to selectively attend to important image regions in generating an answer to the visual question about **subject1**.

eling objective, we can predict the next *next token*  $x_{L+1}$  using the representation  $\mathbf{z}_{N+L}$ .

#### 4.3.1 Prioritization Embeddings from Topological Sort

Transformers require *position embeddings* to be added to each image region and word – enabling it to distinguish which region is which. We supplement the position embeddings of the regions  $\{\mathbf{e}_1 \dots \mathbf{e}_N\}$  in the edited image  $I_E$  with the result of a topological sort.

**Graph definition.** We define the graph over image regions in the edited image as follows. We begin by sourcing a seed region  $\mathbf{s} \in \{\mathbf{e}_1 \dots \mathbf{e}_N\}$ . Let  $G = (V, E)$ , where each  $v \in V$  represents metadata of some  $\mathbf{r}_i \in \phi(I_E)$ , defined as  $v_i \in m(I_E)$  for simplicity, s.t.:

$$v_i = \{x_1, y_1, x_2, y_2, s_i, l_i\} \quad (3)$$

where  $x_1, y_1, x_2, y_2$  represents the bounding box of  $\mathbf{r}_i$ ,  $s_i \in \{1, 0\}$  denoting if  $\mathbf{r}_i$  is a subject of  $I_E$ , and  $l_i \in \{\text{introduced}, \text{altered}, \text{missing}\}$  denoting the label of  $\mathbf{r}_i$ .

We build the graph iteratively: for each iteration, we define an edge  $\mathbf{e} = \{v, u\}; u \in V$  s.t.:

$$\forall v \in m(I_E), \forall u \in V, E = E \cup (u, v) \in E' \quad (4)$$

We define  $E'$  as the set of edges  $(u, v)$  in which  $u$  and  $v$  are *notationally similar*. We define three cases in which this is true: if  $s_i \in u_i \wedge s_j \in v_j$ , if  $l_i \in u_i = l_j \in v_j$ , and if  $x_1, y_1, x_2, y_2 \in u_i$  and  $x_3, y_3, x_4, y_4 \in u_i$  overlaps, in which the percentage overlap is defined by standard intersection-over-union:

$$\frac{\min\{x_4, x_2\} - \max\{x_3, x_1\}}{\min\{y_4, y_2\} - \max\{y_3, y_1\}} \quad (5)$$

We cap the number of outgoing edges at 3, and prevent cycles by allowing edges only to unseen image regions. In cases where there are more than three possible edges, we add edges in the order defined in the previous paragraph, and break overlap ties via maximum overlap.

To produce embeddings, we run topological sort over the directed graph to assign each image region an embedding, then assign an embedding to each image region based on the ordered index. The embedding is zeroed out for image regions that are missing from the DAG, and from the source image (which are unlabeled). We include bounding box and class labels. To generate text and classification labels, we attach the embeddings onto the input for an encoder-decoder structure.

## 5. Experimental Results on EMU

In this section, we evaluate a variety of strong vision-and-language generators on EMU. We split our dataset into 80% training, with 10% for validation and testing respectively;. We perform this split on the image level (so image pairs from the training set do not leak into the test set). We use three metrics for evaluation on EMU. The first is classification – we task models with giving a label for if the response to  $q$  is positive or negative. Next, we evaluate the BLEU score of the generated labels (it is important to note, however, that endings are highly open ended, which has issues w.r.t BLEU correlation [35]). Finally, we provide two human evaluation metrics – head-to-head, in which generated responses are compared to human responses, and accuracy, in which humans are asked to label if generated responses are accurate in regards to the given edit. We also report perplexity

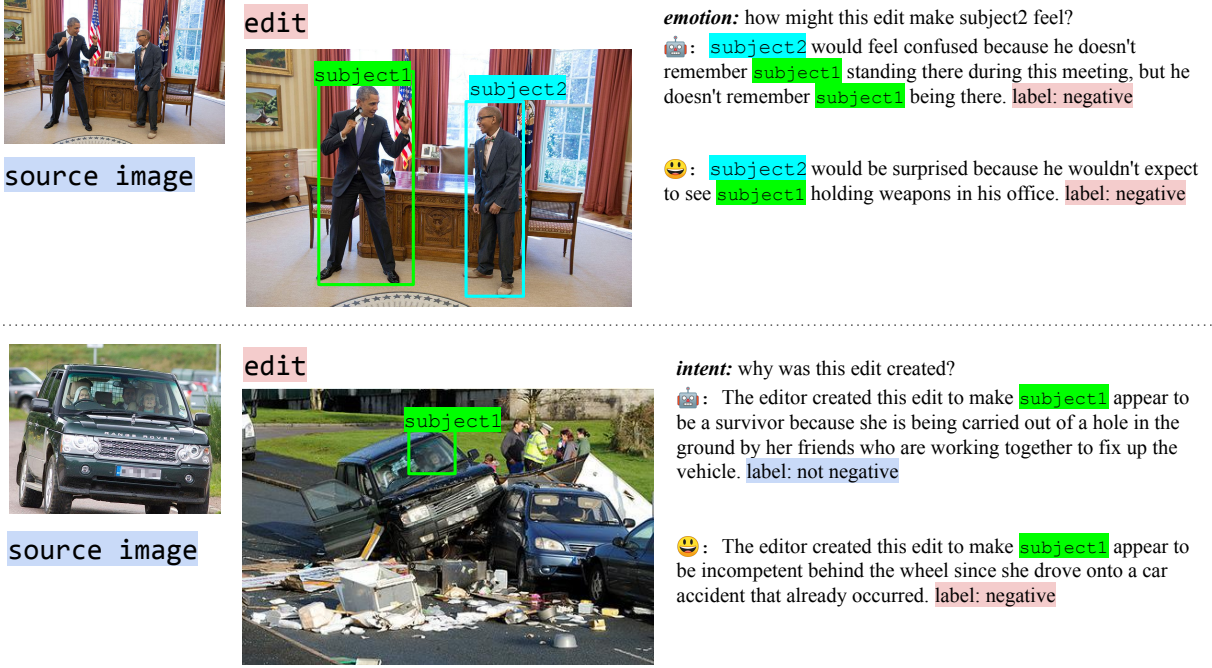


Figure 4. Generation examples from PELICAN, marked with results from human evaluation. PELICAN is able to correctly reference marked figures and is able to infer intent accordingly across each question type.

per model, which is most comparable between VLP and PELICAN due to use of the same vocab.

### 5.1. Baselines

In addition to evaluating PELICAN, we compare and evaluate the performance of various potentially high-performing baselines on our task.

**a. Retrieval.** For a retrieval baseline, which generally performs well for generation-based tasks, we use features from ResNet-158 [19], defined as  $\phi$ , to generate vectors for each  $I_E$  in the test set. We then find the most similar edited image  $I_T$  in the training set  $\mathbf{T}$  via cosine similarity:

$$\operatorname{argmax}_{I_T \in \mathbf{T}} \frac{\phi(I_E) \cdot \phi(I_T)}{\|\phi(I_E)\| \times \|\phi(I_T)\|} \quad (6)$$

We use the captions associated with the most similar image in the training set.

**b. GPT-2 [34].** As a text-only baseline, we use the 117M parameter model from GPT-2, fine-tuned on the captions from our dataset. Since the images are not taken into consideration, we generate from the seeds associated with each question type and use the same captions for all images in the test set.

**c. Cross-Modality GPT-2.** We test a unified language-and-vision model on our dataset. Similar to [2], we append the visual features  $\phi(I_S)$  and  $\phi(I_E)$  to the beginning of the token embeddings from GPT-2 (117M). For the questions

involving a subject, we append an additional vector  $\phi(r)$ , where  $r$  is the region defined by the bounding box for that subject.

**d. Dynamic Relational Attention [40].** We test the best model from previous work on image edits on our task, Dynamic Relational Attention. We train the model from scratch on our dataset, using the same procedure as [40].

**e. VLP [45].** We test VLP, a pre-trained vision-and-language transformer model. For image captioning, VLP takes a single image as input and uses an off-the-shelf object detector to extract regions, generation a caption using sequence-to-sequence decoding and treating the regions as a sequence of input tokens.

To generate a caption for a particular question type, we fix the first few generated tokens to match the prefix for that question type. We fine-tune VLP starting from weights pre-trained on Conceptual Captions (3.3m image-caption pairs) [37] and then further trained on COCO Captions (413k image-caption pairs) [27].

### 5.2. Quantitative Results and Ablation Study

We present our results in Table 2. Generations from PELICAN are preferred over human generations 14.0% of the time, with a 0.86 drop in perplexity compared to the next best model. Our model also improves in BLEU and classification accuracy. To investigate the performance of the model, we run an ablation study on various modeling

Model	Automated metrics			Human evaluation	
	Perplexity ↓	BLEU@4 ↑	Accuracy ↑	Head-to-Head ↑	Accurate % ↑
Humans	n/a	n/a	89.08	50.00	93.56
Retrieval Baseline	n/a	7.91	51.93	4.22	18.23
GPT-2 [34]	27.14	6.55	50.00	0.00	4.67
Cross-Modality GPT-2	22.71	7.12	51.01	4.36	10.35
Dynamic Relational Attention [40]	24.40	8.38	51.84	0.00	19.55
VLP [45]	12.44	9.01	53.19	10.82	27.80
PELICAN (ours)	<b>11.58</b>	<b>9.96</b>	<b>55.40</b>	<b>14.03</b>	<b>40.35</b>

Table 2. Experimental results on EMU. We compare our model, PELICAN, with several strong baseline approaches. We evaluate primarily using a human evaluation – Head-to-Head tests if model-written answers are chosen over reference answers; Accurate % measures if an open-ended answer was rated as at least ‘Slightly Accurate’. Accuracy refers to the classification accuracy on a balanced (50/50) test set, and other metrics evaluate quality of generated explanations.

Model	Automatic Eval	
	Perplexity ↓	Accuracy ↑
PELICAN	<b>11.58</b>	<b>55.40</b>
without pretraining	12.14	53.47
without annotated features	11.98	54.44
without directed graph	11.70	54.91
without source image	11.61	55.35

Table 3. Ablation study for PELICAN. Removing the indices from the topologically sorted graph reduces accuracy by 0.5%, and ignoring the annotated regions entirely reduces accuracy by another 0.5%. Excluding the source image harms accuracy only by 0.05%, suggesting that for today’s models, it is not as important as carefully modeling the edited image.

attributes, detailed in Table 3. First, we investigate the effect of pretraining (on Conceptual Captions [37, 45]). We find that performance drops without pretraining (53.47%), but surprisingly still beats other baselines. This suggests that the task requires more pragmatic inferences than the semantic learning typically gained from pre-training tasks. Second, we ablate the importance of including annotated features from the dataset when creating the directed graph (54.44%). We also ablate our use of topological sort and a directed graph by suggesting a simple (but consistent) order for image regions (54.91%). Finally, we ablate including the visual regions from the source image. The performance is similar (55.35%), suggesting that PELICAN would be able to perform in real-world settings in which only the edited image is present (e.g. social media posts).

### 5.3. Qualitative Results

Last, we present qualitative examples in Figure 4. PELICAN is able to correctly understand image pairs which require mostly surface level understanding - for example, in the top example, it is able to identify that the gun and action implies negative context, but misunderstands the response

with regards to the situation. In the bottom example, we show that PELICAN is able to refer to **subject1** correctly, but misinterprets the situation to be non-negative.

## 6. Related Work

**Language-and-Vision Datasets** Datasets involving images and languages cover a variety of tasks, including visual question answering [1, 15], image caption generation [27, 43, 24], visual story telling [31], machine translation [12], visual reasoning [22, 20, 38], and visual common sense [44].

**Two-image tasks** Though most computer vision tasks involve single images, some work has been done on exploring image pairs. The NLVR2 dataset [38] involves yes-no question answering over image pairs. Neural Naturalist [14] tests fine-grained captioning of bird pairs; [21] identifies the difference between two similar images.

**Image Edits** There has been some computer vision research studying image edits. Unlike our EMU dataset, however, much of this work has focused on modeling lower-level image edits wherein the *cultural implications* do not change significantly between images. For example, [40] predicts image editing requests (generate ‘change the background to blue’ from a pair of images). Past work has also studied learning to perform image adjustments (like colorization and enhancement) from a language query [6, 42].

## 7. Conclusion

We present Edited Media Understanding – a language-and-vision task requiring models to answer open-ended questions that capture the intent and implications of an image edit. Our dataset, EMU, is the first of its kind and is 4.8x the annotation size of the next largest image edit dataset – containing 48k question-answer pairs written in rich natural language about a variety of edited images. Our model, PELICAN, kickstarts progress on our dataset – beating all previous models and with humans rating its answers as accurate 40.35%



of the time. At the same time, there is still much work to be done – humans prefer human-annotated captions 93.56% of the time – and we provide analysis that highlights areas for further progress.

## Acknowledgements

We thank the Mechanical Turk workers for doing such an outstanding job with dataset creation - this dataset and paper would not exist without them. Thanks also to Michael Schmitz for helping with the dataset split and Jen Dumas for legal advice. This work was supported by the National Science Foundation through a Graduate Research Fellowship (DGE-1256082) and NSF grants (IIS-1524371, 1637479, 165205, 1703166), the DARPA CwC program through ARO (W911NF-15-1-0543), the IARPA DIVA program through D17PC00343, the Sloan Research Foundation through a Sloan Fellowship, the Allen Institute for Artificial Intelligence, the NVIDIA Artificial Intelligence Lab, and gifts by Google and Facebook. The views and conclusions contained herein are those of the authors and should not be interpreted as representing endorsements of IARPA, DOI/IBC, or the U.S. Government.

## References

- [1] Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh, and Dhruv Batra. Vqa: Visual question answering. *International Journal of Computer Vision*, 123:4–31, 2015. 5, 8
- [2] Chris Alberti, Jeffrey Ling, Michael Collins, and David Reitter. Fusion of detected objects in text for visual question answering. *ArXiv*, abs/1908.05054, 2019. 7
- [3] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018. 5
- [4] Jawadul H Bappy, Amit K Roy-Chowdhury, Jason Bunk, Lakshmanan Nataraj, and BS Manjunath. Exploiting spatial structure for localizing manipulated image regions. In *Proceedings of the IEEE international conference on computer vision*, pages 4970–4979, 2017. 1
- [5] Tiziano Bianchi and Alessandro Piva. Image forgery localization via block-grained analysis of jpeg artifacts. *IEEE Transactions on Information Forensics and Security*, 7(3):1003–1017, 2012. 1
- [6] Jianbo Chen, Yelong Shen, Jianfeng Gao, Jingjing Liu, and Xiaodong Liu. Language-based image editing with recurrent attentive models. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8721–8729, 2017. 8
- [7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Learning universal image-text representations. *arXiv preprint arXiv:1909.11740*, 2019. 5
- [8] Bobby Chesney and Danielle Citron. Deep fakes: A looming challenge for privacy, democracy, and national security. *Calif. L. Rev.*, 107:1753, 2019. 1
- [9] Davide Cozzolino and Luisa Verdoliva. Noiseprint: a cnn-based camera model fingerprint. *IEEE Transactions on Information Forensics and Security*, 15:144–159, 2019. 2
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5
- [11] J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 5
- [12] Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30k: Multilingual english-german image descriptions. *ArXiv*, abs/1605.00459, 2016. 8
- [13] Hany Farid. A survey of image forgery detection. *IEEE Signal Processing Magazine*, 26(2):16–25, 2009. 1
- [14] Maxwell Forbes, Christine Kaeser-Chen, Piyush Sharma, and Serge J. Belongie. Neural naturalist: Generating fine-grained image comparisons. In *EMNLP/IJCNLP*, 2019. 8
- [15] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–6334, 2017. 8
- [16] Ulrike Gretzel. The visual turn in social media marketing. 2017. 1
- [17] David Güera, Sriram Baireddy, Paolo Bestagini, Stefano Tubaro, and Edward J Delp. We need no pixels: Video manipulation detection using stream descriptors. *arXiv preprint arXiv:1906.08743*, 2019. 2
- [18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. 4
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 7
- [20] Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6693–6702, 2019. 8
- [21] Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of similar images. In *EMNLP*, 2018. 8
- [22] Johanna E. Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1988–1997, 2017. 8
- [23] Jan Kietzmann, Linda W Lee, Ian P McCarthy, and Tim C Kietzmann. Deepfakes: Trick or treat? *Business Horizons*, 63(2):135–146, 2020. 1

- [24] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123:32–73, 2016. 8
- [25] Liunian Harold Li, Mark Yatskar, Da Yin, C. Hsieh, and Kai-Wei Chang. What does bert with vision look at? In *ACL*, 2020. 5
- [26] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 2, 5
- [27] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. *ArXiv*, abs/1405.0312, 2014. 5, 7, 8
- [28] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23, 2019. 2
- [29] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 5
- [30] Britt S. Paris and Joan M. Donovan. Deepfakes and cheap fakes. Technical report, Data and Society, 2019. 1
- [31] Cesc C. Park and Gunhee Kim. Expressing an image stream with a sequence of natural sentences. In *NIPS*, 2015. 8
- [32] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, 2015. 5
- [33] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. Technical report, 2018. 5
- [34] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 7, 8, 11
- [35] Ehud Reiter. A structured review of the validity of bleu. *Computational Linguistics*, Just Accepted:1–8, 2018. 6
- [36] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 5
- [37] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 7, 8
- [38] Alane Suhr, Stephanie Zhou, Iris D. Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *ACL*, 2019. 8
- [39] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *EMNLP/IJCNLP*, 2019. 2, 5
- [40] Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. Expressing visual relationships via language. In *ACL*, 2019. 7, 8, 11
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *ArXiv*, abs/1706.03762, 2017. 5
- [42] Hai Wang, Jason D. Williams, and SingBing Kang. Learning to globally edit images with textual description. *ArXiv*, abs/1810.05786, 2018. 8
- [43] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. 8
- [44] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6713–6724, 2019. 5, 8
- [45] Luowei Zhou, Hamid Palangi, Lefei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. *ArXiv*, abs/1909.11059, 2019. 5, 7, 8, 11

Distribution across subject counts

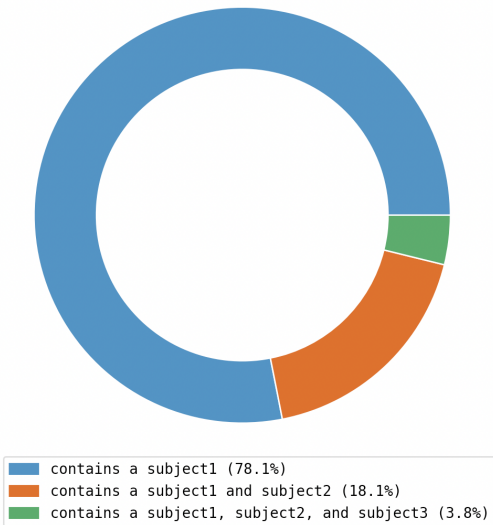


Figure 5. Subject distribution. To highlight our decision for a 3 subject limit, we show that the majority of images contains 1-2 subjects.

## A. Appendices

### A.1. Reproducibility of Experiments

We provide downloadable source code of all scripts, and experiments, at to-be-provided. We use two Titan X GPUs to train and evaluate all models, except Dynamic Relational Attention [40], which was trained on a single Titan Xp GPU. For GPT-2 [34], we use the 117M parameter model, taking 5 hours to train. Our configuration of VLP [45] has 138,208,324 parameters, taking 6 hours to train. Our model, PELICAN, has 138,208,324 parameters, taking 6 hours to train. Our Dynamic Relational Attention model has 55,165,687 parameters, taking 10 hours to train.

### A.2. Reproducibility of Hyperparameters

For models using GPT-2 as their underlying infrastructure, we use a maximum sequence length of 1024, 12 hidden layers, 12 heads for each attention layer, and 0.1 dropout in all fully connected layers. For Dynamic Relational Attention [40], we use a batch size of 95, hidden dimension size of 512, embedding dimension size of 256, 0.5 dropout, Adam optimizer, and a 1e-4 learning rate. We used early stopping based on the BLEU score on the validation set at the end of every epoch; the test scores reported are for a model trained for 63 epochs. For all models relying on VLP as their underlying infrastructure, we use 30 training epochs, 0.1 warmup proportion, 0.01 weight decay, 64 batch size.

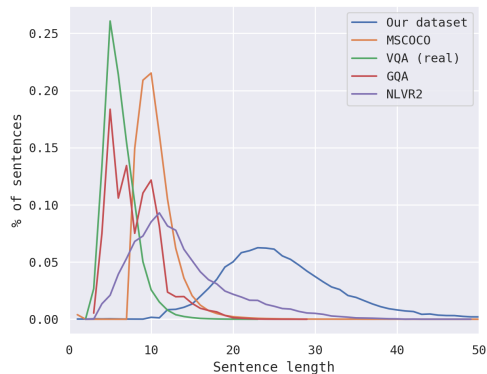


Figure 6. Language sentence length distribution, measured in words, across other language-and-vision datasets. The natural language answers in EMU show a high degree of complexity, with an average sentence length of 26.45 words.

### A.3. Reproducibility of Datasets

Our dataset has 39338 examples in the training set and 4268 and 3992 examples in the development and test sets respectively. All training on additional datasets (e.g. Zhou2019UnifiedVP) matches their implementation exactly. Our train/val/test splits were chosen at random, during the annotation period. No data was excluded, and no additional pre-processing was done. A downloadable link is available at to-be-provided after publication.

### A.4. Data Collection

For reference and reproducibility, we show the full template used to collect data in Figure 7.

We also show our human evaluation process in Figure 8.

### A.5. Additional Annotation Details

For an image pair (consisting of an image edit and a source image), we 1) ask the annotator to identify and index the main subjects in the image, 2) prime the annotator by asking them to describe the physical change in the image, 3) ask a series of questions for each main person they identified, and 4) ask a series of questions about the image as a whole. For each question we require annotators to provide both an answer to the question and a rationale (e.g. the physical change in the image edit that alludes to their answer). This is critical, as the rationales prevent models from guessing a response such as “would be harmful” without providing the proper reasoning for their response. We ask annotators to explicitly separate the rationale from the response by using the word “because” or “since” (however, we find that the vast majority of annotators naturally do this, without being explicitly prompted). For the main subjects, we limit the number of subjects to 3. This also mitigates a large variation in workload between image pairs, which was gathered as

Image edit:

Titles, Edit with Boxes, Original Image ([expand/collapse](#))

Original image title:  
\$(image\_o\_title)

Edited image title:  
\$(image\_e\_title)

Edit with Boxes:

Original Image (if this is missing, it's ok!):

Prompt 1: What physically changed in the image edit?

Focus on the important change(s) only - how you would describe this edit in a sentence to someone w

Caption The important physical change in this image is

Prompt 2: Is there a main person in this image (MP1)? ☐ Yes ☒ No

Prompt 3: Is there a second main person in this image (MP2)? ☐ Yes ☒ No

Prompt 4: Is there a third main person in this image (MP3)? ☐ Yes ☒ No

Prompt 5: Write 1 - 3 captions that answer: Why did Editor create this edit? What physical part of t

For example: Editor created this edit to **make MP1 appear to be a survivor of the end of the world b**

Caption 1 Editor created this edit to

Caption 2 Editor created this edit to

Caption 3 Editor created this edit to

Prompt 6: Write 1 - 3 captions that answer: What are the possible implications of this edit? What pl

For example: This edit could potentially be used to **make MP1 seem like a public hero because it loo**

Caption 1 This edit could potentially be used to

Caption 2 This edit could potentially be used to

Caption 3 This edit could potentially be used to

Prompt 7: Write 0 - 3 captions that answer: How might this edit mislead the viewer? What physical

For example: In regards to the edit as a whole, this edit might mislead someone into believing that **the**



**You only need to answer Prompt 7 if it applies to the given edit.** Otherwise, you can skip.

Caption 1 In regards to the edit as a whole, this edit might mislead someone into believing that

Caption 2 In regards to the edit as a whole, this edit might mislead someone into believing that

Caption 3 In regards to the edit as a whole, this edit might mislead someone into believing that

Figure 7. Example of our annotation process.

Original Image:  Edited Image, with Boxes: 

Caption A: \$(human)

Caption B: \$(machine)

Which caption gives a better analysis of the edit?

☐ Definitely A

☐ Slightly A

☐ Slightly B

☐ Definitely B

How accurate is the worse caption?

☐ Slightly Accurate

☐ Not Accurate

Figure 8. Example of our evaluation process.

potentially problematic from annotator feedback. We limit the number of captions per type to 3. We find that a worker chooses to provide more than one label for a type in only a small proportion of cases, suggesting that usually, one caption is needed to convey all the information about the image edit relating to that type .

## A.6. Lexical Analysis

**Word-Level Statistics** We analyze the lexical statistics of this dataset. We remove stop words as words such as “him”. We show that different types require different language in their response. In addition, we highlight that many of the rationales involve people, suggesting that understanding social implications is critical to solving this task.



A: subject1 would ...

B: subject1 might ...

[1] Which caption gives a better analysis of the edit?

- ☐ Definitely A
- ☐ Slightly A
- ☐ Slightly B
- ☐ Definitely B

[2] How accurate is the worse caption?

- ☐ Slightly Accurate
- ☐ Not Accurate

Figure 9. Our template for human evaluations. Each annotator is shown an edited image, the source image, and is asked to compare a human annotated captions and a machine annotated caption.



<b>Rationales</b>		<b>Responses</b>									
		intent		implication		disinformation		emotion		deception	
holding	4.21%	fun	4.83%	public	3.07%	movie	2.93%	confused	7.62%	likes	3.00%
face	4.09%	powerful	1.13%	think	2.12%	woman	2.12%	amused	4.38%	hates	2.21%
wearing	3.17%	funny	1.09%	man	1.75%	new	1.92%	embarrassed	3.88%	loves	1.36%
man	2.64%	hero	1.02%	fun	1.68%	game	1.23%	upset	3.50%	wants	1.35%
appears	2.41%	movie	1.01%	disgrace	1.25%	real	1.23%	proud	2.61%	doesn't	1.31%

Table 4. Lexical statistics. Statistics for each dimension represent omit the rationale, and statistics for the rationale are reported separately.