

000
001
002003

HAA300: A New Benchmark for Human-Centric Atomic Action Recognition

004
005
006
007
008
009
010
011

012 Anonymous CVPR submission

013
014

015 Paper ID

016

017

Abstract

018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034

We propose HAA300: human-centric atomic action video dataset for fine-grained action recognition. Unlike existing and emerging atomic action datasets [15, 5, 13, 34] emphasizing data size and videos captured “in the wild,” HAA300 advocates to human-centric actions with an average of 77% detectable joints for the relevant human poses. Our dataset contains 300 human action classes with 6K action labels in total. Each clip (taken from YouTube) lasts no more than 3 seconds. The collected human actions cover a broad range of classes including human-only actions such as backflip, human-object interactions such as playing instruments, as well as human-human interactions such as shaking hands. We describe the dataset statistics, collection methodology, and compare quantitative performance for baseline neural network architectures trained and tested on human action classification task using this dataset. We demonstrate fine-grained action recognition by evaluating the AVA dataset with models pretrained with HAA300, and transfer learning results on UCF101 and HMDB51 after fine-tuning our pretrained model.

035
036
037

1. Introduction

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Large-scale image datasets such as ImageNet [7] and COCO [25] have greatly contributed to the recent breakthrough in image understanding especially detection and recognition. However, we have not witnessed the same level of impact on video understanding contributed by large-scale datasets in particular action recognition. Observe the coarse annotation provided by commonly-used action recognition datasets such as [28, 21, 50], where the same action label was assigned to a given complex video action sequence (e.g. “Play Football,” “Play Rugby”) typically lasting 10 seconds and thus 300 frames, thus introducing a lot of ambiguities during training as two or more action categories may contain the same *atomic action* (e.g., “Run” is one of the atomic actions for both “Play Football” and “Play Rugby”). In this paper we are interested in human-centric atomic actions (HAAs):

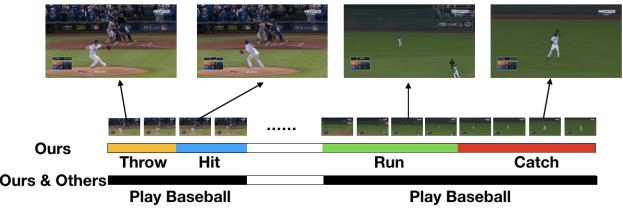
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Figure 1. Recognizing “Play baseball” with its atomic actions using our approach vs others.

Human-centric Atomic Action (HAA) Unlike [15], our human-centric atomic action (HAA) videos were curated with an average of 77% detectable joints and 18% human coverage, while similar to [15] they have clear visual signatures, typically independent of the interacted objects. Such atomic actions usually last no longer than 3 seconds, which are generally enough for capturing a fine-scale temporal movement.

To address the issue of fine-grained action recognition, a number of emerging *atomic action* datasets [15, 5, 13, 34] have been recently introduced. Google’s AVA actions dataset [15] provides dense annotations of 80 atomic visual actions in 430 fifteen-minute video clips, where actions are localized in space and time. Later, AVA spoken activity dataset [35] contains temporally labeled face tracks in videos, where each face instance is labeled as speaking or not, and whether the speech is audible. The something-something dataset [13] contains clips of humans performing pre-defined basic actions with everyday objects.

While the general direction is commendable in decomposing complex actions into a sequence of atomic actions, observed that the videos collected in the first versions of these datasets are almost “in the wild” and the reported performance was very low [15]. Current state-of-the-art performance on AVA is only 34% mAP. We believe AVA’s limited advance in video recognition with poor model performance is due to its action categorization in the presence of a large amount of noise, e.g., class “open” contains both actions such as “open a door” and “open a carpet” with totally different visuals, which poses a deep neural network an almost insurmountable challenge to properly learn the pertinent atomic action.

In this paper, we propose to collect curated action videos where the relevant human poses are detectable (2D pose skeletons as structure), and carefully construct our Human-centric Atomic Action datasets. HAA300 dataset provides fine-grained human-centric atomic visual action clips with clean temporal annotations. An example of the collected atomic actions is shown in Figure 1. The clips are class-balanced and contain clear visual signals with little occlusion. This clean, class-balanced dataset provides an effective alternative to advance research in atomic visual actions recognition and thus video understanding.

2. Related Works

2.1. Representative Action Recognition Dataset

Table 1 compares existing action recognition datasets. Representative action recognition datasets, such as KTH [36], HMDB51 [23], UCF101 [39], Weizmann [1], Hollywood-2 [26] and Kinetics [21], consist of short clips, which are manually trimmed to capture a single action. These datasets are ideally suited for training fully-supervised, whole-clip, forced-choice video classifiers. A few datasets used in action recognition research, such as CMU [22], MSR Actions [49], UCF Sports [33] and JHMDB [17], provide spatio-temporal annotations in each frame for short videos, but they only contain few actions. Aside from the subcategory of shortening the video length, recent extensions such as UCF101 [39], DALY [47] and Hollywood2Tubes [27] evaluate spatio-temporal localization in untrimmed videos, resulting in a performance drop due to the more difficult nature of the task. One common issue on these aforementioned datasets is their actions are restricted to a limited number of composite actions (e.g., tennis), thus are unsuitable for fine-grained action analysis and recognition. Google AVA [15] is the first large-scale atomic action dataset which provides 80 atomic action labels collected from movie clips each lasts for no longer than 3 seconds. However, for the action detection task in the AVA dataset, the state-of-the-art model only achieves mAP of 34% [11]. Other specialized datasets such as Moment in Times [28], HACS [50], Something-Something [13], and Charades-Ego [37] provide classes for atomic actions but none of them are HAA. Our dataset differs from all the aforementioned datasets as we provide over 300 fine-grained atomic human action classes with videos in higher quality (higher resolution), which is carefully constructed to have a better balanced class distribution, a finer granularity in video class selection, and more comprehensive action labeling.

2.2. Action Recognition Architectures

Current action recognition architectures can be categorized into two major approaches, using either 2D image-

Dataset	Videos	Actions	Atomic	
KTH	600	6		162
UCF Sports	182	10		163
HMDB51	6,766	51		164
UCF101	13,320	101		165
Weizmann	81	10		166
Hollywood-2	1,707	12		167
Daily	36,000	10		168
YouTube-8M	8,264,650	4,800		169
Kinetics-400	306,245	400		170
HACS	1,550,000	200	✓	171
Moment in Times	1,000,000	339	✓	172
AVA	57,600	80	✓	173
HAA300	6,000	300	✓	174

Table 1. Summary of representative action recognition datasets. HAA300 is human-labeled and human-centric (18% human coverage and 77% detectable joints) unlike other atomic action datasets.

based or 3D video-based kernels for their convolutional and layer operators. [9, 16]. While 3D video-based methods achieve state-of-the-art performance on most of the action recognition tasks, the model size and total amount of parameters are considerably larger than 2D image-based methods.

2D Image-Based Approaches Temporally recurrent layers such as LSTMs and feature aggregation over time are two popular techniques for 2D image-based methods for propagating temporal information across frames. Without compromising the temporal structure and exploiting image classification networks [20, 24, 29, 30], ConvNets with LSTMs are considered as top performers among current models in this category. However, LSTMs suffer from unsatisfactory performance on feature extraction from low-level motion and expensive training cost, as backpropagation requires unrolling the network across multiple frames through time. Given this bottleneck, researchers have turned to the 3D video-based approaches for more complex action analysis tasks.

3D Video-Based Approaches By incorporating spatio-temporal filters to standard convolutional networks, 3D ConvNets are regarded as a natural approach to video modeling. However, a considerably large number of parameters needs to be optimized [18, 43, 44, 46, 8, 19, 40]. The more practical 3D video-based approach [38] is a two-stream network with both RGB frames and pre-computed optical flow as input to a ConvNet pretrained using ImageNet. Later, the two-stream inflated 3D ConvNets (I3D) [3] was proposed which combines 3D ConvNets and two-stream networks. This approach achieves a very competitive performance on existing benchmarks, and thus most of the state-of-the-art methods for action recognition have adopted this architecture [19]. While one line of research effort has been made to further improve the precision for the 3D video-

216 based approach [31, 6], the other focuses on finding the
 217 trade-off between precision and speed for this architecture,
 218 involving more effectively exploiting the spatial-channels
 219 and temporal-channels correlation information throughout
 220 network layers [45, 48, 51, 52, 42, 32].
 221

222 2.3. Atomic Action Recognition

223 While there exist many action recognition datasets, most
 224 of them consist of composite actions (e.g. long-jump, high-
 225 jump, play-tennis, play-football) thus missing the possibility
 226 of utilizing the inherently hierarchical nature of a given
 227 action. In [10] the authors highlighted the hierarchical
 228 nature of human activities, where at the finest level, the
 229 actions consist of atomic body movements or object manip-
 230 ulation, while at the coarser level the most natural descrip-
 231 tions are in terms of intentionally and goal-directed behav-
 232 iors. To model finer-level events, the AVA dataset [15]
 233 was introduced to provide person-centric spatio-temporal
 234 annotations on atomic actions, which includes 1.58M an-
 235 notations in 80 atomic visual action categories. Although
 236 Google AVA is a large-scale dataset, in practice the dataset
 237 is too noisy for training purpose which is evidenced by
 238 the large number of submissions to the AVA challenge re-
 239 quiring pre-training on ImageNet or Kinetics-400 [12] to
 240 improve their model performance, while many entries still
 241 suffer from class imbalance and non-apparent actions due
 242 to low-resolution videos in the AVA dataset. On the other
 243 hand HAA300 provides a clean and class-balanced dataset
 244 on fine-grained actions [4, 14].
 245

246 3. HAA300

247 3.1. Data Collection

248 The annotation of HAA300 consists of two stages: vo-
 249 cabulary collection and video clip selection. Unlike Google
 250 AVA’s bottom-up approach which generates action classes
 251 based on the selected videos, we aim at building a clean
 252 and fine-grained dataset for atomic action recognition, thus
 253 the videos are collected based on pre-defined atomic actions
 254 following a top-down approach.
 255

256 **Vocabulary Collection** To make the dataset as clean as
 257 possible and useful for recognizing fine-grained atomic ac-
 258 tions, we narrowed down the scope of our super classes into
 259 sports, instrument-playing, and daily routines. Next, fol-
 260 lowing the definition of atomic actions, our dataset has a
 261 hierarchical structure: a composite action as root and dif-
 262 ferent atomic actions decomposed from its composite cat-
 263 egory as leaves. For example, *playing tennis* is a compo-
 264 site action containing leaves such as *serving*, *swinging fore-*
 265 *hand*, *swinging backhand*, *volleying forehand* and *volleying*
 266 *backhand*. Such design poses a challenge on action recog-
 267 nition as we introduce both generality and diversity among
 268 269

270 the actions. Our dataset has atomic actions that vary within
 271 its root action while overlapping across different general ac-
 272 tions (e.g., *forehand* and *serve* are significantly different un-
 273 der their root action of playing tennis, but *forehand* itself ap-
 274 pears similarly across general actions of playing tennis and
 275 playing table tennis). We chose our classes to be human-
 276 centric and pose-apparent because these choices not only
 277 result in relatively less occlusion but also enable an alter-
 278 native approach to recognizing actions with human-pose-
 279 based information. Consequently, we ended up with 300
 280 atomic action classes, where 178 are sports, 43 are instru-
 281 ment playing and the rest are daily routines. The appendix
 282 in Figure 9 shows some samples classes in HAA300.
 283

284 **Video Clip Selection** To ensure our dataset is clean
 285 and class-balanced, all the video clips are collected from
 286 YouTube with the majority having a resolution of at least
 287 720p. Each class of atomic action contains at least 16
 288 training clips. We manually selected the clips with appar-
 289 ent human-centric actions for each class from the top-50
 290 YouTube videos returned by searching using the class name.
 291 Clips are properly trimmed to cover the desired actions. In
 292 general, each class contains clips from at least 3 different
 293 sources so that the action recognition task cannot be triv-
 294 ially reduced to identifying the corresponding background.
 295

296 **Training/Validation/Test Sets** They are split at the class
 297 level, and since the clips in different classes are mutually
 298 exclusive, all clips appear only in one split. The 6000 clips
 299 are split roughly as 16:1:3, resulting in segments of 4800
 300 training, 300 validation, and 900 test clips.
 301

302 3.2. Dataset Analysis

303 HAA300 includes 300 atomic action classes where each
 304 class contains around 20 clips. Its human-centric and fine-
 305 grained nature are the two key features differentiating HAA
 306 from other existing datasets.
 307

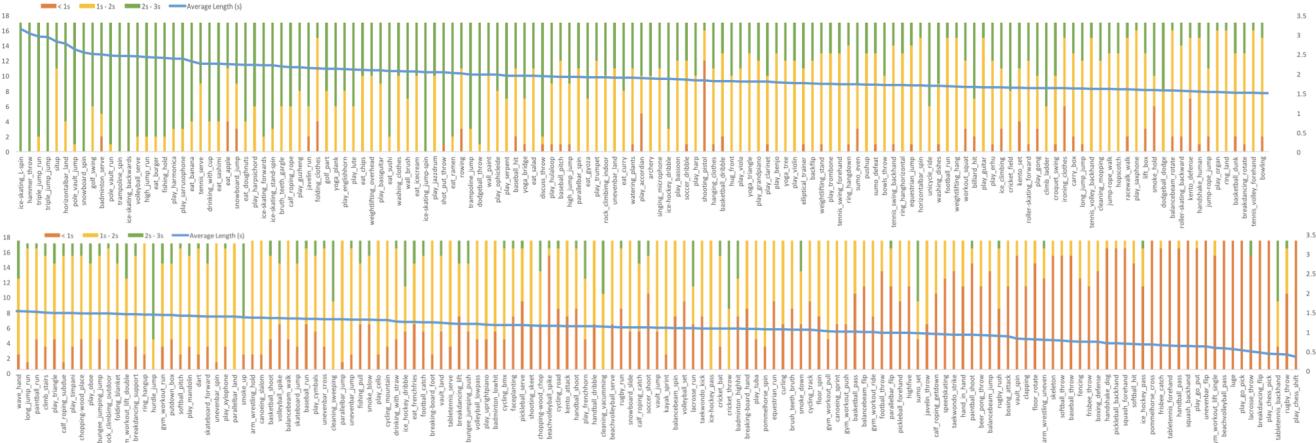
308 3.2.1 Statistics of HAA300

309 Table 2 summarizes the HAA300 statistics. The dataset
 310 contains around 6000 clips with an average length of 1.49
 311 second. Figure 2 shows the distribution of clip lengths
 312 (< 1s, 1 – 2s, or 2 – 3s) and the average clip length
 313 (blue line) for each class. In addition to action category la-
 314 bels, each clip was annotated with several meta-information
 315 tags to allow a more precise and thorough evaluation. The
 316 meta-information contains the following fields: the number
 317 of persons (single or multiple), the presence of occlusion
 318 (human-verified), and the consistency of camera angle. Fig-
 319 ure 3 shows the distribution of the meta tags for the dataset.
 320

action	clips	mean clip length	total duration
300	6000	1.46s	8760s

321 Table 2. Summary of HAA300

324



325



326



327



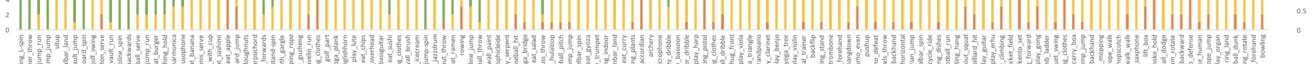
328



329



330



331



332



333



334



335



336



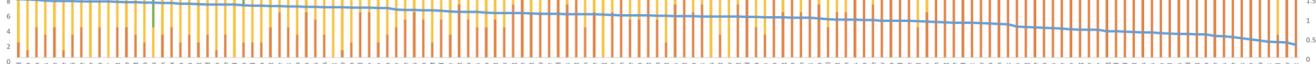
337



338



339



340



341



342



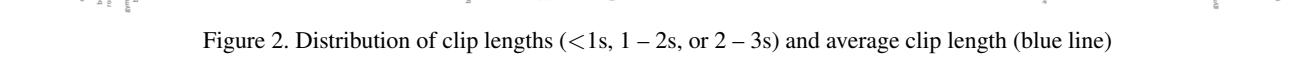
343



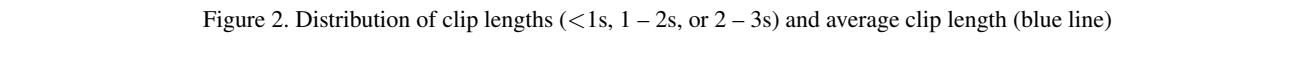
344



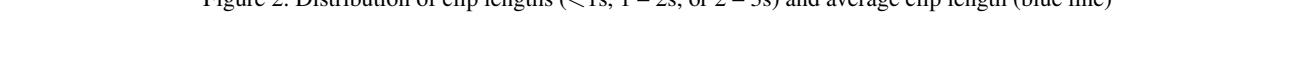
345



346



347



348

Figure 2. Distribution of clip lengths (<1s, 1 – 2s, or 2 – 3s) and average clip length (blue line)

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

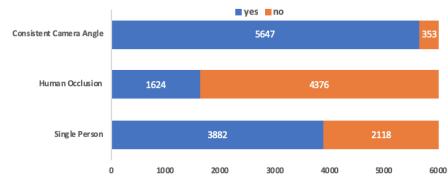


Figure 3. Meta-information for the whole dataset .

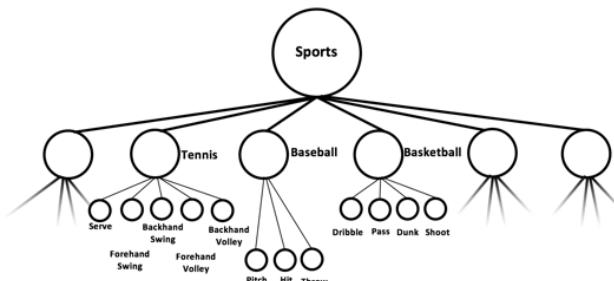


Figure 4. Tree-diagram of classes: sports branch as an example. The full hierarchical structure is shown in Figure 10 in the Appendix.

3.2.2 Characteristics of HAA300

Fine-grained Each class under the HAA300 hierarchy is mutually exclusive, with only the leaf nodes counted as the 300 classes. In general, the first level of the hierarchy is categorized into “sport”, “play instrument”, and “daily routine”. The second level consists of, for example, “tennis”, “badminton” and “kendo” categories, where each sport type will then be further broken into the corresponding atomic actions. For example, tennis is composed of *serving*, *swing forehand*, *swing backhand*, *volley forehand* and *volley backhand* in the third level. Figure 4 illustrates part of the class hierarchical structure.

Human-centric Every clip in HAA300 is human-centric with detectable human poses and thus conducive to fine-grained action recognition. The average of the human coverage in each clip across the whole training set is 18%, where single-person clips tend to have a higher coverage. We also analyze the percentage of detectable joints (each skeleton has a maximum 18 joints) in each clip in the training set and compute the average number of detectable joints for different classes with a keypoint detection model [2]. The number of detectable joints and human coverage do not correlate, while generally sport classes have a larger number of detectable joints. Figure 5 shows the distribution of average human coverage and number of detectable joints, with an average of 77% detectable joints in our dataset.

4. Experiments

4.1. Baseline Benchmarks

We have evaluated the baseline performance on our HAA300 with three different architectures. We adopted 2D-based CNNs and 3D-based CNNs including the I3D network and the SlowFast Network. Our experiments were implemented with PyTorch.

4.1.1 Implementation Details

We adopted the pre-trained ResNet-101 [16] architecture for the 2D models. Figure 6 shows the model architecture. One frame in each clip is randomly selected as the input to the 2D model. In the 3D models, we extracted RGB frames, optical flows, and human-pose estimation for each frame in the clips as inputs. We randomly selected 64 consecutive frames from the clip. The frames were resized with the smallest dimension of 256 pixels with bilinear interpolation while preserving the aspect ratio. Pixel values were

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

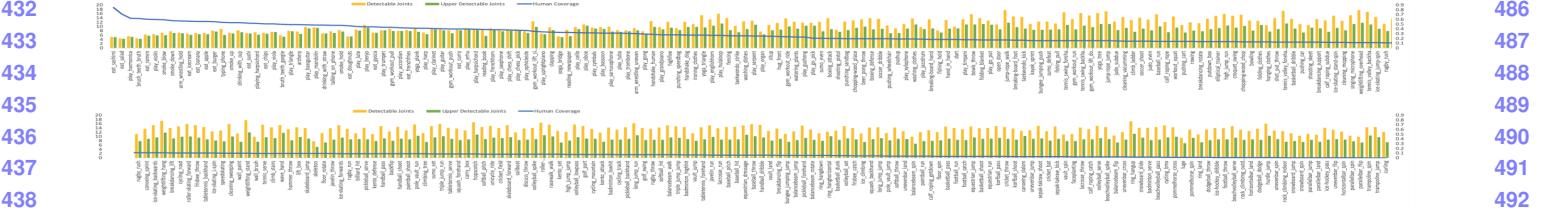


Figure 5. Average human coverage and number of detectable joints of different classes. Zoom in for details.

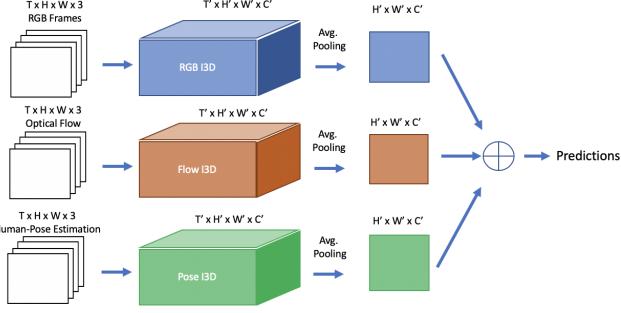


Figure 6. Baseline Model

rescaled to $[-1, 1]$. During training, we randomly selected a 224×224 image crop, while selecting the center 224×224 image crop from the video during testing. For the optical-flow stream, we converted the videos to grayscale and applied the PWC-Net [41] to extract the optical flows. The networks were trained with Softmax loss over 300 classes. Table 3 tabulates the classification results.

4.1.2 Results

	Top-1	Top-5	# of parameters
2D-CNN (RGB)	25.1	41.8	42M
I3D (RGB)	28.1	47.1	12M
I3D (Flow)	32.4	54.2	12M
I3D (Two-Stream)	36.8	59.5	12M
SlowFast (RGB)	29.2	45.5	62M

Table 3. Classification results on HAA300.

There are several noteworthy observations. First, the 3D-based model outperforms the 2D-based models. This suggests that temporal information is indeed a crucial factor for recognizing actions, since multiple actions may share similar human-poses at particular frames or contain similar background. Second, the contribution from the flow stream is slightly higher than that from the RGB stream. We believe this results from our vocabulary collection methods, where multiple action classes are collected from similar scenes/backgrounds containing similar RGB information (See Figure 7). Optical flow, which captures motion features, can be helpful for discriminating actions among

the same composite action.



(a) Balance-beam: Jump (b) Balance-beam: Rotate

Figure 7. Classes with similar backgrounds.

4.2 Human Pose Attention

With our human-centric dataset where human poses are detectable to a high extent, we evaluate the effectiveness of human pose estimations as an attention mechanism on HAA300.

4.2.1 Implementation Details

We extract human poses for each frame with OpenPose [2]. Outputs of OpenPose include a confidence map for body part detection and a map of Part Affinity Fields (PAFs), which is a set of 2D vector fields that encode the location and orientation of limbs over the image domain parts. We stack the confidence map together with the x and y components of the PAFs to form a three-dimensional representation map for each frame. We train the I3D model with the stacked feature maps, following the spatial and temporal cropping settings in our baseline model. Table 4 tabulates the classification results.

4.2.2 Results

Compared to the RGB-Stream I3D, the two-stream model with an additional human-pose input stream has improved the baseline performance for 10.5%. This suggests that human-pose information can be helpful for exploiting the human-centric nature of HAA300.

4.3 Fine-grained Action Labels

With the two-level action class label hierarchy of HAA300, we evaluate the effects of fine-grained labels to the performance of the coarser-grained actions.

	Top-1	Top-5	# of parameters
I3D (RGB)	28.1	47.1	12M
I3D (Human-Pose)	34.8	57.1	12M
I3D (Two-Stream)	39.6	63.2	12M

Table 4. Classification Results on HAA300

4.3.1 Fine-grained: HAA300

HAA300 contains 300 fine-grained action classes, each of which belongs to one of the 129 composite action classes. We compare the performance of models trained with and without an additional classification loss on the fine-grained action classes. We add an additional output stream to predict the fine-grained classes and sum up the classification loss on both the composite action classes and fine-grained action classes. We compare the performance on recognizing the 129 composite actions. We follow same settings in the baseline models for training. Table 5 tabulates the classification results.

4.3.2 Results

Method	Results
Composite classes	44.8
Composite + Fine-grained classes	48.5

Table 5. Results on two-level action class classification on HAA300

4.3.3 Fine-grained: AVA

We evaluate the effects of fine-grained action labels on the AVA dataset. We select a subset of the AVA dataset, containing the labels “Play musical instrument”. We then classify each of the labels into two levels of finer-grained actions. The first level groups actions with similar visual signals, including “string instruments”, “guitar instruments”, “brass instruments”, “woodwind instruments”, “keyboard instruments”, and “percussion instruments”. The second level contains a further classification where each class represents playing a specific musical instrument, such as “play violin”. This results in 6 coarser-level actions and 16 fine-grained actions. Table 6 shows the class structure. We follow the training settings in the previous section. Table 7 tabulates the classification results.

4.3.4 Results

Compared to models trained only with the composite labels, the model trained with an additional fine-grained classification loss has improved the performance of recognizing

Composite Action	Atomic Action
string instruments	play cello, play violin
guitar instruments	play guitar, play mandolin, play guitar, play guzheng
brass instruments	play trumpet, play trombone, play tuba, play frenchhorn, play saxophone
woodwind instruments	play clarinet, play harmonica
keyboard instruments	play piano, play accordion
percussion instruments	play drum

Table 6. Subset of composite actions in AVA

Method	Results
Composite classes	71.2
Composite + Fine-grained classes	74.4

Table 7. Results on two-level action class classification on AVA

composite actions. This shows that such fine-grained labels enabled the model to learn better features in order to differentiate the subtle details that distinguishes each action from others.

4.4 Decomposition of Composite Actions: UCF101

We apply HAA300 to decompose the subset of composite actions in UCF 101 into different atomic actions. This experiment shows the capability of HAA300 as a fine-grained model can decompose complex composite actions which have different visual expressions in one clip into different atomic actions.

4.4.1 Implementation Details

We begin by selecting a subset of UCF101 that the subset will only contain composite action classes (eg. frisbee catch, baseball pitch, and volleyball spiking are several counterexamples). Table 8 shows the full subset of UCF101 we picked for this experiment. We select a subset of HAA300 that contains the corresponding actions in the UCF101 subset. We follow the training settings in the previous section. We evaluate the performance by feeding different segments of clips the chosen subset from UCF101.

4.4.2 Results

Among all the classes in the subset, HAA300 can successfully decompose classes Floor Gymnastic, Long Jump, High Jump, Parallel Bars, Pole Vault, Sumo Wrestling, and Uneven Bars into their corresponding atomic actions, while

648	Composite Action	Atomic Action	702
649	Balance Beam	spin, walk	703
650	Floor Gymnastic	spin, rotate	704
651	Long Jump	run, jump	705
652	High Jump	run, jump	706
653	Jump Rope	jump, walk	707
654	Sumo Wrestling	even, defeat	708
655	Parallel Bars	spin, jump, land	709
656	Uneven Bars	spin, jump, land	710
657	Pole Vault	run, jump	711
658	Trampoline Jump	jump, spin	712

Table 8. Subset of composite actions in UCF1010

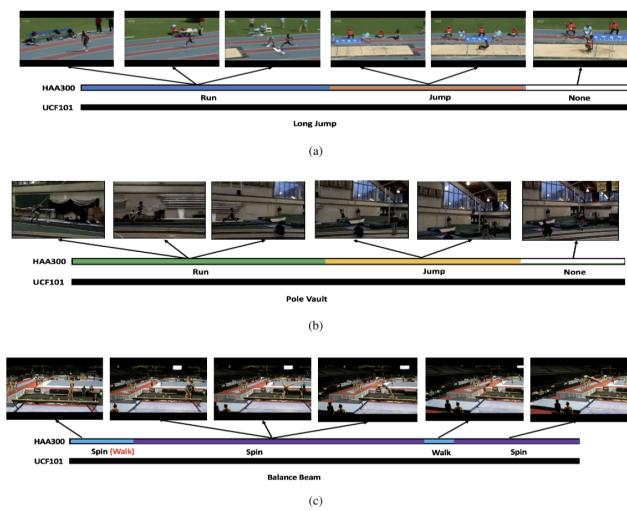


Figure 8. (a) and (b) are examples of successful classes that can be decomposed by HAA300. (c) is an example of failed class that HAA300 wrongly classifies atomic action walk to spin.

HAA300 fails to decompose classes Jump Rope, Balance Beam, and Trampoline Jump. Figure 8 shows the result of several successful and failed examples. There are 2 observations we identify for explaining the experiment results. Firstly, for the successful classes, we found out that their composite action can all be broken down into atomic actions that are in absolute sequence and are significantly different in visual expression. For example, composite action Floor Gymnastic is always composed by two atomic actions spinning then rotating. Same cases for High Jump, Long Jump and Pole Vault, all of them are always composed by two atomic actions running then jumping. While for the failed classes Jump Rope, Balance Beam, and Trampoline Jump, they are not composed by several atomic actions in absolute order. For example, atomic actions walking and spinning can appear alternatively during the whole composite action Balance Beam. This difference between successful classes and failed classes indicates that a crucial component for HAA300 to decompose a composite action is the tempo-

ral information. Secondly, we found out that the difference of visual expression between each atomic actions for composite actions in failed class is not clear and can sometime overlap with each other. For example, atomic actions walking sometime will be included in the beginning of spinning for composite action Balance Beam, thus make the model be confused with the boundary of different atomic actions.

4.5. AVA Action Detection

We evaluate the generalization performance of models trained with HAA300 with a subset of the AVA dataset. We show that the fine-grained nature of HAA300 can improve performance on composite action recognition.

4.5.1 Implementation Details

We begin by selecting a subset of the AVA dataset, containing the frames with labels “Run”, “Play Instruments”, and “Eat” from both the training set and the validation set. We selected these action categories since these actions are typically exclusive of each other, e.g. a person would normally not be “running” and “playing instrument” at the same time. To better assess the model’s ability to recognizing human-centric actions, we removed frames from the subset that contains excessive occlusion. This results in a subset of 7659 training clips and 766 validation clips, which we dub as **miniAVA**. We then select a subset of actions from HAA300 under categories “Run”, “Play Instruments”, and “Eat”, which contains 31 actions to produce **miniHAA**. We train a model with miniHAA following the implementation details in the baseline model. We compare the performance with a model trained from scratch with miniAVA and report the mean Average Precision (mAP) over the three classes. Table 9 tabulates the results.

4.5.2 Results

Compared to the model trained on miniAVA, the model trained on miniHAA has improved the mAP for 4.9%. The result seems even more surprising considering the fact that the size of miniHAA (~ 600) is much smaller than the training set in miniAVA (~ 7000). This suggests that a class-balanced dataset of human actions in HAA300 helps the model to learn the features in a more efficient manner even HAA300 is much smaller. Thus, HAA300 is easily extensible to cover more classes with fewer videos required while achieving better performance.

4.6. Transfer Learning Results

Models trained on HAA300 can be finetuned on other action datasets. By comparing finetuned models with models trained from scratch, we can assess the generalization performance of spatial-temporal features learned on our

	miniAVA	miniHAA
Play Instrument	36.8	34.7
Eat	45.9	55.6
Run	86.4	93.5
mAP	56.4	61.3

Table 9. Results on miniAVA

dataset. We evaluate the transfer learning results on 2 action classification benchmarks, UCF101 and HMDB51.

4.6.1 Implementation Details

We train and test on the I3D model with RGB inputs and optical flow inputs. We begin by evaluating the performance of models trained from scratch with UCF101 and HMDB51. We then evaluate the models pretrained with HAA300 followed by finetuning on the two datasets. We follow the spatial and temporal cropping settings in the baseline models. Table 10 tabulates the results.

4.6.2 Results

The results showed that models pretrained on HAA300 has an improved performance on the HMDB51 dataset, and a marginal gain difference on the performance of UCF101. The gains from pre-training with HAA300 on HMDB51 may be higher due to the nature of the dataset, where the average clip length (~ 3 seconds) is closer to that of HAA300 compared to UCF101 (~ 7 seconds).

Model	Pretrain	UCF101 Top-1	HMDB51 Top-1
I3D (RGB)	None	58.5	22.2
I3D (RGB)	HAA300	70.9	39.1
I3D (Flow)	None	85.2	50.7
I3D (Flow)	HAA300	82.8	53.5
I3D (Two-Stream)	None	85.5	53.1
I3D (Two-Stream)	HAA300	85.6	58.6

Table 10. Performance on the UCF-101 and HMDB-51 test sets (split 1 of both) starting with / without HAA300 pretrained weights.

5. Conclusion

This paper introduces HAA300, a new human action dataset with fine-grained atomic action labels and human-centric clip annotations. We have provided baseline benchmarks on this dataset and demonstrated the effectiveness of human-pose estimations on improving the performance of recognizing human actions. We have also shown how fine-grained action labels can help recognizing composite actions decompose composite actions into atomic atomic ac-

tions. We hope HAA300 will inspire new methods and architectures for modeling higher-complexity human actions.

6. Appendix

Figure 9 shows some sample classes in HAA300.



Figure 9. Sample classes in HAA300.

Figure 10 shows the full hierarchy of HAA300.

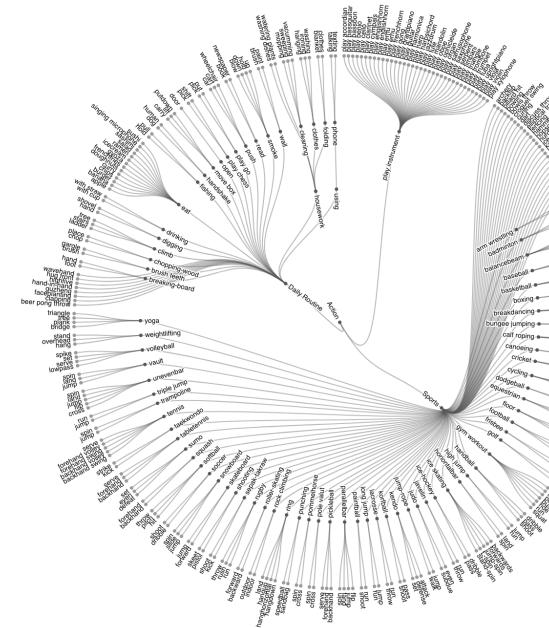


Figure 10. Full hierarchical structure of all classes in HAA300.

864

References

- [1] Moshe Blank, Lena Gorelick, Eli Shechtman, Michal Irani, and Ronen Basri. Actions as space-time shapes. In *2005 IEEE International Conference on Computer Vision (ICCV 2005)*, 2005. 2
- [2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7291–7299, 2017. 4, 5
- [3] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*, 2017. 2
- [4] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiaxuan Wang, and Jia Deng. HICO: A benchmark for recognizing human-object interactions in images. In *2015 IEEE International Conference on Computer Vision (ICCV 2015)*, 2015. 3
- [5] Sourish Chaudhuri, Joseph Roth, Daniel P. W. Ellis, Andrew C. Gallagher, Liat Kaver, Rebecca Marvin, Caroline Pantofaru, Nathan Reale, Loretta Guarino Reid, Kevin W. Wilson, and Zhonghua Xi. Ava-speech: A densely labeled dataset of speech activity in movies. In *INTERSPEECH*, 2018. 1
- [6] Nieves Crasto, Philippe Weinzaepfel, Karteek Alahari, and Cordelia Schmid. Mars: Motion-augmented rgb stream for action recognition. In *2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2019)*, 2019. 3
- [7] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009. 1
- [8] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Mohammad Mahdi Arzani, Rahman Yousefzadeh, Juergen Gall, and Luc Van Gool. Spatio-temporal channel correlation networks for action classification. In *2018 European Conference on Computer Vision (ECCV 2018)*, 2018. 2
- [9] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, and Kate Saenko. Long-term recurrent convolutional networks for visual recognition and description. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2015)*, 2015. 2
- [10] H. Warren Dunham. Midwest and its children: The psychological ecology of an american town. roger g. barker , herbert f. wright. *American Journal of Sociology*, 62(2), 1956. 3
- [11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. *CoRR*, abs/1812.03982, 2018. 2
- [12] Rohit Girdhar, João Carreira, Carl Doersch, and Andrew Zisserman. A better baseline for AVA. *CoRR*, abs/1807.10066, 2018. 3
- [13] Raghad Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The “something something” video database for learning and evaluating visual common sense. In *ICCV*, volume 1, page 3, 2017. 1, 2
- [14] Raghad Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, Florian Hoppe, Christian Thurau, Ingo Bax, and Roland Memisevic. The “something something” video database for learning and evaluating visual common sense. In *2017 IEEE International Conference on Computer Vision (ICCV 2017)*, 2017. 3
- [15] Chunhui Gu, Chen Sun, David A. Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2018)*, 2018. 1, 2, 3
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2016)*, 2016. 2, 4
- [17] Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J. Black. Towards understanding action recognition. In *2013 IEEE International Conference on Computer Vision (ICCV 2013)*, 2013. 2
- [18] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3d convolutional neural networks for human action recognition. In *2010 International Conference on Machine Learning (ICML 2010)*, 2010. 2
- [19] Boyuan Jiang, Mengmeng Wang, Weihao Gan, Wei Wu, and Junjie Yan. STM: spatiotemporal and motion encoding for action recognition. *CoRR*, abs/1908.02486, 2019. 2
- [20] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Fei-Fei Li. Large-scale video classification with convolutional neural networks. In *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2014)*, 2014. 2
- [21] Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *CoRR*, abs/1705.06950, 2017. 1, 2
- [22] Yan Ke, Rahul Sukthankar, and Martial Hebert. Efficient visual event detection using volumetric features. In *2005 IEEE International Conference on Computer Vision (ICCV 2005)*, 2005. 2
- [23] Hildegard Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso A. Poggio, and Thomas Serre. HMDB: A large video database for human motion recognition. In *2011 IEEE International Conference on Computer Vision (ICCV 2011)*, 2011. 2
- [24] Ivan Laptev, Marcin Marszalek, Cordelia Schmid, and Benjamin Rozenfeld. Learning realistic human actions from movies. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008)*, 2008. 2
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In

- 972 David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–
973 755, Cham, 2014. Springer International Publishing. 1 1026
974 1027
975 [26] Marcin Marszalek, Ivan Laptev, and Cordelia Schmid. Actions in context. In *2009 IEEE Computer Society Conference
976 on Computer Vision and Pattern Recognition (CVPR 2009)*, 2009. 2 1028
977 1029
978 [27] Pascal Mettes, Jan C. van Gemert, and Cees G. M. Snoek. Spot on: Action localization from pointly-supervised proposals. In *2016 European Conference on Computer Vision (ECCV 2016)*, 2016. 2 1030
979 1031
980 [28] Mathew Monfort, Bolei Zhou, Sarah Adel Bargal, Alex Andonian, Tom Yan, Kandan Ramakrishnan, Lisa M. Brown,
981 Quanfu Fan, Dan Gutfreund, Carl Vondrick, and Aude Oliva. Moments in time dataset: one million videos for event under-
982 standing. *CoRR*, abs/1801.03150, 2018. 1, 2 1032
983 1033
984 [29] Juan Carlos Niebles, Hongcheng Wang, and Fei-Fei Li. Un-
985 supervised learning of human action categories using spatial-
986 temporal words. *International Journal of Computer Vision*,
987 79(3), 2008. 2 1034
988 1035
989 [30] Maxime Oquab, Léon Bottou, Ivan Laptev, and Josef Sivic.
990 Learning and transferring mid-level image representations
991 using convolutional neural networks. In *2014 IEEE Conference
992 on Computer Vision and Pattern Recognition (CVPR 2014)*, 2014. 2 1036
993 1037
994 [31] Bowen Pan, Jiankai Sun, Wuwei Lin, Limin Wang, and
995 Weiyao Lin. Cross-stream selective networks for action
996 recognition. In *2017 IEEE Conference on Computer Vision
997 and Pattern Recognition (CVPR) Workshops (CVPR 2017)*,
998 2019. 3 1038
999 1039
1000 [32] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-
1001 temporal representation with pseudo-3d residual networks.
1002 In *IEEE International Conference on Computer Vision (ICCV 2017)*, 2017. 3 1040
1003 1041
1004 [33] Mikel D. Rodriguez, Javed Ahmed, and Mubarak Shah. Action
1005 MACH a spatio-temporal maximum average correlation
1006 height filter for action recognition. In *2008 IEEE Computer
1007 Society Conference on Computer Vision and Pattern Recog-
1008 nition (CVPR 2008)*, 2008. 2 1042
1009 1043
1010 [34] Joseph Roth, Sourish Chaudhuri, Ondrej Klejch, Rad-
1011 hika Marvin, Andrew Gallagher, Liat Kaver, Sharadh
1012 Ramaswamy, Arkadiusz Stopczynski, Cordelia Schmid,
1013 Zhonghua Xi, and Caroline Pantofaru. Ava-activespeaker:
1014 An audio-visual dataset for active speaker detection, 2019. 1 1044
1015 1045
1016 [35] J. Roth, S. Chaudhuri, O. Klejch, R. Marvin, A. Gallagher,
1017 L. Kaver, S. Ramaswamy, A. Stopczynski, C. Schmid, Z.
1018 Xi, and C. Pantofaru. Ava-activespeaker: An audio-visual
1019 dataset for active speaker detection. *arXiv:1901.01342*,
1020 2019. 1 1046
1021 1047
1022 [36] Christian Schüldt, Ivan Laptev, and Barbara Caputo. Recog-
1023 nizing human actions: A local SVM approach. In *2004 IEEE
1024 International Conference on Pattern Recognition (CVPR
1025 2004)*, 2004. 2 1048
1026 1049
1027 [37] Gunnar A. Sigurdsson, Gü̈l Varol, Xiaolong Wang, Ali
1028 Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in
1029 homes: Crowdsourcing data collection for activity under-
1030 standing. In *European Conference on Computer Vision*,
1031 2016. 2 1051
1032 1053
1033 1054
1034 1055
1035 1056
1036 1057
1037 1058
1038 1059
1039 1060
1040 1061
1041 1062
1042 1063
1043 1064
1044 1065
1045 1066
1046 1067
1047 1068
1048 1069
1049 1070
1050 1071
1051 1072
1052 1073
1053 1074
1054 1075
1055 1076
1056 1077
1057 1078
1058 1079