
VERIFIED: A Video Corpus Moment Retrieval Benchmark for Fine-Grained Video Understanding (Dataset Document)

1 A Datasheets

2 1. What is the name of the dataset?

3 The names of the datasets are Charades-FIG, DiDeMo-FIG, and ActivityNet-FIG. They contain fine-grained video moment captions for videos collected in previous Charades-STA [1],
4 DiDeMo [2], and ActivityNet Captions [3]. They have the same organizational structure. The
5 following descriptions are suitable for each of the datasets.
6

7 2. For what purpose was the dataset created?

8 The dataset was created to advance the research for fine-grained video understanding, including
9 but not limited to video (corpus) moment retrieval, dense video captioning, and video question
10 answering. It should be beneficial to controllable video generation as well and more work
11 exploring this is highly encouraged.

12 3. What do the instances that comprise the dataset represent?

13 Each instance consists of a fine-grained text annotation with a confidence score (having selected
14 the annotation with the highest score) for its corresponding video moment.

15 4. How many instances are there in total?

16 Charades-FIG contains 16128 instances (12408 in the train split and 3720 in the test split);
17 DiDeMo-FIG contains 41091 instances (32903 in the train split, 4175 in the validation split, and
18 4013 in the test split); ActivityNet-FIG contains 71656 instances (37261 in the train split, 17429
19 in the validation 1 split, and 16966 in the validation 2 split).

20 5. How was the data collected?

21 The dataset is annotated on the videos collected by the previous dataset.

22 6. How was the data labeled?

23 The dataset is labeled by the proposed VERIFIED pipeline.

24 7. Was the raw data saved?

25 The access to the raw videos should comply with the previous video collector, but we will provide
26 ResNet152 features extracted by ourselves, which are accessible.

27 8. What are the intended uses of the dataset?

28 The dataset is primarily intended for training and evaluating the fine-grained video corpus moment
29 retrieval task but is also suitable for other video tasks.

30 9. Will the dataset be distributed to third parties outside of the entity on behalf of which the 31 dataset was created?

32 Yes, the dataset will be available under a Creative Commons license.

33 **10. Who is supporting/hosting/maintaining the dataset?**

34 The dataset is maintained by the research lab that created it.

35 **B Metadata**

36 The dataset is organized by JSONL, where each instance occupies a JSON body. For each instance,
37 the field descriptions contain:

- 38 1. "desc_id": a unique ID for each instance.
- 39 2. "video": the URL of the video.
- 40 3. "duration": the duration of the video.
- 41 4. "time": the start and end timestamps(unit: second) of the annotated video moment.
- 42 5. "cog_desc": the previous coarse-grained annotation text.
- 43 6. "fig_desc": the fine-grained annotation text generated by our proposed VERIFIED pipeline.
- 44 7. "fig_desc_score": the confidence score for the fine-grained annotation text.

45

46

47

48

49

50

51

52

53

54

55

57 **D Data License**

58 This work is licensed under a CC BY 4.0 license.

59 **References**

- 60 [1] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization
61 via language query. In *Proceedings of the IEEE international conference on computer vision*,
62 pages 5267–5275, 2017.
- 63 [2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan
64 Russell. Localizing moments in video with natural language. In *Proceedings of the IEEE*
65 *international conference on computer vision*, pages 5803–5812, 2017.
- 66 [3] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning
67 events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages
68 706–715, 2017.

¹<https://prior.allenai.org/projects/charades>

²<https://github.com/LisaAnne/LocalizingMoments>

³<http://activity-net.org/>