

592 **EgoSchema Datasheet**

593 **Motivation**

594

595 **For what purpose was the dataset created?** Was there a specific task in mind? Was
596 there a specific gap that needed to be filled? Please provide a description.

597 EgoSchema is a diagnostic benchmark for assessing very long-form video-language understanding
598 capabilities of modern multimodal systems. While some prior works have proposed video datasets
599 with long clip lengths, we posit that merely the length of the video clip does not truly capture the
600 temporal difficulty of the video task that is being considered. To remedy this, we introduce temporal
601 certificate sets, a general notion for capturing the intrinsic temporal understanding length associated
602 with a broad range of video understanding tasks & datasets. Please see Section 3.2 in the main paper
603 for more details.

604 **Who created this dataset (e.g., which team, research group) and on behalf of which
605 entity (e.g., company, institution, organization)?**

606 The authors created the dataset within the Malik Group at Berkeley AI Research, UC Berkeley. The
607 authors created it for the public at large without reference to any particular organization or institution.

608 **Who funded the creation of the dataset?** If there is an associated grant, please provide
609 the name of the grantor and the grant name and number.

610 EgoSchema creation is funded by the ONR MURI award number N00014-21-1-2801

611 **Any other comments?**

612 No

613 **Composition**

614

615 **What do the instances that comprise the dataset represent (e.g., documents, photos,
616 people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings;
617 people and interactions between them; nodes and edges)? Please provide a description.

618 Each instance in the dataset represents a 3-minute video and text that contains a question and five
619 answer options.

620 **How many instances are there in total (of each type, if appropriate)?**

621 EgoSchema has a total of 5063 instances each containing one video, one question, and five answer
622 options. You can see further statistics on the whole data on our website EgoSchema.github.io.

623 **Does the dataset contain all possible instances or is it a sample (not necessarily
624 random) of instances from a larger set?** If the dataset is a sample, then what is the
625 larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so,
626 please describe how this representativeness was validated/verified. If it is not representative
627 of the larger set, please describe why not (e.g., to cover a more diverse range of instances,
628 because instances were withheld or unavailable).

629 The video component of our dataset derives from the broader Ego4D dataset. For our research,
630 we selectively extracted non-overlapping three-minute segments from the Ego4D video data, each
631 segment consisting of a minimum of 30 human-annotated narrations (where each narration refers to a
632 timestamped sentence). Detailed statistic of the number of viable clips for different possible length
633 and narration density choices is discussed in Supplementary Section 7. The selected subset is very
634 diverse in human behavior as can be seen by the activity statistics presented on EgoSchema.github.io.

635

636 **What data does each instance consist of? “Raw” data (e.g., unprocessed text or**
637 **images) or features?** In either case, please provide a description.

638 Each instance in our dataset comprises raw mp4 video data, captured at a rate of 30 frames per second
639 and with a high resolution. Accompanying this video data, there are six text elements - one question
640 and five corresponding answer options one of which is marked as the correct answer to the question.
641

642 **Is there a label or target associated with each instance?** If so, please provide a
643 description.

644 Each instance is associated with a label ranging from 1 to 5 that indicates which of the five answer
645 options is correct.

646 **Is any information missing from individual instances?** If so, please provide a description,
647 explaining why this information is missing (e.g. because it was unavailable). This does not
648 include intentionally removed information but might include, e.g., redacted text.

649 All instances are complete.

650 **Are relationships between individual instances made explicit (e.g., users’ movie**
651 **ratings, social network links)?** If so, please describe how these relationships are made
652 explicit.

653 Some instances may have the same video but different questions and answers. It will be indicated by
654 a clip unique identifier in the final dataset.

655 **Are there recommended data splits (e.g., training, development/validation, testing)?** If
656 so, please provide a description of these splits, explaining the rationale behind them.

657 EgoSchema is designed specifically for zero-shot testing. Its primary purpose is to be able to asses
658 the out of the box long-term video-language understanding capabilities of modern multimodal models.
659

660 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please
661 provide a description.

662 The dataset was very carefully manually curated to mitigate any incidence of errors within the
663 questions and answers. Although different questions may be posed for the same clip, it is ensured
664 that there is no overlap between any two distinct clips. Further related details are also discussed in
665 the limitations section in the main paper.

666 **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
667 **(e.g., websites, tweets, other datasets)?** If it links to or relies on external resources, a)
668 are there guarantees that they will exist, and remain constant, over time; b) are there official
669 archival versions of the complete dataset (i.e., including the external resources as they
670 existed at the time the dataset was created); c) are there any restrictions (e.g., licenses,
671 fees) associated with any of the external resources that might apply to a future user? Please
672 provide descriptions of all external resources and any restrictions associated with them, as
673 well as links or other access points, as appropriate.

674 Entirety of the dataset will be made publicly available at our project website EgoSchema.github.io.
675 We will also provide a download tool for preprocessing all the videos such as cutting clips, associating
676 the question/answer text etc. Text will be released in a JSON format, hosted on our [github repository](#).
677 EgoSchema will be publicly released under the Ego4D license, which allows public use of the video
678 and text data for both research and commercial purposes.

679 **Does the dataset contain data that might be considered confidential (e.g., data that is**
680 **protected by legal privilege or by doctor-patient confidentiality, data that includes the**
681 **content of individuals non-public communications)?** If so, please provide a description.

682 No

683 **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
684 **threatening, or might otherwise cause anxiety?** If so, please describe why.

685 No

686 **Does the dataset relate to people?** If not, you may skip the remaining questions in this
687 section.

688 Some videos do contain people. However, the Ego4D authors employed an array of de-identification
689 procedures primarily centered on ensuring a controlled environment with informed consent from
690 all participants, and, where applicable, in public spaces with faces and other personally identifiable
691 information suitably obscured. We strictly import all RGB information from Ego4D without any
692 addition of our own.

693 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please
694 describe how these subpopulations are identified and provide a description of their respective
695 distributions within the dataset.

696 No

697 **Is it possible to identify individuals (i.e., one or more natural persons), either directly**
698 **or indirectly (i.e., in combination with other data) from the dataset?** If so, please
699 describe how.

700 No, Ego4D has employed an array of deidentification procedures in order to obscure any personally
701 identifiable information such as people's faces.

702 **Does the dataset contain data that might be considered sensitive in any way (e.g., data**
703 **that reveals racial or ethnic origins, sexual orientations, religious beliefs, political**
704 **opinions or union memberships, or locations; financial or health data; biometric or**
705 **genetic data; forms of government identification, such as social security numbers;**
706 **criminal history)?** If so, please provide a description.

707 No

708 **Any other comments?**

709 No

710

Collection Process

711

712 **How was the data associated with each instance acquired?** Was the data directly
713 observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or
714 indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses
715 for age or language)? If data was reported by subjects or indirectly inferred/derived from
716 other data, was the data validated/verified? If so, please describe how.

717 The video data, which is directly observable, was procured from the publicly accessible Ego4D
718 dataset. In contrast, the text data was generated through the use of Large Language Models (LLMs)
719 including GPT4, BARD, and Claude. These LLMs employed visual narrations from each video
720 within the Ego4D dataset to generate the corresponding text.

721 **What mechanisms or procedures were used to collect the data (e.g., hardware appa-**
722 **paratus or sensor, manual human curation, software program, software API)?** How were
723 these mechanisms or procedures validated?

724 The video and narration data were downloaded in accordance with the official Ego4D guidelines
725 for data access: <https://ego4d-data.org/docs/start-here/download-data>. For the generation of the text
726 data within our dataset, we utilized API access for GPT4 via OpenAI, for BARD via Google, and
727 for Claude via Anthropic. This allowed us to generate three distinct questions for each video clip
728 sampled from the Ego4D dataset. Upon the generation of these questions for each sampled video
729 clip, we implemented a series of filtering procedures including Rule-based filtering, Blind filtering,
730 and Manual curation. See Section 3.1.2 in the main paper for a more detailed explanation.

731 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,**
732 **deterministic, probabilistic with specific sampling probabilities)?**

733 The video component of our dataset derives from the broader Ego4D dataset. For our research,
734 we selectively extracted non-overlapping three-minute segments from the Ego4D video data, each
735 segment consisting of a minimum of 30 human-annotated narrations (where each narration refers to a
736 timestamped sentence). Detailed statistic of the number of viable clips for different possible length
737 and narration density choices is discussed in Supplementary Section 7.

738 **Who was involved in the data collection process (e.g., students, crowdworkers,**
739 **contractors) and how were they compensated (e.g., how much were crowdworkers**
740 **paid)?**

741 Our research utilized the services of Quantigo, a specialized data labelling company. The teams of
742 Quantigo employees that were based in Bangladesh were compensated at a rate of 5 dollars per hour,
743 at a wage significantly higher than the market hourly rate in Bangladesh. This was done to ensure
744 fair compensation for the complex tasks performed while also contributing to the highest quality
745 of the work delivered. It's important to note that our collaboration with Quantigo followed ethical
746 guidelines, with the fair treatment of all employees involved and the appropriate respect for their
747 expertise and labor. For exact instructions for human curation, see Supplementary Section 11.

748 **Over what timeframe was the data collected? Does this timeframe match the creation**
749 **timeframe of the data associated with the instances (e.g., recent crawl of old news**
750 **articles)?** If not, please describe the timeframe in which the data associated with the
751 instances was created.

752 The original videos within the Ego4D dataset were collected across various occasions spanning from
753 2019 to 2021. As for the EgoSchema, the textual information was collected over several sprints
754 during the first half of 2023 based on the Ego4D narrations.

755 **Were any ethical review processes conducted (e.g., by an institutional review board)?**
756 If so, please provide a description of these review processes, including the outcomes, as
757 well as a link or other access point to any supporting documentation.

758 No

759 **Does the dataset relate to people?** If not, you may skip the remaining questions in this
760 section.

761 Yes

762 **Did you collect the data from the individuals in question directly, or obtain it via third**
763 **parties or other sources (e.g., websites)?**

764 The video and narration data were acquired in accordance with the official Ego4D guidelines for data
765 access: <https://ego4d-data.org/docs/start-here/download-data>. The Ego4D authors had in turn ensured
766 consent of the people involved.

767 **Were the individuals in question notified about the data collection?** If so, please
768 describe (or show with screenshots or other information) how notice was provided, and
769 provide a link or other access point to, or otherwise reproduce, the exact language of the
770 notification itself.

771 Ego4d paper followed several procedures to ensure the preservation of privacy and the upholding of
772 ethical standards. Notably, these procedures included obtaining informed consent from those wearing
773 the cameras and adhering to de-identification requirements for personally identifiable information
774 (PII). Given that the video collection was conducted by Ego4D, we are not in a position to provide
775 specific instructions that were given to the camera wearers. The Ego4D privacy statement is available
776 at <https://ego4d-data.org/pdfs/Ego4D-Privacy-and-ethics-consortium-statement.pdf>

777 **Did the individuals in question consent to the collection and use of their data?** If so,
778 please describe (or show with screenshots or other information) how consent was requested
779 and provided, and provide a link or other access point to, or otherwise reproduce, the exact
780 language to which the individuals consented.

781 Ego4d paper privacy procedures have included obtaining informed consent from those wearing the
782 cameras. Given that the video collection was conducted by Ego4D, we are not in a position to provide
783 specific instructions that were given to the camera wearers. See [Ego4D privacy statement](#).

784 **If consent was obtained, were the consenting individuals provided with a mechanism
785 to revoke their consent in the future or for certain uses?** If so, please provide a
786 description, as well as a link or other access point to the mechanism (if appropriate).

787 Ego4d paper privacy procedures have included allowing camera users to ask questions and withdraw
788 at any time. Additionally, they were free to review and redact their own video. Given that the
789 video collection was conducted by Ego4D, we are not in a position to provide specific instructions
790 that were given to the camera wearers. You can find the Ego4D privacy statement at [https://ego4d-
791 data.org/pdfs/Ego4D-Privacy-and-ethics-consortium-statement.pdf](https://ego4d-data.org/pdfs/Ego4D-Privacy-and-ethics-consortium-statement.pdf).

792 **Has an analysis of the potential impact of the dataset and its use on data subjects
793 (e.g., a data protection impact analysis) been conducted?** If so, please provide a
794 description of this analysis, including the outcomes, as well as a link or other access point
795 to any supporting documentation.

796 While we recognize the importance of this topic, we would, once more, refer to the Ego4D paper
797 for an in-depth discussion. Ego4D acknowledges the potential privacy risks associated with the
798 use of wearable devices in data collection and has taken several steps such as depersonalizing any
799 sensitive information, blurring out faces and bodies, etc. towards maintaining privacy. The same
800 carries over to the video data in EgoSchema as well. Broadly, very long-form video understanding is
801 a core capability for agents that are to perceive the natural visual world. Hence, developing datasets
802 such as EgoSchema will be critical to unlocking this key AI capability. Additionally, according to
803 [Ego4D privacy statement](#), all videos from Ego4D were reviewed by an approved member of one of
804 the participant's universities or institutes to identify and assess potential privacy concerns.

805 **Any other comments?**

806 No

807

808 **Preprocessing/cleaning/labeling**

809

810 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or
811 bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of**

812 **instances, processing of missing values)?** If so, please provide a description. If not, you
813 may skip the remainder of the questions in this section.

814 The set of generated questions and answers from output was filtered by those LLMs and finally
815 curated by humans. A detailed description can be found in Section 3. There was no preprocessing
816 done on the video clips sampled from Ego4D.

817 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g.,**
818 **to support unanticipated future uses)?** If so, please provide a link or other access point
819 to the “raw” data.

820 Human curation was employed to rectify errors in the question-answer sets, particularly cases where
821 the identified correct answer was wrong or a wrong answer was actually correct. Given the crucial
822 role of this step in ensuring the accuracy of our dataset, we do not find it necessary to release a version
823 of the dataset prior to human curation. However, all the discarded "raw" data is indeed also saved.

824 **Is the software used to preprocess/clean/label the instances available?** If so, please
825 provide a link or other access point.

826 The APIs for the Large Language Models (LLMs) are publicly accessible. The prompts for filtering
827 and instructions for human curation are provided in Supplementary Section 6 and Supplementary
828 Section 11 respectively. Additionally all necessary code for generation, filtering etc. is provided in
829 the supplementary materials.

830 **Any other comments?**

831 No

832

833

Uses

834

835 **Has the dataset been used for any tasks already?** If so, please provide a description.

836 We have used our dataset as a benchmark for several video question-answering models (please see
837 §4.2 for more details). Additionally, we will be providing the code for zero-shot evaluation on our
838 project github repository.

839 **Is there a repository that links to any or all papers or systems that use the dataset?** If
840 so, please provide a link or other access point.

841 It will be made public on our [website](#) once more papers will start to use our dataset.

842 **What (other) tasks could the dataset be used for?**

843 Besides multiple-choice video question-answering tasks our dataset also can be used for open-ended
844 video question-answering. There might be several more video-language tasks that could also be
845 perhaps benchmarked, but currently we leave that exploration for future work.

846 **Is there anything about the composition of the dataset or the way it was collected**
847 **and preprocessed/cleaned/labeled that might impact future uses?** For example, is
848 there anything that a future user might need to know to avoid uses that could result in unfair
849 treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other
850 undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is
851 there anything a future user could do to mitigate these undesirable harms?

852 No

853 **Are there tasks for which the dataset should not be used?** If so, please provide a
854 description.

855 No

856 **Any other comments?** No

857

Distribution

858

859

860 **Will the dataset be distributed to third parties outside of the entity (e.g., company,**
861 **institution, organization) on behalf of which the dataset was created?** If so, please
862 provide a description.

863 The dataset will be made publicly available and can be used for both research and commercial
864 purposes under the Ego4D license.

865 **How will the dataset be distributed (e.g., tarball on website, API, GitHub) Does the**
866 **dataset have a digital object identifier (DOI)?**

867 The dataset will be distributed as a JSON file describing the unique identifier for each clip, the
868 associated question, the five answer options, the label, and additional clip information that facilitates
869 the tracing of the clip back to the original Ego4D data, such as the Ego4D video identification of the
870 clip's source video, among other details. In addition, download tools to acquire and pre-process the
871 video RGB data will also be provided on our website.

872 **When will the dataset be distributed?**

873 The full dataset will be made available upon the acceptance of the paper before the camera-ready
874 deadline.

875 **Will the dataset be distributed under a copyright or other intellectual property (IP)**
876 **license, and/or under applicable terms of use (ToU)?** If so, please describe this license
877 and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant
878 licensing terms or ToU, as well as any fees associated with these restrictions.

879 EgoSchema will be publicly released under the Ego4D license, which allows direct public use of the
880 video and text data for both research and commercial purposes.

881 **Have any third parties imposed IP-based or other restrictions on the data associated**
882 **with the instances?** If so, please describe these restrictions, and provide a link or other
883 access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees
884 associated with these restrictions.

885 No

886 **Do any export controls or other regulatory restrictions apply to the dataset or to**
887 **individual instances?** If so, please describe these restrictions, and provide a link or other
888 access point to, or otherwise reproduce, any supporting documentation.

889 No

890 **Any other comments?**

891 No

892

Maintenance

Who will be supporting/hosting/maintaining the dataset?

The authors of the paper will be maintaining the dataset, pointers to which will be hosted on github repo <https://github.com/egoschema/EgoSchema> along with the code for download and preprocessing tool, with the actual data hosted either on Amazon AWS as an S3 bucket or as a google drive folder.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

We will post the contact information on our website. We will be available through github issues as well as through email.

Is there an erratum? If so, please provide a link or other access point.

We will host an erratum on the Github repo in the future, to host any approved errata suggested by the authors or the video research community.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?

Yes, we plan to host an erratum publicly. There are no specific plans for a v2 version, but there does seem plenty oppurtunities for exciting future dataset work based on EgoSchema.

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.

No.

Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to users.

N/A There are no older versions at the current moment. All updates regarding the current version will be communicated via our website.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description.

Contributions will be made possible using standard open-source tools, submitted as pull requests to the relevant GitHub repository. Moreover, we will provide information on how to trace sampled clips back to their original source within the Ego4D dataset. This will enable users to access additional Ego4D data, such as narrations, summaries, and object detections, as applicable.

Any other comments?

No