

The authors bear all responsibility in case of violation rights. Upon acceptance, the dataset will be publicly released on GitHub under the CC-BY 4.0 license.

Datasheet (uses format introduced by Gebru et al.)¹

Some question answers are omitted for the sake of anonymity.

Motivation

For what purpose was the dataset created? To explore human uncertainty judgments for ambiguous images and to produce more robust visual event classification models and uncertainty quantification techniques.

Who created the dataset, and on behalf of which entity? Omitted

Who funded the creation of the dataset? Omitted

Composition

What do the instances that comprise the dataset represent? Each instance is an image from a video depicting a specific event type.

How many instances are there in total? 12,000 images

Does the dataset contain all possible instances or is it a sample of instances from a larger set? The dataset does not contain all possible instances – it is a small sample of the possible instances in the domain.

What data does each instance consist of? Each instance consists of an image (denoted as a video ID and corresponding video frame number), a ground truth event type label, and human uncertainty scores associated with it (if applicable).

Is there a label or target associated with each instance? Yes, the ground truth event type and any human annotations are included.

Is any information missing from individual instances? Yes, a variety of additional metadata and annotations could potentially be included.

Are relationships between individual instances made explicit? Images belonging to the same video are labeled as such.

Are there recommended data splits? Due to the relatively small number of human annotations, it is recommended to use the set of images with such annotations as a validation or test set for most tasks.

Are there any errors, sources of noise, or redundancies in the dataset? Limitations are listed in Section 6 of the paper.

Is the dataset self-contained, or does it link to or otherwise rely on external resources? The dataset links to YouTube videos and videos from the Extended UCF Crime dataset.

Does the dataset contain data that might be considered confidential? No.

Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? Yes, the videos have not been checked for offensive or otherwise harmful content.

Does the dataset identify any subpopulations? No.

Is it possible to identify individuals, either directly or indirectly from the dataset? It may be possible. Videos may identify individual people, but the owners of the videos chose to upload

¹ Gebru, Timnit, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. "Datasheets for datasets." *Communications of the ACM* 64, no. 12 (2021): 86-92.

this content publicly and have the choice to take it down at any time (consequently removing it from the dataloader as well).

Does the data contain data that might be considered sensitive in any way? Yes, see above.

Collection process

How was the data associated with each instance acquired? Data collection details are provided in Section 3 of the paper.

What mechanisms or procedures were used to collect the data? See above.

If the dataset is a sample from a larger set, what was the sampling strategy? Once YouTube queries were made, the top videos according to the YouTube search algorithm were retrieved. Afterwards, irrelevant videos were manually removed.

Who was involved in the data collection process and how were they compensated? Amazon Mechanical Turk annotators produced the human uncertainty judgments and were paid \$0.20 per six judgments (approximately \$16/hr based on estimations).

Over what timeframe was the data collected? The data was collected between November 2021 and June 2022.

Were any ethical review processed conducted? No.

Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources? Human annotations were collected via Amazon Mechanical Turk.

Were the individuals in question notified about the data collection? Not explicitly.

Did the individuals in question consent to the collection and use of their data? Yes, via the Amazon Mechanical Turk participation agreement.

If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? N/A

Has an analysis of the potential impact of the dataset and its use on data subjects been conducted? No.

Preprocessing/cleaning/labeling

Was any preprocessing/cleaning/labeling of the data done? Ground truth and human judgment labels were added, and padding was added to images with unusual aspect ratios.

Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data? The video IDs to the source material are included in the dataset.

Is the software that was used to preprocess/clean/label the data available? No, but necessary steps are detailed in the paper.

Uses

Has the dataset been used for any tasks already? Yes, it has been used for the experiments detailed in the paper.

Is there a repository that links to any or all papers or systems that use the dataset? No

What (other) tasks could the dataset be used for? The dataset could be used for a variety of tasks including learning human uncertainty scoring functions, assessing uncertainty quantification approaches, developing more robust vision models, etc.

Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? Yes, dataset collection limitations are detailed in section 6 of the paper.

Are there tasks for which the dataset should not be used? Due to limitations listed in the paper, judgment should be used for tasks that may directly effect real people.

Distribution

Will the dataset be distributed to third parties outside of the entity on behalf of which the dataset was created? Yes, it will be publicly available.

How will the dataset be distributed? Via GitHub

When will the dataset be distributed? The data will be released once the paper is accepted and before the camera-ready version is submitted.

Will the dataset be distributed under a copyright or other intellectual property license, and/or under applicable terms of use? The dataset will be licensed under CC-BY 4.0

Have any third parties imposed IP-based or other restrictions on the data associated with the instances? The YouTube videos belong under the YouTube license.

Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? No

Maintenance

Who will be supporting/hosting/maintaining the dataset? The first author will maintain the dataset, and it will be hosted on GitHub.

How can the owner/curator/manager of the dataset be contacted? They can be contacted via email.

Is there an erratum? An erratum will be included in the repository if necessary.

Will the dataset be updated? The dataset will be updated to correct any errors that are found.

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances? Yes, all of the visual data will be retained as long as the individual creators of the videos used keep their content publicly available.

Will older versions of the dataset continue to be supported/hosted/maintained? No

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? Yes, they are free to fork the GitHub repository.