

# A Large-scale Study of Spatiotemporal Representation Learning with a New Benchmark on Action Recognition (*Supplementary Material*)

Andong Deng\*    Taojiannan Yang\*    Chen Chen  
Center for Research in Computer Vision  
University of Central Florida, USA

andong.deng@ucf.edu, taoyang1122@knights.ucf.edu, chen.chen@crcv.ucf.edu

## 1. Datasets

### 1.1. Sports datasets

Sports-related videos are one of the most numerous videos in human visual records. Sports videos could contain a variety of different action patterns; they can be as simple as jogging and jumping, or as more complicated as some professional actions like cross-over in basketball games, all of which could be qualified learning samples for various action recognition models. In our sports datasets section, we selected three representative datasets: Sports1M [11], MOD20 [16], and FineGym[21]. **Sports1M** is one of the biggest sports video datasets in the vision community, which contains 487 categories and has been well-annotated. Considering the fact that some of its URLs are no longer available as well as its huge amount (the original version contains over 1 million videos), we construct a mini version of it, which only includes 50 samples per class, 40 for training, and 10 for testing. **MOD20** is a multi-viewpoint outdoor dataset collected from both YouTube videos and a drone camera, which alleviates the dataset scarcity in terms of viewpoint. Specifically, there are totally four types of views in MOD20, three of which are above the person and the fourth is an elevation view. Furthermore, these two datasets only contain coarse sports categories, ignoring the analysis of sub-actions within a sports event, which may weaken the scope of our benchmark. Considering this, we include a recently released fine-grained dataset, **FineGym**. In practice, we use one of its versions – FineGym99, which is composed of 99 fine-grained gymnastic actions from top-level world competitions.

### 1.2. Daily datasets

We construct our daily action datasets based on two rules: abundant first-person-video data and diverse activity categories. Studies on egocentric video analysis could help us to forge an in-depth understanding of the interac-

tions between humans and surroundings, which is essential for many cutting-edge AI technologies such as embodied AI. **CharadesEgo** [22] is a large-scale dataset with paired first- and third-person videos to facilitate the investigation of the intrinsic correspondence between different views for the same action. In our benchmark, we choose to only carry out the evaluation using its 1st person part, since we do not focus on the correlations between two different views. Based on its official temporal annotations, we manage to segment the original videos into 43,594 short clips. **HACS** [28], human action clips and segments, is a large-scale dataset for both action recognition and temporal action localization. We only leverage HACS clips for our action recognition studies; likewise, considering its scale, we randomly sample 50 videos per class to form mini-HACS, 40 for training and 10 for testing. In addition, we also include two more self-collected datasets based on videos captured from real daily events participants, Toyota Smarthome [3] and MPII Cooking [19]. **Toyota Smarthome** is a 3rd view dataset containing videos from different cameras deployed in an apartment, whose subjects are 18 senior people. The videos are collected from 7 cameras in the dining room, kitchen, and living room. We use the cross-subject train-test split in our evaluation, i.e., the training data are from 11 subjects and the rest are utilized for testing. **MPII Cooking** dataset is a fine-grained cooking activity dataset, which is originally built for action detection. In our benchmark, we obtain action clips containing one action label given the official temporal annotation. The raw videos are recorded based on 12 participants, and we use videos from 10 subjects as the training set.

### 1.3. Anomaly datasets

Action analysis for anomaly or crime-related videos is an important real-world application, and one of the objectives of our proposed benchmark is to provide practical guidance for the application scenario of human action recognition models. Hence, we build the anomaly track with three representative datasets to evaluate the performance of various

---

\*Equal Contribution.

models in such a realistic case. **UCF-Crime** [23] is a challenging anomaly video dataset collected from surveillance cameras. We select 12 human-related crime categories from its original 14-class recognition version as we only focus on human actions. **XD-violence** [27] is another video anomaly dataset that includes data from various sources such as action movies, sports videos, and CCTV cameras. This results in a more extensive collection of video samples, enriching our anomaly datasets track. Specifically, the original XD-violence has a multi-label text set with temporal annotation, according to which we segment the test videos into single-label clips. We also include a recently released fall detection dataset **MUVIM** [5] (Multi Visual Modality Fall Detection Dataset), which consists of visual data from multiple sensors: infrared, depth, RGB, and thermal cameras. Considering the data consistency, we only utilize their RGB version.

### 1.4. Instructional datasets

Instructional videos are captured in order to guide people to accomplish particular operations, e.g., assembling some objects with the components, operating a clinical surgery, or other necessary tasks which require additional knowledge. Accurate video analysis for such instructional videos is an important and irreplaceable phase for many practical applications, such as intelligent robots for industrial or medical usage. **COIN** [24] dataset is a large-scale dataset built for comprehensive instructional video analysis based on videos collected from YouTube. It consists of 180 tasks in 12 different domains related to tasks about daily living, e.g., 'change the car tire' and 'replace the door knob'. **INHARD** [2], Industrial Human Action Recognition Dataset, is collected in a human-robot collaboration scenario. 16 distinct subjects are invited to finish an assembly task with the guidance of a robotic arm. The classes contain the specific actions during this operation, such as 'put down measuring rod' and 'put down component'. Similarly, we include another instructional dataset **MECCANO** [17], which is also related to an assembly operation but collected with wearable cameras. The target task is to build a toy motorbike given all the components and the booklet, and the whole assembly process is precisely divided into 61 action steps. To cover scenarios as much as possible, we add two medical instructional datasets **MISAW** [10] (Micro-Surgical Anastomose Workflow) and **PETRAW** [9] (PEg TRAnSfer Workflow). Both two datasets are collected in simulated environments. The whole process is constructed by step-wise professional clinical operations, such as 'suturing' and 'knot tying'. Both two datasets provide frame-wise annotation in terms of phase, step, and action labels for the left hand and right hand and MISAW additionally provides the target and tool annotations. To generate segment-level action annotation, for MISAW, we view the action label of the left hand and the corresponding target as a whole, when any of them

Table 1: The training details of our supervised pre-training.

HyperParams	TSN	TSM	I3D	NL	TimeSformer	VideoSwin
Batch Size	64	64	64	64	64	64
lr	0.05	0.05	0.01	0.01	5e-3	5e-4
lr policy	StepLR	StepLR	StepLR	StepLR	StepLR	CosineLR
lr step	[20, 40]	[20, 40]	[20, 40]	[20, 40]	[10, 20]	/
# Epoch	50	50	50	50	30	30
# WarmUp	/	/	10	10	/	2.5
Optimizer	SGD	SGD	SGD	SGD	SGD	AdamW

changes, we change the segment annotation. For instance, if the annotation of the action and the target for the current frame are 'Hold' and 'Left artificial vessel', we annotate the segment where it belongs as 'Hold Left artificial vessel'; if the annotation changes in the next frame into 'Catch' and 'Needle', we start a new segment and annotate it as 'Catch Needle' until the next change. Similarly, we segment PE-TRAW videos based on the change of the left-hand action.

### 1.5. Gesture datasets

Gesture recognition is critical for application in human-computer interfaces and has become an appealing topic in recent years. In this track, we view all datasets whose data semantic originated from symbols constructed by human body parts as the generalized gesture datasets, e.g., gestures, sign language, and other body language or pose. **Jester** [15] is collected from 1,376 actors based on 27 gesture classes. The categories contained in Jester include gestures that usually appear in interactions between humans and some smart devices, such as "Zoom in with two fingers". **WLASL** [12], short for World-Level American Sign Language, is built for sign language understanding, which could make progress for the communications of the blind and deaf. Its original version contains 2000 categories of common sign language; in our benchmark, we use its subset WLASL100 for our evaluation. Besides, **UAV Human** dataset [14] contains videos captured from unmanned aerial vehicles, which also includes body sign language, e.g., the victory sign posed by two arms. For completeness, we also keep other regular classes of this dataset.

## 2. Training Details

### 2.1. Supervised pre-training

We pre-train all of the 6 models on Kinetics400. Specifically, the total training epochs for CNNs and transformers are 50 and 30, respectively. For VideoSwin, we utilize AdamW as the optimizer, while we use SGD for the rest of the models. In testing, we adopt single-view evaluation for all datasets. The frame sampling strategy is sampling 8 frames in total and sampling 1 frame per 16 frames. The weight decay is 1e-4 and momentum is 0.9 for all models. The complete training details are shown in Table 1.

Table 2: The training details of our self-supervised pre-training.

HyperParams	TSN	TSM	I3D	NL	TimeSformer	VideoSwin
Batch Size	64	64	64	64	64	64
lr	0.05	0.05	0.1	0.01	5e-4	5e-4
lr policy	CosineLR	CosineLR	CosineLR	CosineLR	CosineLR	CosineLR
# Epochs	50	50	50	50	50	50
# WarmUp	10	10	10	10	10	10
Optimizer	SGD	SGD	SGD	SGD	AdamW	AdamW

## 2.2. Self-supervised pre-training

We utilize  $\rho$ MoCo [6], which is an extension of MoCo [7] in the video domain, as our self-supervised pre-training method. Specifically,  $\rho$ MoCo, where  $\rho$  stands for the number of temporal views for contrastive learning, aims at learning an encoder that could generate invariant features for different clips of the same video. In our pre-training, we set  $\rho=2$ . The weight decay is 1e-4 and momentum is 0.9. The complete training details are presented in Table 2.

## 2.3. Standard Finetuning

Similarly, we adopt the same training settings in supervised pre-training for our standard finetuning experiments except for the frame interval. Considering the average frame number of the clips could vary across datasets, we select different frame intervals for each dataset according to its average frame number. The details are shown in Table 3. Other training settings are the same as standard finetuning.

Table 3: Sampling frame interval of different datasets.

Dataset	Avg. # Frames	Frame Interval
XD-Violence	517	16
UCF-Crime	402	16
MUVIM	194	16
WLASL100	68	8
Jester	36	4
UAV Human	133	16
CharadesEgo	298	16
Toyota Smarthome	248	16
MPII Cooking	78	8
Mini-Sports1M	711	16
FineGym99	74	8
MOD20	216	16
COIN	160	16
MECCANO	21	2
INHARD	71	8
PETRAW	65	8
MISAW	117	16

## 2.4. Few-shot learning

Since different training sample choices can have a large impact on few-shot learning, we randomly generate 3 train-

ing splits for all datasets and report the average performance to reduce such variation. And there are some datasets that contain categories that only have 1 or 2 samples, we ignore the training number shortage. The training setting is the same as the standard finetuning.

## 3. Standard finetuning

We showcase the complete finetuning results in Table 4 and Table 5, which include both top-1 and top-5 accuracy based on both supervised pre-training and self-supervised pre-training. We only provide top-1 results in the main paper. Besides, we also attach existing SoTA results, if any, for each dataset in BEAR for reference in Table 4. For MUVIM, MPII-Cooking, and InHARD, there are no existing reported results for action recognition. For XD-Violence and CharadesEGO, the original data contains multiple labels for one long video, but we segment them into single-label clips via their temporal annotation. Besides, the SoTA results of Mini-HACS and Mini-Sports1M are actually from the complete version, we only include them for reference.

## 4. Few-shot learning

We showcase the complete few-shot learning results in Table 6, Table 8, Table 10, Table 7, Table 9 and Table 11, which include both top-1 and top-5 accuracy for all the 3 training split based on both supervised pre-training and self-supervised pre-training.

Table 4: Top-1 and top-5 accuracy of finetuning based on supervised pre-training and SoTA results for each dataset.

Dataset	TSN	TSM	I3D	NL	TimeSFormer	VideoSwin	SoTA
<b>XD-Violence</b>	<b>85.54/NA</b>	82.96/NA	79.93/NA	79.91/NA	82.51/NA	82.40/NA	–
<b>UCF-Crime</b>	35.42/77.78	<b>42.36/79.17</b>	31.94/77.08	34.03/81.94	36.11/76.39	34.72/77.78	28.4[23]
<b>MUVIM</b>	79.30/NA	<b>100/NA</b>	97.80/NA	98.68/NA	94.71/NA	<b>100.00/NA</b>	–
<b>WLASL</b>	29.63/62.96	43.98/77.31	49.07/78.70	<b>52.31/78.24</b>	37.96/73.61	45.37/75.46	83.30[8]
<b>Jester</b>	86.31/99.66	<b>95.21/99.77</b>	92.99/99.68	93.49/99.66	93.42/99.61	94.27/99.68	98.15[26]
<b>UAV-Human</b>	27.89/50.82	<b>38.84/61.47</b>	33.49/59.74	33.03/54.21	28.93/51.14	38.66/61.42	37.98 [1]
<b>CharadesEGO</b>	8.26/30.38	8.11/29.49	6.13/21.86	6.42/22.03	<b>8.58/29.96</b>	8.55/29.86	–
<b>Toyota Smarthome</b>	74.73/95.95	<b>82.22/96.74</b>	79.51/95.60	76.86/94.52	69.21/93.30	79.88/97.13	71.0 [4]
<b>Mini-HACS</b>	84.69/98.04	80.87/96.48	77.74/95.17	79.51/95.17	79.81/96.48	<b>84.94/97.58</b>	95.5 [13]
<b>MPII Cooking</b>	38.39/71.17	46.74/74.96	<b>48.71/74.05</b>	42.19/70.11	40.97/68.44	46.59/80.88	–
<b>Mini-Sports1M</b>	54.11/80.74	50.06/76.57	46.90/72.85	46.16/72.77	51.79/77.15	<b>55.34/80.18</b>	75.5 [25]
<b>FineGym</b>	63.73/94.60	<b>80.95/98.49</b>	72.00/96.14	71.21/95.94	63.92/93.88	65.02/92.89	80.4 [21]
<b>MOD20</b>	<b>98.30/99.86</b>	96.75/100	96.61/100	96.18/100	94.06/99.72	92.64/99.72	74.0 [16]
<b>COIN</b>	81.15/96.19	78.49/95.24	73.79/92.58	74.30/92.07	<b>82.99/96.70</b>	76.27/93.53	88.02 [24]
<b>MECCANO</b>	<b>41.06/75.20</b>	39.28/70.88	36.88/67.45	36.13/66.63	40.95/75.17	38.89/72.19	42.85 [17]
<b>InHARD</b>	84.39/98.99	<b>88.08/98.99</b>	82.06/98.63	86.31/98.99	85.16/99.23	87.60/99.11	–
<b>PETRAW</b>	94.30/99.92	95.72/99.97	94.84/99.92	94.54/99.85	94.30/99.87	<b>96.43/99.90</b>	88.51 [9]
<b>MISAW</b>	61.44/94.34	<b>75.16/97.17</b>	68.19/96.08	64.27/95.64	71.46/96.65	69.06/97.17	63.4 [10]

Table 5: Top-1 and top-5 accuracy of finetuning based on self-supervised pre-training.

Dataset	TSN	TSM	I3D	NL	TimeSFormer	VideoSwin
<b>XD-Violence</b>	80.49/NA	<b>81.73/NA</b>	80.38/NA	80.94/NA	77.47/NA	77.91/NA
<b>UCF-Crime</b>	<b>37.50/83.33</b>	35.42/81.94	34.03/80.56	34.72/83.33	36.11/77.78	34.03/80.56
<b>MUVIM</b>	99.12/NA	<b>100/NA</b>	66.96/NA	66.96/NA	99.12/NA	<b>100/NA</b>
<b>WLASL</b>	27.01/50.22	27.78/53.70	29.17/65.74	<b>30.56/62.50</b>	25.56/59.44	28.24/65.28
<b>Jester</b>	83.22/99.23	<b>95.32/99.78</b>	87.23/99.46	93.89/99.68	90.33/99.41	90.18/99.40
<b>UAV-Human</b>	15.70/35.89	30.75/55.07	31.95/56.69	26.28/51.09	21.02/44.28	<b>35.12/59.47</b>
<b>CharadesEGO</b>	6.29/23.52	6.59/23.81	6.24/22.25	6.31/22.74	7.59/27.81	<b>7.65/27.41</b>
<b>Toyota Smarthome</b>	68.71/91.61	<b>81.34/96.63</b>	77.82/95.53	76.16/93.58	61.64/91.44	80.18/97.00
<b>Mini-HACS</b>	64.60/90.38	63.24/90.38	70.24/92.20	60.57/86.05	73.92/95.73	<b>75.58/95.62</b>
<b>MPII Cooking</b>	34.45/66.16	<b>50.08/75.42</b>	42.79/72.08	40.36/70.56	35.81/64.34	47.19/76.33
<b>Mini-Sports1M</b>	43.02/71.23	43.59/70.86	46.28/73.12	45.56/72.07	44.60/72.44	<b>47.60/73.88</b>
<b>FineGym</b>	54.62/91.21	<b>75.87/97.84</b>	69.62/95.57	68.79/95.70	47.60/85.76	58.94/92.56
<b>MOD20</b>	91.23/98.87	92.08/99.29	91.94/99.58	92.08/99.58	90.81/99.43	<b>92.36/99.58</b>
<b>COIN</b>	61.48/88.52	64.53/89.85	71.57/92.20	<b>72.78/92.58</b>	67.64/89.97	68.78/89.28
<b>MECCANO</b>	32.34/65.92	35.10/65.99	34.86/66.99	33.62/66.03	33.30/67.87	<b>37.80/72.30</b>
<b>InHARD</b>	75.63/97.68	<b>87.66/99.46</b>	82.54/98.87	80.81/98.63	71.28/97.38	80.10/98.87
<b>PETRAW</b>	93.18/99.87	<b>95.51/99.92</b>	95.02/99.92	94.38/99.87	85.56/99.38	91.46/99.90
<b>MISAW</b>	59.04/90.20	<b>73.64/97.82</b>	70.37/96.73	64.27/94.55	60.78/97.39	68.85/97.39

Table 6: Top-1 and top-5 accuracy of few-shot learning based on supervised pre-training on training split. NL means NonLocal network and TSF means TimeSformer.

#Shot	XD-Violence	UCF-Crime	MUVM	WLASL	Jester	UAV Human	CharadesEGO	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
16	65.81/NA	38.89/84.03	90.75/NA	35.05/64.35	31.16/68.39	17.04/40.07	3.80/15.69	39.79/80.03	79.86/96.17	14.87/43.10	45.52/74.52	9.87/34.41	94.34/99.86	77.28/94.86	3.68/14.28	16.69/60.13	33.63/92.82	28.54/74.07
8	61.21/NA	33.37/56.69	55.30/NA	15.74/43.52	20.91/52.08	9.56/26.89	2.75/12.53	27.37/70.75	76.08/95.52	15.02/39.96	39.86/69.70	6.54/25.62	95.19/99.86	71.07/92.45	1.77/8.68	15.32/52.44	27.12/83.38	19.17/58.17
4	58.86/NA	24.31/74.31	41.41/NA	5.09/15.28	12.06/38.57	4.34/14.45	2.67/10.24	14.38/46.36	70.64/92.09	11.53/30.96	30.86/60.72	4.87/21.65	93.21/99.86	62.88/89.66	1.49/8.01	15.38/47.79	21.73/77.42	13.94/48.58
2	57.96/NA	20.85/68.06	37.44/NA	1.89/6.94	9.50/24.68	1.82/6.67	1.27/8.11	10.79/37.59	62.59/89.63	4.86/29.76	22.85/50.76	2.28/8.50	85.71/99.72	51.52/83.06	0.74/7.01	10.67/41.84	23.99/78.68	10.02/35.95
16	65.47/NA	42.36/79.84	99.56/NA	47.69/78.24	36.86/76.49	20.20/43.08	4.26/16.44	46.90/83.23	76.38/94.41	21.40/50.23	42.94/70.00	15.43/45.98	93.92/99.72	72.08/92.20	4.71/16.12	22.47/73.78	61.24/97.82	37.47/76.69
8	59.19/NA	29.17/79.17	68.28/NA	26.85/58.80	21.82/55.45	12.84/31.51	2.82/11.47	38.14/75.06	67.59/94.21	19.73/40.67	37.06/64.52	8.90/27.29	94.20/99.86	66.94/88.90	3.75/10.52	22.05/69.55	51.67/93.79	29.19/76.91
4	54.60/NA	25.00/70.83	64.32/NA	9.72/25.93	12.36/37.49	6.08/19.36	2.31/9.48	24.15/56.99	67.52/92.30	8.80/24.89	27.70/56.58	6.62/25.41	88.26/99.58	57.17/85.74	2.44/10.52	18.71/66.09	41.61/80.53	14.60/51.63
2	50.22/NA	22.92/71.53	61.23/NA	4.17/12.50	9.92/24.89	2.82/9.63	2.08/7.77	13.80/43.94	60.67/88.52	7.74/31.11	21.58/48.28	6.13/21.71	83.31/98.59	45.56/76.84	2.90/10.17	16.09/58.76	40.46/83.79	9.89/36.17
16	58.07/NA	33.33/80.56	85.46/NA	49.07/79.63	35.22/74.06	19.13/43.10	2.90/11.44	45.32/82.74	76.69/95.17	17.91/43.55	39.36/67.21	14.48/46.29	94.20/100	66.88/90.10	4.18/14.74	29.68/77.29	48.36/94.07	39.22/83.44
8	51.57/NA	27.08/70.14	65.29/NA	29.63/66.67	25.11/60.70	11.69/30.10	2.55/10.28	35.93/74.32	74.27/93.81	11.53/40.52	34.15/61.42	9.44/33.62	94.34/100	61.17/86.55	2.50/10.24	21.33/62.57	42.18/85.99	21.79/67.54
4	56.17/NA	21.53/67.36	32.16/NA	10.65/33.33	13.79/42.29	5.76/19.05	2.06/8.93	20.12/56.38	69.94/90.89	13.96/34.14	27.21/54.62	7.41/27.26	88.68/99.72	53.17/82.11	1.52/8.40	20.98/57.33	28.50/76.19	13.51/56.21
2	48.26/NA	18.06/66.67	37.00/NA	4.17/16.20	10.46/44.38	3.02/9.83	2.05/8.12	18.28/48.09	62.39/86.92	6.22/24.58	20.13/34.91	5.67/23.13	80.87/99.87	41.12/72.97	2.66/12.65	19.01/60.85	32.94/79.01	6.75/32.68
16	58.52/NA	31.97/81.25	95.39/NA	51.39/83.33	38.16/76.95	18.65/41.09	2.93/11.09	42.38/78.65	76.28/94.01	19.58/47.34	38.97/66.02	15.03/44.06	94.34/99.86	68.46/90.55	5.32/17.43	25.27/72.94	56.77/96.05	33.12/69.93
8	54.04/NA	25.00/72.22	54.24/NA	32.87/68.52	26.84/63.62	11.12/29.52	2.24/9.50	34.68/70.02	74.67/93.10	17.45/42.64	33.31/60.99	9.02/27.52	90.38/100	60.47/87.06	5.01/16.10	20.56/65.38	46.79/90.58	22.44/64.80
4	51.01/NA	17.36/67.36	54.19/NA	9.72/34.72	15.29/44.98	5.53/17.49	2.10/8.73	20.80/55.02	60.03/91.59	10.77/29.29	26.51/54.07	6.90/25.57	85.15/99.86	52.60/81.35	4.53/15.69	21.22/61.28	35.89/78.32	11.34/49.02
2	44.39/NA	17.36/64.58	45.37/NA	4.63/15.28	10.70/35.80	2.85/8.75	1.87/7.46	14.23/46.88	60.83/88.22	9.10/30.52	20.18/45.77	5.65/22.38	81.19/99.91	41.05/73.35	4.29/16.33	13.83/48.03	34.12/79.30	7.41/28.76
16	67.71/NA	32.64/77.78	65.56/NA	40.74/72.69	30.47/68.33	14.19/32.95	3.47/11.52	33.37/74.01	74.67/94.66	21.09/48.56	45.05/72.61	11.84/39.33	90.81/99.92	79.19/95.11	4.50/12.22	14.78/53.81	44.54/94.20	38.56/78.43
8	60.76/NA	31.25/68.75	57.27/NA	22.22/55.09	13.53/41.70	6.51/19.11	2.99/10.45	18.74/52.59	67.72/91.94	13.66/35.81	39.88/67.60	7.63/28.89	88.53/99.86	75.63/93.15	2.27/8.25	12.34/56.44	32.84/82.35	14.60/47.81
4	53.48/NA	23.61/58.33	38.33/NA	7.41/21.30	8.89/31.70	2.76/9.95	2.45/8.69	9.96/47.73	60.07/86.81	7.74/27.77	31.70/59.98	5.89/20.78	84.44/99.29	68.40/90.23	1.66/7.05	12.99/53.28	32.04/82.94	14.16/44.88
2	47.65/NA	15.28/68.06	48.90/NA	2.78/6.02	6.66/27.00	1.37/5.21	2.98/7.59	11.03/35.21	49.65/76.89	5.31/33.69	44.99/20.14	4.49/20.14	72.56/99.71	51.90/82.49	1.56/7.62	9.83/44.70	31.04/73.91	16.12/42.70
16	59.30/NA	34.03/77.78	93.36/NA	43.06/75.93	41.89/83.89	22.74/47.72	3.47/11.52	46.94/87.43	81.22/96.93	21.70/51.59	44.35/70.92	11.79/39.44	82.89/98.16	65.23/89.15	4.68/15.41	23.96/68.00	37.20/87.53	30.50/71.24
8	60.65/NA	31.94/72.92	65.20/NA	31.02/61.11	30.32/70.60	15.59/37.57	2.83/10.94	36.96/77.21	76.54/96.02	14.26/38.54	37.68/66.18	8.38/31.63	77.79/99.17	60.22/83.38	2.90/12.43	18.41/53.16	25.99/79.71	27.67/70.15
4	57.96/NA	14.98/65.97	48.90/NA	10.19/20.56	16.36/46.92	8.68/24.70	2.47/9.24	28.18/67.23	70.85/93.71	7.59/29.14	30.80/58.77	6.73/24.00	72.14/96.89	46.13/73.67	1.59/6.67	15.61/59.65	27.91/73.65	15.47/45.21
2	43.39/NA	19.94/72.22	67.84/NA	2.31/10.19	15.04/45.65	4.05/13.30	1.84/7.93	16.29/44.14	63.34/90.43	6.68/28.38	22.73/48.81	5.38/23.52	72.70/95.62	37.31/64.40	1.56/11.97	7.51/29.81	23.40/77.63	7.63/29.85

Table 7: Top-1 and top-5 accuracy of few-shot learning based on self-supervised pre-training on training split. NL means NonLocal network and TSF means TimeSformer.

#Shot	XD-Violence	UCF-Crime	MUVM	WLASL	Jester	UAV Human	CharadesEGO	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
16	42.94/NA	27.78/76.39	57.71/NA	3.24/8.33	5.87/26.69	1.76/6.33	2.53/10.22	15.55/45.00	57.20/85.80	9.56/31.71	34.35/63.80	5.35/21.94	79.21/97.60	53.43/84.45	2.55/8.47	11.20/53.81	25.91/78.48	17.86/57.95
8	38.45/NA	27.78/70.83	44.05/NA	1.85/2.78	4.42/21.45	0.75/5.74	1.63/6.61	11.91/26.06	51.26/82.32	5.92/21.40	26.37/56.08	4.90/19.38	73.27/97.60	39.40/75.00	1.06/6.20	23.90/47.32	13.49/67.09	5.01/43.57
4	32.62/NA	20.83/69.03	40.09/NA	0.93/2.31	4.52/21.40	0.73/3.28	1.13/5.93	3.87/26.85	41.69/74.22	4.10/21.18	17.45/43.84	4.98/21.17	68.46/96.89	24.87/60.72	3.72/11.02	3.40/41.66	19.04/65.57	5.01/38.13
2	39.35/NA	18.06/63.89	42.73/NA	0.93/2.78	3.44/17.89	0.69/3.42	1.10/4.57	6.59/25.23	29.81/60.83	9.01/13.96	9.73/28.38	3.14/14.24	55.16/92.50	13.01/35.85	1.88/14.28	12.22/51.79	17.55/76.37	4.36/13.73
16	58.30/NA	31.94/77.17	79.74/NA	22.22/55.09	25.62/66.22	8.67/26.31	2.49/10.29	41.71/80.95	84.38/94.44	16.54/47.04	34.87/63.55	10.03/36.49	81.19/98.73	57.17/84.64	2.62/10.91	11.86/49.34	51.00/95.33	30.50/74.07
8	56.73/NA	34.03/79.86	68.72/NA	6.94/24.07	8.16/30.74	2.56/9.35	1.60/7.95	28.88/69.83	48.44/80.21	9.10/32.78	27.21/55.24	6.60/23.13	77.09/99.01	47.53/78.24	1.70/8.25	20.20/57.33	41.46/93.12	15.90/60.78
4	34.30/NA	29.86/77.78	64.84/NA	1.85/5.24	4.59/21.74	0.84/3.79	1.00/6.25	11.34/33.72	42.50/75.83	5.16/20.94	20.21/46.18	5.63/20.29	68.60/96.61	34.84/68.02	2.98/9.39	6.32/29.80	20.55/83.02	9.59/37.47
2	38.45/NA	15.97/61.81	58.15/NA	0.93/2.31	4.32/21.33	0.76/3.99	1.24/4.82	7.18/33.00	34.04/67.37	4.10/16.08	14.05/35.95	4.54/17.48	65.91/96.61	24.24/55.84	1.91/8.11	9.77/31.94	18.91/80.37	12.64/55.29
16	56.73/NA	28.47/77.78	93.27/NA	31.48/60.19	29.49/72.09	7.04/22.01	2.71/11.28	40.20/77.93	62.03/89.58	13.35/41.27	37.17/65.20	9.37/32.08	79.35/98.44	64.85/88.90	3.01/10.38	10.79/39.15	46.33/96.33	30.72/71.68
8	54.04/NA	29.86/75.09	51.68/NA	10.65/34.26	10.72/27.28	1.77/6.59	2.09/8.65	33.13/69.61	56.24/86.86	9.41/33.08	29.94/57.74	5.98/24.10	77.37/98.44	57.87/84.31	1.63/8.43	13.29/29.80	41.20/90.12	18.95/63.62
4	49.55/NA	26.39/70.83	46.21/NA	2.78/5.56	5.77/24.25	0.95/3.75	1.65/6.66	15.33/48.00	47.73/78.65	4.25/20.18	22.40/48.87	4.61/18.84	68.88/97.74	42.45/76.02	2.20/8.27	20.80/68.89	27.35/77.42	9.37/37.04
2	43.05/NA	24.31/69.44	45.88/NA	1.89/5.18	4.51/21.64	0.81/3.99	1.25/4.94	8.04/33.92	29.76/57.96	4.59/21.70	12.01/33.18	4.68/15.78	60.96/93.21	24.81/57.17	1.31/8.94	10.31/44.40	33.30/84.09	8.28/33.55
16	55.83/NA	34.72/77.47	95.59/NA	30.56/60.05	28.75/70.20	5.71/19.00	2.56/10.64	40.36/79.02	60.17/87.81	13.05/37.33	35.63/64.54	8.75/28.95	81.00/99.01	65.36/90.16	4.43/11.96	10.85/31.17	42.30/96.36	29.19/70.15
8	52.03/NA	27.78/76.69	51.54/NA	10.19/28.70	11.81/38.63	1.42/3.18	1.94/7.74	29.28/68.14	52.72/84.14	9.10/31.87	29.20/57.91	6.37/23.65	73.82/97.88	57.17/85.96	2.05/9.64	10.61/28.31	35.53/83.66	18.30/55.77
4	49.78/NA	25.69/68.75	47.58/NA	1.85/5.09	6.36/24.93	0.85/3.82	1.30/5.52	10.95/36.78	46.12/77.69	6.37/24.58	21.60/47.35	5.68/19.83	60.96/9					

Table 8: Top-1 and top-5 accuracy of few-shot learning based on self-supervised pre-training on training split2. NL means NonLocal network and TSF means TimeSformer.

#	Shot	XD-Violence	UCF-Ctime	MUVM	WLASL	Jester	UAV Human	CharadesEgo	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
TSN	16	62.00NA	36.1181.25	74.45NA	32.4163.89	32.6670.59	17.6940.94	4.5717.65	33.8376.37	80.1696.88	24.5861.31	45.2874.62	16.1649.44	96.3210.00	76.3395.05	13.3542.30	41.8488.08	40.1294.72	38.3483.44
	8	63.90NA	27.7877.78	71.81NA	25.0045.83	20.7153.50	10.0727.49	3.8014.68	31.0970.99	76.4496.22	18.0653.72	38.5869.06	9.2253.83	95.0510.00	70.4991.75	7.0826.74	26.1675.74	34.5690.23	32.2471.24
	4	50.67NA	24.3175.00	66.52NA	14.3530.56	10.0835.88	4.0212.67	2.7010.63	14.3846.36	70.3993.66	14.8745.98	31.7061.66	7.4128.41	90.3171.00	62.5089.02	6.4122.81	24.4964.36	23.1480.25	18.5252.29
	2	42.83NA	20.8371.53	53.30NA	6.0211.11	6.7428.53	2.1447.66	1.6314.97	2.4899.07	62.1989.02	8.0435.66	25.2665.70	5.8722.59	90.8199.72	52.3563.06	4.8516.47	16.2162.51	25.5865.17	15.0345.53
TSM	16	65.47NA	42.3679.86	99.56NA	47.6978.24	36.8676.49	20.2043.08	4.2616.44	46.9083.23	76.3894.41	21.4050.23	42.9470.00	15.4345.98	93.9299.72	72.0892.20	47.116.12	22.4773.78	61.2497.82	37.4776.69
	8	59.19NA	29.1779.17	68.32NA	26.8558.80	18.2255.45	12.4821.51	2.8211.47	38.1475.06	67.5294.21	19.7340.62	36.0064.52	8.9095.29	94.0699.72	66.9488.90	3.7515.20	22.0569.55	51.6793.79	29.1976.91
	4	54.60NA	25.0070.83	64.32NA	9.7225.93	12.3657.49	6.8819.36	2.319.48	24.1556.99	67.5294.21	8.8024.89	7.7656.22	6.6225.41	88.2699.58	57.1785.74	2.4470.52	17.6166.09	41.6180.53	14.6051.63
	2	50.22NA	22.9271.53	61.23NA	4.1712.50	9.9234.89	2.829.63	2.308.77	13.8043.94	60.6788.52	7.7431.11	21.5848.28	6.1321.71	83.3198.59	45.5676.84	2.9010.17	16.0915.76	40.4663.79	9.5936.49
I3D	16	54.26NA	34.0375.69	65.64NA	48.13580.09	37.5076.91	18.4842.31	3.2712.34	41.8694.24	77.6495.02	32.0266.46	39.4967.27	19.9758.71	94.6371.00	67.0790.36	13.7140.45	55.1390.58	59.0896.74	47.0686.49
	8	55.94NA	35.4270.83	67.40NA	37.5073.15	24.0759.93	11.3670.24	3.3811.54	34.0180.14	73.4194.16	21.2453.41	33.4991.66	13.8444.78	93.9299.72	10.8032.55	14.8085.10	48.3392.97	37.0484.10	
	4	43.16NA	25.6968.06	53.74NA	19.9142.13	14.9343.29	5.4691.68	2.479.47	23.2125.10	69.3992.33	14.4244.92	27.1354.66	8.2590.94	90.9599.86	53.5682.42	7.1224.09	27.7171.45	35.1586.45	27.4665.58
	2	41.26NA	16.6763.89	55.07NA	8.3919.91	8.0031.01	2.629.45	1.977.97	18.6851.39	62.2488.42	11.9937.48	22.6497.07	6.4922.49	86.4298.73	41.3174.05	4.1817.25	24.3162.46	32.4391.12	11.3346.19
NL	16	56.39NA	35.4270.86	77.53NA	51.3989.09	40.8681.78	18.3591.01	3.2212.65	38.8780.66	77.9095.17	27.7764.04	39.1266.53	19.6756.14	95.1910.00	67.1390.29	14.2841.12	55.7293.15	66.4797.15	49.4683.44
	8	50.22NA	25.0072.93	66.96NA	39.3572.69	25.7263.19	3.1810.89	3.0472.21	73.2694.31	73.2694.31	26.1055.39	32.6160.25	13.1744.59	92.9310.00	62.1886.93	11.3233.46	24.1284.74	48.5990.94	31.8174.95
	4	45.74NA	22.9268.75	68.72NA	20.8345.83	15.0444.16	5.3116.76	2.619.63	20.2858.62	70.1991.84	18.5147.34	26.6563.37	8.5940.21	90.5299.86	52.0381.66	8.7825.93	24.9770.44	36.0487.25	24.6266.23
	2	39.01NA	21.5362.50	47.58NA	7.4125.00	8.7434.48	2.859.49	2.227.91	18.7650.29	60.9887.56	8.8043.10	22.5547.66	5.9522.14	87.2799.72	42.5873.03	3.7515.80	20.509.18	33.9488.30	15.2548.58
TSF	16	66.82NA	31.5277.47	92.95NA	42.5975.93	45.2584.82	18.7540.89	4.2714.97	38.6981.37	81.3296.53	35.0565.55	46.8074.35	22.4000.81	95.0299.86	79.1295.49	14.5642.47	53.0493.33	64.7599.41	30.3990.41
	8	60.99NA	35.4279.86	84.14NA	34.2665.74	30.4169.20	12.4920.71	3.7013.72	31.6876.35	78.5096.48	26.7165.82	42.1470.84	15.5749.95	96.8999.72	76.4694.48	10.5632.77	42.6186.23	55.1194.66	35.0881.26
	4	56.17NA	26.3969.44	81.50NA	17.5944.44	16.9349.71	6.7829.24	3.3311.27	24.0661.64	72.9693.61	17.0050.53	34.9564.66	11.8641.10	91.2399.86	70.1191.81	8.1824.48	21.3965.55	38.2589.92	23.9762.96
	2	48.32NA	19.3663.89	59.47NA	6.9423.61	10.3335.71	3.3310.76	2.8170.05	17.5852.24	62.9488.27	14.1136.72	28.3256.57	8.3295.90	97.9899.43	57.8784.64	7.6022.32	19.5557.57	35.8699.76	20.4862.31
Swin	16	51.01NA	29.8679.17	92.95NA	43.0675.93	34.2373.22	13.5844.47	4.1014.87	43.0982.27	81.9296.83	34.9167.37	48.9774.58	18.7855.43	85.7199.15	61.2985.28	14.3543.18	44.8284.51	23.5585.63	36.1784.97
	8	52.58NA	23.6171.53	79.74NA	37.0465.73	23.6858.92	6.8720.88	3.5813.33	36.1582.22	77.0996.02	24.5862.06	41.1768.44	12.4343.24	84.8698.73	48.1678.93	12.6133.62	7.5751.64	30.4582.20	26.3675.08
	4	39.69NA	22.9265.28	52.42NA	24.5445.00	13.8943.32	2.8910.20	2.9011.01	22.0757.30	71.8093.20	10.7950.22	33.0060.27	9.3054.81	76.1097.45	35.9865.80	6.8023.66	19.1353.58	27.0479.45	22.6667.08
	2	33.74NA	22.2265.97	54.53NA	9.2627.31	7.7230.08	1.335.27	2.108.23	17.6748.89	61.7389.73	13.3535.96	26.2465.42	6.3225.22	80.7698.30	27.9253.11	5.5320.33	19.7955.36	16.7568.83	8.0642.70

Table 9: Top-1 and top-5 accuracy of few-shot learning based on self-supervised pre-training on training split2. NL means NonLocal network and TSF means TimeSformer.

#	Shot	XD-Violence	UCF-Ctime	MUVM	WLASL	Jester	UAV Human	CharadesEgo	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
TSN	16	45.74NA	34.7285.42	74.01NA	2.7897.72	9.5693.35	1.6567.79	2.6810.78	17.0477.82	58.0186.46	16.3943.40	33.6363.88	10.1634.98	85.5798.87	54.7083.12	7.4728.27	5.0737.19	25.6074.17	29.4181.92
	8	34.75NA	20.8361.11	46.70NA	0.4657.00	0.835.79	1.988.22	1.988.22	10.0929.23	51.1181.77	11.9940.21	25.5655.93	5.9272.24	80.7699.15	39.9774.56	7.4723.66	11.1443.08	18.0379.60	7.6352.51
	4	36.66NA	9.7255.56	47.58NA	0.4657.00	4.2519.81	0.955.69	1.565.99	7.347.82	42.3576.23	2.7317.60	18.5643.49	5.1317.61	68.6099.01	24.8158.31	4.9217.14	7.6341.42	19.0968.80	6.3240.31
	2	42.60NA	9.7255.56	43.61NA	1.3923.31	4.4020.63	0.653.45	1.054.22	3.4226.63	29.7661.33	3.038.04	11.2529.61	5.5717.57	57.1699.23	15.4238.07	1.4516.65	7.3347.20	17.3477.09	7.4177.23
TSM	16	56.28NA	39.8677.08	82.82NA	20.8351.85	26.1264.21	8.4826.25	2.7111.63	36.7881.13	56.2463.50	30.3368.74	35.2866.59	16.2956.13	86.2329.15	56.0384.58	10.0232.06	34.0386.63	56.0894.74	49.0289.98
	8	50.45NA	29.8677.78	68.72NA	8.8030.56	9.6650.81	2.429.47	2.529.78	32.4174.73	48.6481.87	20.0399.03	27.5255.24	10.0238.60	82.3299.15	46.7077.47	8.0828.48	16.6947.14	43.8293.02	35.0881.05
	4	43.61NA	11.8163.19	40.09NA	0.4643.63	4.8622.09	0.633.34	1.325.06	7.2730.94	35.4067.77	5.4618.51	14.9937.52	5.4119.41	66.3497.17	36.2968.15	5.9925.12	2.9231.35	31.8180.94	13.7352.51
	2	40.34NA	16.6765.97	61.23NA	1.8560.02	6.4920.46	0.843.62	1.566.13	5.7629.23	30.6163.14	10.1729.14	14.9933.98	4.9019.25	64.9294.63	28.3658.06	2.1317.07	8.7038.44	25.0489.12	11.7644.01
I3D	16	50.34NA	36.8178.47	92.51NA	30.5659.72	31.1671.43	6.7522.13	3.3212.55	38.2781.37	63.1989.78	29.7464.92	37.3965.98	14.1947.43	89.2599.29	64.9189.47	13.4340.31	33.3777.59	43.8494.00	43.5788.02
	8	51.35NA	27.7867.39	33.04NA	17.1441.67	12.7038.79	1.286.06	2.6810.01	33.8575.28	56.7086.91	27.6295.52	30.2358.03	8.5433.90	85.1599.58	53.1678.57	9.1430.18	30.6376.16	39.2890.38	36.8286.71
	4	41.26NA	27.7867.39	54.19NA	3.2411.11	5.6424.36	1.007.95	1.837.52	20.1049.23	47.9481.47	15.6349.01	22.0948.30	6.5725.14	78.6498.73	40.4873.73	6.5423.85	13.5955.24	29.1281.48	24.4062.75
	2	42.94NA	16.6765.97	61.23NA	1.8560.02	6.4920.46	0.843.62	1.566.13	5.7629.23	30.6163.14	10.1729.14	14.9933.98	4.9019.25	64.9294.63	28.3658.06	2.1317.07	8.7038.44	25.0489.12	11.7644.01
NL	16	47.98NA	35.4277.08	95.15NA	28.7063.43	31.9873.35	5.1178.07	3.1612.64	37.3379.35	61.5358.09	32.4765.10	36.3965.11	13.1645.15	87.3199.29	65.0489.28	15.1242.59	36.2384.21	39.6487.02	42.4887.88
	8	48.21NA	29.8677.69	77.09NA	12.5036.11	10.657.27	2.649.65	3.1907.182	55.7448.30	44.4178.80	27.1154.61	21.1146.67	5.6374.67	81.4399.29	39.7272.65	7.8124.57	12.5148.81	33.4388.56	30.0783.88
	4	39.8NA	25.0072.92	63NA	3.2483.33	4.7125.97	1.033.71	1.637.12	19.2051.41	44.4178.80	28.0559.02	14.0253.84	6.9832.63	60.8293.49	28.3356.73	2.3810.12	12.5757.09	30.5881.73	20.0454.90
	2	43.61NA	22.2268.75	54.63NA	0.4644.17	4.8220.92	0.753.43	1.365.34	12.3938.06	44.4178.80	28.0559.02	14.0253.84	6.9832.63	60.8293.49	28.3356.73	2.3810.12	12.5757.09	30.5881.73	20.0454.90
TSF	16	54.15NA	35.4277.78	65.64NA	6.0220.83	14.5879.99	4.3514.35	3.3112.59	23.6062										

Table 10: Top-1 and top-5 accuracy of few-shot learning based on supervised pre-training on training split3. NL means NonLocal network and TSF means TimeSformer.

#	Shot	XD-Violence	UCF-Crime	MUVM	WLASL	Jester	UAV Human	CharadesEGO	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
TSN	16	61.66NA	38.1975/50	74.89NA	37.0467/59	16.437/39.77	16.437/39.77	4.081/5.90	34.387/8.04	79.979/6.93	26.566/2.06	45.137/4.74	16.895/0.97	95.901/00	75.769/5.05	14.28/42.47	42.798/7.19	42.899/5.74	35.297/6.91
	8	58.52NA	24.317/2.92	67.84NA	22.225/40.46	20.88/54.01	8.61/24.40	3.921/3.56	30.156/4.90	76.38/95.30	20.33/56.30	38.99/69.51	9.78/36.08	96.75/99.86	70.94/91.81	10.06/31.53	30.81/80.93	30.32/89.89	25.63/66.88
	4	52.13NA	22.926/6.67	67.84NA	5.09/22.22	1.12/17.76	3.74/12.74	3.04/11.74	22.66/2.52	69.13/93.61	16.54/43.85	30.47/61.13	6.41/25.25	93.64/99.72	62.82/88.58	7.93/24.94	29.56/68.71	20.60/82.09	25.53/65.04
	2	38.00NA	15.977/0.83	48.02NA	2.31/12.04	0.68/30.20	2.10/7.18	2.35/8.88	15.41/44.82	63.24/90.18	7.13/34.29	22.48/51.27	5.10/21.87	92.64/99.01	52.79/92.74	3.86/20.97	20.14/67.63	23.37/83.68	23.31/63.81
TSM	16	56.73NA	34.037/5.00	100NA	48.617/5.93	36.507/5.97	19.32/42.72	3.34/13.15	76.64/94.41	33.84/63.73	41.707/10.10	20.005/7.30	94.23/10.00	70.81/92.83	16.19/44.56	52.62/88.44	66.37/97.20	47.93/81.92	
	8	54.82NA	24.317/3.61	99.12NA	21.50/5.82	11.65/29.98	3.00/11.51	31.36/68.30	73.41/93.61	29.59/58.25	35.79/64.64	11.62/40.14	93.94/99.86	63.96/88.07	56.35/83.76	8.54/27.06	33.37/81.59	40.57/92.38	33.33/65.80
	4	51.46NA	16.677/6.67	97.36NA	7.41/31.94	1.12/8.58	5.71/18.55	2.45/9.11	21.76/53.15	68.98/91.99	22.46/52.40	26.96/56.24	7.76/28.79	91.94/99.86	63.51/83.76	3.74/21.75	20.44/69.48	37.11/87.53	37.11/87.53
	2	37.78NA	14.58/60.42	56.39NA	3.24/16.20	9.20/30.85	1.83/7.86	1.83/7.86	18.13/46.40	59.32/87.06	14.87/40.67	20.92/46.43	20.24/57.21	89.67/99.01	46.51/76.52	7.44/21.75	20.44/69.48	37.11/87.53	37.11/87.53
BD	16	58.30NA	28.477/4.31	84.14NA	49.54/80.36	38.747/7.70	18.97/41.71	3.13/11.82	44.23/82.61	77.44/95.77	33.08/69.50	39.14/67.25	20.24/57.21	94.91/99.72	67.70/90.42	17.36/43.64	55.42/91.00	62.21/96.02	44.88/84.10
	8	56.05NA	25.007/1.53	71.37NA	35.65/67.59	25.16/62.01	10.40/29.54	2.90/10.62	34.29/72.37	74.07/95.23	22.91/53.26	33.04/65.03	13.49/43.27	93.94/99.86	59.96/86.10	12.29/33.86	46.07/93.92	34.42/70.15	
	4	52.80NA	13.89/59.72	70.48NA	13.42/40.74	3.00/41.60	5.02/17.33	2.49/9.10	25.36/55.95	69.54/93.25	22.91/54.32	26.88/65.94	7.44/29.68	92.22/99.86	53.17/81.22	7.83/25.61	32.48/73.24	32.56/85.27	26.80/62.09
	2	40.22NA	9.72/57.64	48.46NA	3.24/18.06	8.83/30.15	2.20/9.23	1.83/7.84	17.04/45.44	60.78/80.27	14.26/33.69	20.04/40.08	6.37/23.65	88.54/99.01	42.31/73.92	8.42/20.30	21.28/59.59	30.30/87.94	19.83/44.01
NL	16	55.83NA	23.617/1.53	93.83NA	51.85/81.02	39.14/79.26	17.92/40.41	3.23/11.82	41.27/77.45	33.84/60.55	38.40/65.69	20.19/56.78	95.47/93.27	68.53/99.86	42.32/74.43	10.82/25.26	51.22/63.95	32.27/84.99	16.56/52.44
	8	53.14NA	18.75/64.58	84.58NA	40.74/72.22	25.94/63.34	10.09/28.42	3.04/10.86	32.21/70.13	74.47/94.11	26.10/55.84	32.24/61.09	12.32/41.00	95.19/10.00	58.95/86.23	10.98/33.65	43.68/85.04	46.87/93.97	30.50/66.45
	4	51.01NA	11.11/64.58	87.22NA	12.04/38.43	14.62/44.05	4.74/16.50	2.47/9.72	25.58/57.98	70.59/92.65	22.61/53.26	25.59/53.31	7.54/29.73	92.22/99.86	52.22/81.35	9.92/29.58	39.27/78.19	34.94/87.17	20.32/61.87
	2	40.36NA	9.03/57.64	48.46NA	4.17/18.06	11.09/54.59	2.80/9.43	1.96/7.52	16.12/46.73	61.88/88.77	17.45/39.76	19.86/44.46	6.14/24.95	86.32/99.86	42.32/74.43	10.82/25.26	51.22/63.95	32.27/84.99	16.56/52.44
TSF	16	63.23NA	34.037/7.08	90.75NA	48.617/8.24	42.71/84.61	18.60/40.70	3.68/13.00	38.67/78.74	80.06/96.78	32.63/63.58	47.45/73.92	21.57/61.38	96.32/99.86	80.14/94.99	15.34/41.91	53.10/90.52	63.83/97.51	48.15/83.44
	8	63.23NA	33.337/6.30	72.69NA	33.336/5.28	31.69/70.08	11.97/31.10	3.81/13.26	32.52/72.72	78.35/96.53	25.19/60.24	42.07/71.33	16.24/49.83	96.75/99.86	74.68/93.08	11.26/34.29	44.76/85.94	53.18/95.46	39.65/75.82
	4	57.29NA	19.44/68.06	69.16NA	10.65/40.28	21.26/55.87	6.22/19.31	3.44/12.43	28.53/63.19	72.16/94.61	23.22/53.57	34.64/64.44	11.57/39.00	95.62/99.72	68.27/90.67	8.50/26.74	34.51/75.27	43.15/94.61	28.32/73.64
	2	39.01NA	25.00/65.28	41.85NA	3.70/17.59	14.51/42.98	3.01/11.32	2.75/10.01	23.17/47.03	62.79/88.02	13.81/38.99	27.02/55.71	8.84/29.33	86.32/99.86	58.38/85.58	4.75/21.11	25.21/65.79	36.07/86.38	20.11/69.11
Swin	16	48.99NA	26.397/2.92	70.48NA	50.007/7.78	33.367/3.46	18.16/41.22	3.68/13.00	44.80/84.01	81.42/96.84	23.60/67.98	43.43/71.21	19.43/54.08	86.56/99.01	66.12/88.26	18.53/45.91	53.99/93.03	41.28/89.82	50.11/89.11
	8	44.73NA	22.262/5.00	60.35NA	30.09/61.11	21.74/59.32	9.86/4.42	3.45/11.80	36.57/75.08	76.79/95.57	27.01/60.55	37.84/65.73	11.98/39.78	84.02/99.43	57.49/82.30	9.42/32.73	6.91/45.65	25.37/87.28	38.13/62.14
	4	46.52NA	16.67/60.42	75.33NA	15.28/45.83	12.77/42.69	5.981/9.07	2.70/9.88	25.05/62.56	72.96/94.41	17.75/46.74	31.29/44.74	8.29/29.54	70.82/97.88	44.07/71.26	8.47/25.22	23.30/62.63	22.37/84.48	24.62/65.58
	2	34.87NA	20.14/68.33	62.56NA	4.63/18.52	11.81/54.77	3.41/13.13	2.06/8.14	22.38/51.02	61.83/88.52	13.20/39.00	23.49/49.86	6.59/24.81	74.68/98.30	38.38/64.78	6.52/22.42	20.56/51.19	19.09/84.99	17.435/4.68

Table 11: Top-1 and top-5 accuracy of few-shot learning based on self-supervised pre-training on training split3. NL means NonLocal network and TSF means TimeSformer.

#	Shot	XD-Violence	UCF-Crime	MUVM	WLASL	Jester	UAV Human	CharadesEGO	Toyota SmartHome	Mini HACs	MPII Cooking	Mini SportsIM	FineGym	MOD20	COIN	MECCANO	InHARD	PETRAW	MISAW
TSN	16	47.53NA	31.257/3.61	58.59NA	2.787/8.87	8.72/30.80	1.697/6.68	2.991/0.98	14.47/45.35	56.55/86.05	18.21/42.34	33.00/63.43	9.25/36.02	83.73/99.15	54.76/83.31	9.03/31.95	9.71/41.78	25.73/69.29	25.937/1.24
	8	33.30NA	20.14/59.72	55.07NA	1.395/9.09	4.58/21.26	0.73/3.54	2.16/8.45	10.90/33.36	49.85/81.92	12.44/38.09	25.11/55.81	5.43/24.71	79.49/98.87	39.72/73.54	6.31/23.91	7.45/52.15	4.21/65.73	8.93/46.62
	4	43.72NA	16.67/57.64	67.14NA	1.392/7.84	4.43/19.94	0.87/3.43	1.70/6.03	13.84/33.26	39.93/73.01	4.25/27.47	16.76/43.04	5.13/21.89	68.60/99.01	27.21/60.47	1.20/18.63	25.21/43.50	0.05/70.83	10.24/28.98
	2	43.72NA	12.50/55.56	65.64NA	0.46/3.70	4.41/19.95	0.73/3.43	1.00/4.56	7.95/20.12	26.33/59.87	6.22/20.03	9.45/28.81	3.86/13.57	63.08/94.48	13.52/38.64	9.78/24.76	5.90/26.38	17.09/61.11	3.70/43.14
TSM	16	56.61NA	33.337/9.17	82.82NA	20.834/9.07	24.37/63.27	8.99/27.43	2.98/11.30	35.73/84.63	54.73/85.60	32.02/69.65	34.02/63.45	16.32/52.97	85.71/99.15	47.74/84.71	13.18/41.02	33.31/83.19	52.08/94.10	47.93/91.86
	8	51.35NA	26.397/9.17	84.14NA	8.39/29.17	10.51/33.29	2.70/10.05	2.33/9.23	26.49/64.74	48.49/81.82	18.97/58.73	27.19/55.46	9.33/37.84	83.03/98.02	56.19/77.54	10.77/32.91	15.49/59.00	34.25/88.97	31.59/73.86
	4	41.14NA	19.44/68.06	67.84NA	2.31/5.56	5.12/21.78	0.88/3.99	2.25/8.00	11.19/38.78	40.94/74.47	11.23/38.85	19.20/45.81	6.10/25.59	75.67/97.97	35.34/67.20	7.58/27.95	7.39/26.82	24.63/80.37	14.86/52.51
	2	36.32NA	15.28/63.19	53.74NA	1.394/3.63	4.79/21.51	0.88/3.98	1.25/4.48	5.61/55.52	32.43/65.86	5.61/15.17	13.39/35.17	4.84/20.10	68.32/97.31	23.73/54.44	8.71/26.96	6.73/37.25	27.50/78.55	19.61/59.91
BD	16	54.82NA	29.867/7.78	87.22NA	27.316/2.50	29.67/72.63	16.02/40.66	3.17/11.89	37.53/81.52	62.84/89.84	19.88/54.17	36.69/65.11	14.08/45.46	88.13/99.58	65.42/89.15	15.80/40.84	45.41/88.14	46.23/92.10	48.37/91.94
	8	54.15NA	19.44/68.06	74.01NA	13.897/3.96	12.45/59.93	1.60/6.40	2.50/9.52	28.79/69.46	57.65/86.05	20.79/56.60	29.88/58.69	7.81/31.00	82.04/99.86	54.31/81.33	9.28/30.82	25.74/78.34	24.42/85.99	35.51/83.01
	4	52.13NA	11.537/0.83	65.2NA	2.787/4.41	6.19/24.24	1.03/3.78	2.10/7.90	19.384/45.3	45.17/71.84	17.60/44.12	21.66/48.49	6.06/23.67	74.82/98.16	42.07/72.53	5.21/21.75	13.23/44.28	23.02/83.63	21.56/71.08
	2	40.47NA	16.67/65.28	33.48NA	1.855/5.56	4.56/21.43	0.73/3.71	1.41/5.68	9.87/37.90	27.19/59.47	8.50/31.91	13.12/33.94	5.29/20.59	63.79/95.47	25.38/59.45	5.14/22.92	2.92/27.41	25.94/83.35	14.81/59.43
NL	16	55.04NA	30.568/1.94	92.07NA	32.876/5.28	32.437/4.07	5.80/19.50	3.32/12.30	32.47/78.63	39.29/65.40	29.29/65.40	36.20/64.66	12.36/42.67	86.14/99.58	66.12/88.83	11.33/38.67	39.39/84.62	40.69/93.38	48.37/88.67
	8	48.54NA	25.697/2.92	81.5NA	9.26/35.19	14.17/41.59	1.00/4.75	2.44/9.39	23.39/60.03	54.78/84.74	20.94/58.27	30.71/66.14	7.17/31.38	78.22/98.30	54.31/83.38	11.26/32.66	10.67/61.62	28.99/84.60	30.72/72.98
	4	47.98NA	21.536/8.75	84.15NA	3.70/6.48	6.57/25.46	0.89/3.62	1.97/7.52											



## 5. Unsupervised domain adaptation

In BEAR, we construct two different types of transfer for unsupervised domain adaptation (UDA): inter-dataset transfer and intra-dataset transfer. We construct paired source-target with different datasets for inter-dataset, while we build paired data within one dataset according to similar or same actions for intra-dataset. The main challenge in our UDA benchmark could be caused by viewpoint change (e.g., ToyotaSmarthome-MPII-Cooking), long-tail problem (PHAV-Mini-Sports1M), etc.

### 5.1. Inter-dataset

**ToyotaSmarthome-MPIICooking** One of the features of our benchmark is that we collect several datasets with obvious viewpoint shifts, and we also leverage this when we build our UDA datasets. As shown in Fig. 1, Toyota Smarthome and MPII-Cooking consists of videos from different viewpoints. Specifically, as shown in Table 12, we select 6 common categories in Toyota Smarthome and MPII-Cooking to construct the new Toyota Smarthome-MPII-Cooking dataset. As shown in Fig. 3, the video numbers can be imbalanced across source data and target one, for action class '*eat(drink)*', there are a total of 3317 videos in Toyota Smarthome, since 7 original classes are merged; while there are only 21 samples from the original class '*taste*' in MPII-Cooking. The number of videos is 5,233 and 943 for Toyota Smarthome and MPII-Cooking, respectively.

**Mini-Sports1M-MOD20** Similarly, as shown in Table 13, for Mini-Sports1M and MOD20, we select 15 categories to build the UDA dataset. The statistic is shown in Fig. 4. In contrast to Toyota Smarthome-MPII-Cooking, the data distribution in Mini-Sports1M-MOD20 is much more balanced. There are 1,650 videos for Mini-Sports1M and 1,767 for MOD20.

**UCF-Crime-XD-Violence** Similarly, UCF-Crime and XD-Violence share three classes: *abuse*, *fighting*, and *shooting*. As shown in Fig. 5, sample numbers of *fighting* and *shooting* showcase an obvious imbalance distribution, which makes the UDA task here much more challenging.

**PHAV-Mini-Sports1M** We also consider the synthetic-to-real transfer and we leverage an existing dataset PHAV [18]. As shown in Table 14, we combine 15 classes from Mini-Sports1M into 6 categories (*playing soccer*, *playing golf*, *playing baseball*, *shooting gun*, *shooting archery* and *running*) existing in PHAV to build the paired dataset. We also illustrate the class-wise distribution of this dataset in Fig. 6. PHAV contains much more samples than Mini-Sports1M due to it is easily generated.

### 5.2. Intra-dataset

**Jester(S-T)** We also include existing Jester(S-T)[20] in BEAR. The category information is shown in Table 15 that each identical action with a contrary direction is merged into one category. For completeness, we also include its class-wise distribution in Fig. 7.

**InHARD(Left-Top-Right)** InHARD naturally contains three different views and each frame contains the top, left, and right views, respectively as shown in Fig. 2. We simply split the frames and the category is the same as the original dataset, and the samples in each category are also the same in Fig. 8.

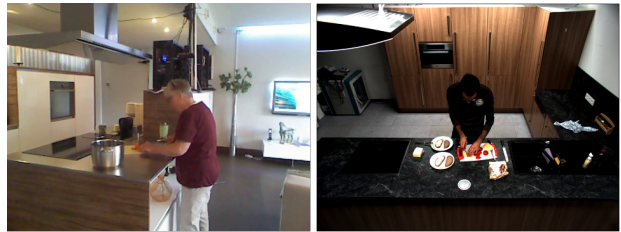


Figure 1: Example frames of Toyota Smarthome and MPII-Cooking. The left frame is from Toyota Smarthome, in which the videos are captured from 7 different cameras; the right one is from MPII-Cooking and is recorded by a fixed down view camera.



Figure 2: Examples of InHARD(Left-Top-Right). The shown frames are from the left, top, and right views, respectively. The left view can be severely occluded, making it much more challenging to transfer knowledge from this view to others.



Table 12: Action classes in Toyota Smarthome-MPII-Cooking.

Toyota Smarthome-MPII-Cooking	Toyota Smarthome	MPII-Cooking
stir	Cook.Stir	stir
wash objects	Cook.Cleandishes	wash objects
cut	Cook.Cut Cutbread	cut out inside cut apart cut in cut dice cut slices cut off ends cut stripes
eat(drink)	Eat.Snack Eat.Attable Drink.Fromcan Drink.Frombottle Drink.Fromcup cut off ends Drink.Fromglass	taste
pour	Pour.Fromkettle Pour.Fromcan Pour.Frombottle Makecoffee.Pourwater Makecoffee.Pourgrains	pour
cleaning up	Cook.Cleanup	wipe clean

Table 13: Action classes in Mini-Sports1M-MOD20.

Mini-Sports1M-MOD20	Mini-Sports1M	MOD20
backpacking	backpacking(wilderness) hiking	backpacking
diving	diving free-diving cuba diving	clif_jumping
cycling	cycling	cycling
boxing	boxing shoot boxing kick boxing	motorbiking
figure skating	figure skating	figure skating
jetsprint	jetsprint	jetskii
kayaking	kayaking	kayaking
motor biking	motorcycle racing grand prix motorcycle racing motorcycle speedway motorcycle drag racing	motorbiking
football	American football Canadian football	nfl_catches
rock climbing	rock_climbing	rock_climbing
running	free running running sprint (running) cross country running	running
skateboarding	freeboard (skateboard) skateboarding	skateboarding
skiing	skiing alpine skiing cross-country skiing freestyle skiing nordic skiing telemark skiing	skiing
surfing	surfing	surfing
windsurfing	windsurfing	windsurfing

Table 14: Action classes in PHAV-Mini-Sports1M.

PHAV-Mini-Sports1M	PHAV	Mini-Sports1M
playing soccer	kick ball	indoor soccer beach soccer
playing golf	golf	golf
playing baseball	swing baseball	baseball
shooting gun	shoot gun	shooting sports practical shooting cowboy action shooting clay pigeon shooting skeet shooting trap shooting
shooting archery	shoot bow	archery
running	run	free running running sprint (running) cross country running

Table 15: Action classes in Jester(S-T).

Jester	Jester Source	Jester Target
Push and Pull	Pushing Hand Away	Pulling Hand Away
	Pushing Two Fingers Away	Pulling Two Fingers Away
Rolling Hand	Rolling Hand Forward	Rolling Hand Backward
Sliding Two Fingers	Sliding Two Fingers Left	Sliding Two Fingers Right
	Sliding Two Fingers Up	Sliding Two Fingers Down
Swiping	Swiping Left	Swiping Right
	Swiping Up	Swiping Down
Thumbs Up and Down	Thumbs Up	Thumbs Down
Zooming In and Out	Zooming Out with Full Hand	Zooming In with Full Hand
	Zooming Out with Two Fingers	Zooming In with Two Fingers
Turning Hand	Turning Hand Counterclockwise	Turning Hand Clockwise

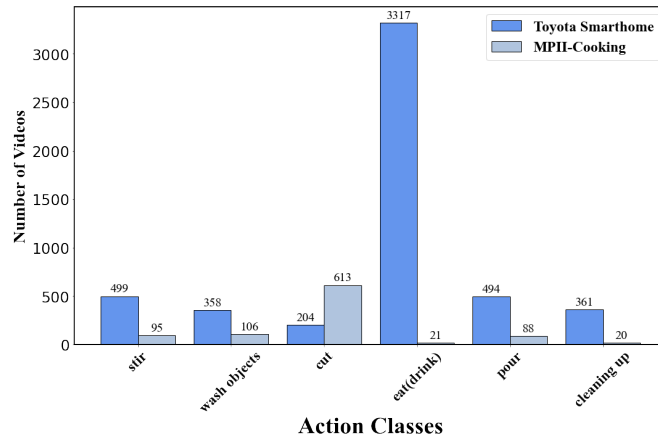


Figure 3: Class-wise distribution of videos in Toyota Smarthome-MPII-Cooking. There is a severe long-tail distribution in Toyota Smarthome.

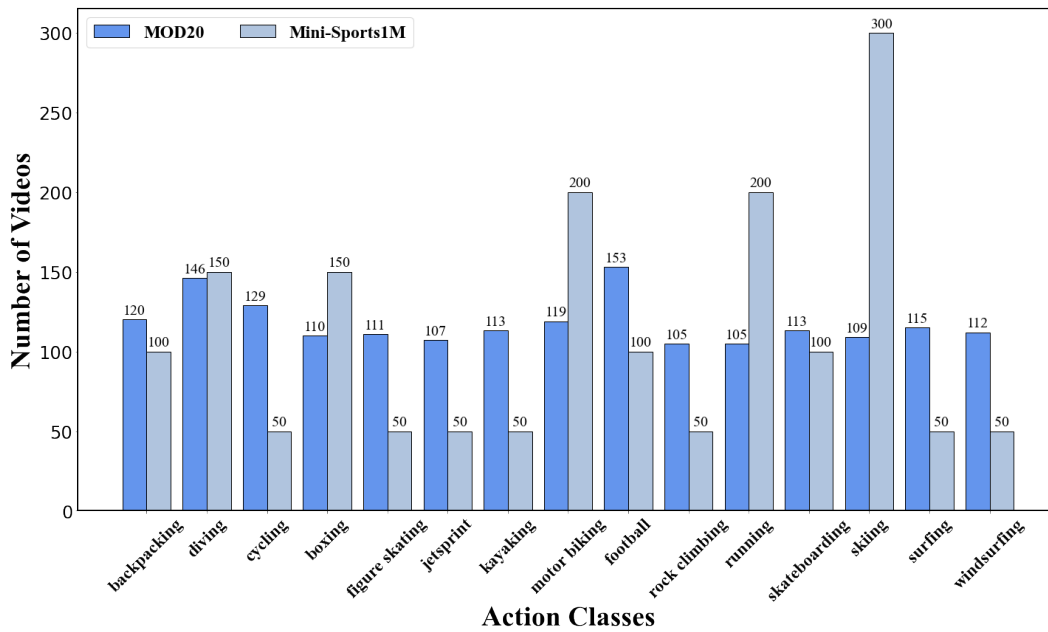


Figure 4: Class-wise distribution of videos in Mini-Sports1M-MOD20. Video numbers in this dataset are much more balanced for source and target.

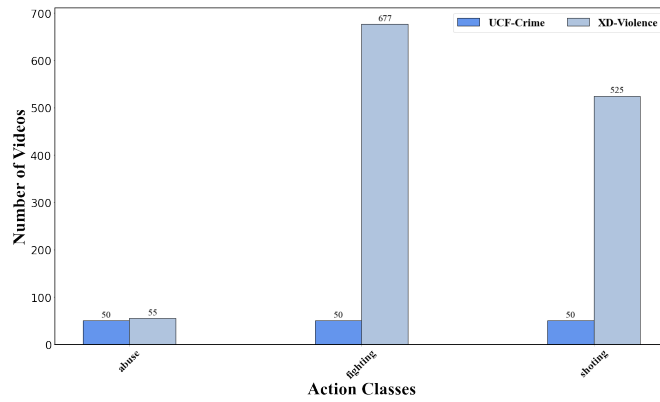


Figure 5: Class-wise distribution of videos in UCF-Crime-XD-Violence.

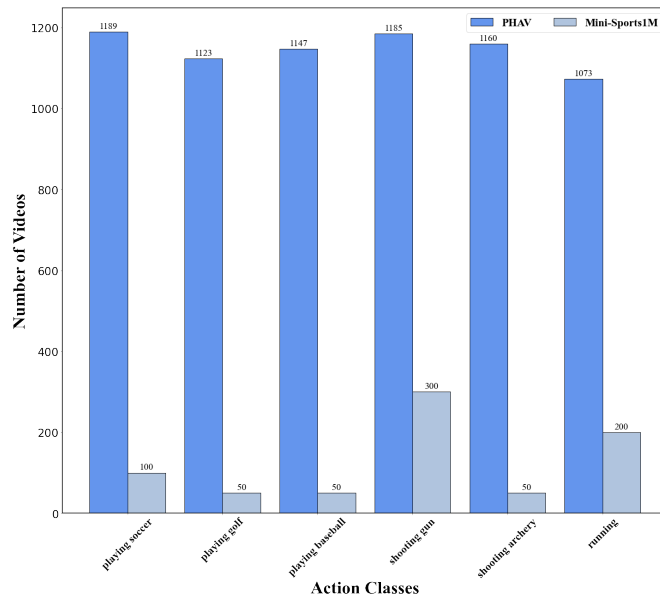


Figure 6: Class-wise distribution of videos in PHAV-Mini-Sports1M.

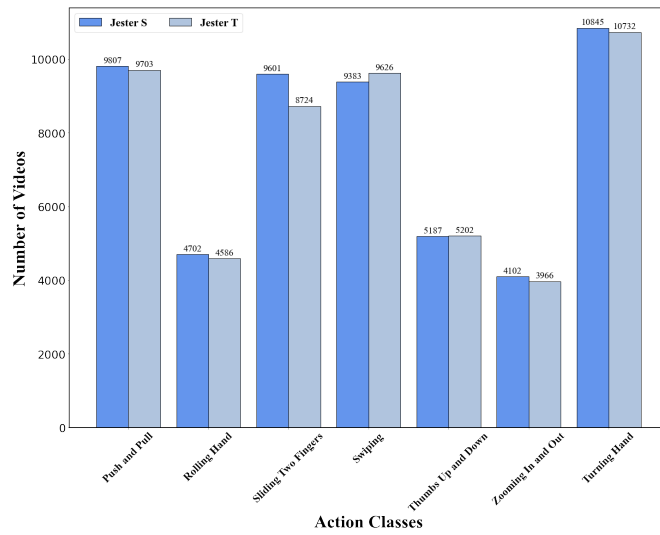


Figure 7: Class-wise distribution of videos in Jester(S-T).

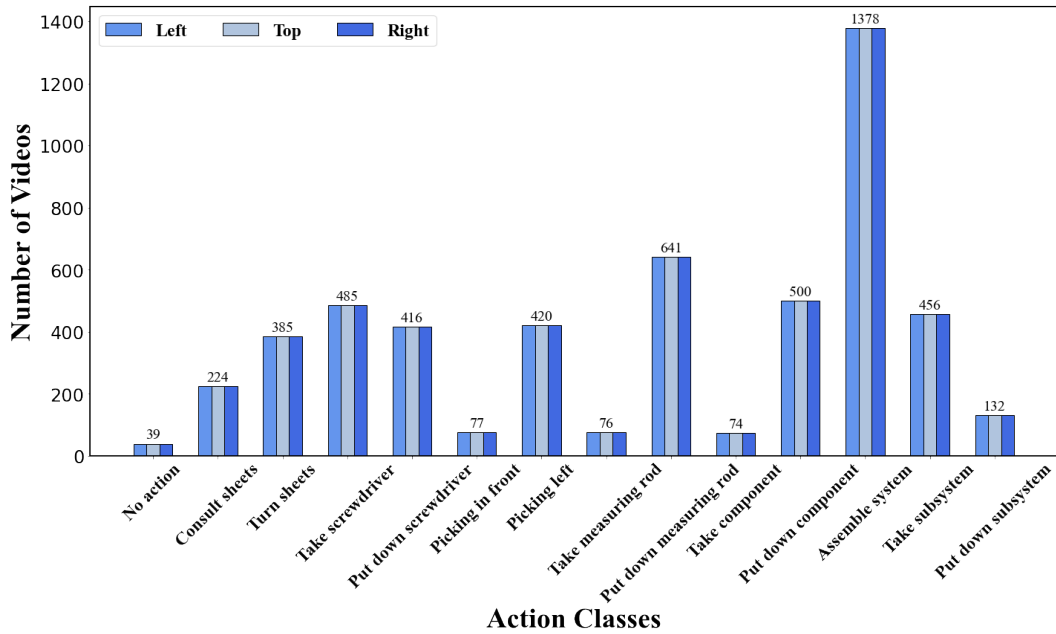


Figure 8: Class-wise distribution of videos in InHARD(Left-Top-Right).

## References

- [1] Ke Cheng, Yifan Zhang, Xiangyu He, Weihan Chen, Jian Cheng, and Hanqing Lu. Skeleton-based action recognition with shift graph convolutional network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 183–192, 2020. 4
- [2] Mejdi Dallel, Vincent Havard, David Baudry, and Xavier Savatier. Inhard-industrial human action recognition dataset in the context of industrial collaborative robotics. In *2020 IEEE International Conference on Human-Machine Systems (ICHMS)*, pages 1–6. IEEE, 2020. 2
- [3] Srijan Das, Rui Dai, Michal Koperski, Luca Minciullo, Lorenzo Garattoni, Francois Bremond, and Gianpiero Francesca. Toyota smarthome: Real-world activities of daily living. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 833–842, 2019. 1
- [4] Srijan Das, Rui Dai, Di Yang, and Francois Bremond. Vpn++: Rethinking video-pose embeddings for understanding activities of daily living. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):9703–9717, 2021. 4
- [5] Stefan Denkovski, Shehroz S Khan, Brandon Malamis, Sae Young Moon, Bing Ye, and Alex Mihailidis. Multi visual modality fall detection dataset. *IEEE Access*, 2022. 2
- [6] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021. 3
- [7] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [8] Hezhen Hu, Weichao Zhao, Wengang Zhou, Yuechen Wang, and Houqiang Li. Signbert: pre-training of hand-model-aware representation for sign language recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11087–11096, 2021. 4
- [9] Arnaud Hualmé, Kanako Harada, Quang-Minh Nguyen, Bogyu Park, Seungbum Hong, Min-Kook Choi, Michael Peven, Yunshuang Li, Yonghao Long, Qi Dou, et al. Peg transfer workflow recognition challenge report: Does multi-modal data improve recognition? *arXiv preprint arXiv:2202.05821*, 2022. 2, 4
- [10] Arnaud Hualmé, Duygu Sarikaya, Kévin Le Mut, Fabien Despinoy, Yonghao Long, Qi Dou, Chin-Boon Chng, Wenjun Lin, Satoshi Kondo, Laura Bravo-Sánchez, et al. Micro-surgical anastomose workflow recognition challenge report. *Computer Methods and Programs in Biomedicine*, 212:106452, 2021. 2, 4
- [11] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. 1
- [12] Dongxu Li, Cristian Rodriguez, Xin Yu, and Hongdong Li. Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1459–1469, 2020. 2
- [13] Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uniformerv2: Spatiotemporal learning by arming image vits with video uniformer. *arXiv preprint arXiv:2211.09552*, 2022. 4
- [14] Tianjiao Li, Jun Liu, Wei Zhang, Yun Ni, Wenqian Wang, and Zhiheng Li. Uav-human: A large benchmark for human behavior understanding with unmanned aerial vehicles. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16266–16275, 2021. 2
- [15] Joanna Materzynska, Guillaume Berger, Ingo Bax, and Roland Memisevic. The jester dataset: A large-scale video dataset of human gestures. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2
- [16] Asanka G Perera, Yee Wei Law, Titilayo T Ogunwa, and Javaan Chahl. A multiviewpoint outdoor dataset for human action recognition. *IEEE Transactions on Human-Machine Systems*, 50(5):405–413, 2020. 1, 4
- [17] Francesco Ragusa, Antonino Furnari, Salvatore Livatino, and Giovanni Maria Farinella. The meccano dataset: Understanding human-object interactions from egocentric videos in an industrial-like domain. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1569–1578, 2021. 2, 4
- [18] Cesar Roberto de Souza, Adrien Gaidon, Yohann Cabon, and Antonio Manuel Lopez. Procedural generation of videos to train deep action recognition networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4757–4767, 2017. 8
- [19] Marcus Rohrbach, Sikandar Amin, Mykhaylo Andriluka, and Bernt Schiele. A database for fine grained activity detection of cooking activities. In *2012 IEEE conference on computer vision and pattern recognition*, pages 1194–1201. IEEE, 2012. 1
- [20] Aadarsh Sahoo, Rutav Shah, Rameswar Panda, Kate Saenko, and Abir Das. Contrast and mix: Temporal contrastive video domain adaptation with background mixing. *Advances in Neural Information Processing Systems*, 34:23386–23400, 2021. 8
- [21] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. Finegym: A hierarchical video dataset for fine-grained action understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2616–2625, 2020. 1, 4
- [22] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Charades-ego: A large-scale dataset of paired third and first person videos. *arXiv preprint arXiv:1804.09626*, 2018. 1
- [23] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488, 2018. 2, 4



- [24] Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019. [2](#), [4](#)
- [25] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5552–5561, 2019. [4](#)
- [26] Thanh-Dat Truong, Quoc-Huy Bui, Chi Nhan Duong, Han-Seok Seo, Son Lam Phung, Xin Li, and Khoa Luu. Direcformer: A directed attention in transformer approach to robust action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20030–20040, 2022. [4](#)
- [27] Peng Wu, Jing Liu, Yujia Shi, Yujia Sun, Fangtao Shao, Zhaoyang Wu, and Zhiwei Yang. Not only look, but also listen: Learning multimodal violence detection under weak supervision. In *European conference on computer vision*, pages 322–339. Springer, 2020. [2](#)
- [28] Hang Zhao, Antonio Torralba, Lorenzo Torresani, and Zhicheng Yan. Hacs: Human action clips and segments dataset for recognition and temporal localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8668–8678, 2019. [1](#)