
Reformulating Zero-shot Action Recognition for Multi-label Actions (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 AVA Dataset Evaluation

2 1.1 Extracting Video Clips

3 Since the AVA dataset consists of multiple actors within one video and ZSAR focuses only on the
4 classification task, we extract clips centered on the ground-truth bounding boxes for each actor in the
5 video. Standard video models expect frame dimensions with the same height and width, so we crop
6 a square region around the actor and resize it to the network specific dimensions (112×112). We
7 present some examples of AVA video frames with their annotations as well as the generated crops in
8 Figure 1. This square crop can cause multiple actors to appear within one clip, as seen in the second
9 example, but it ensures the aspect ratio of the person is not altered, which is necessary as this is the
10 manner in which the video model is trained.

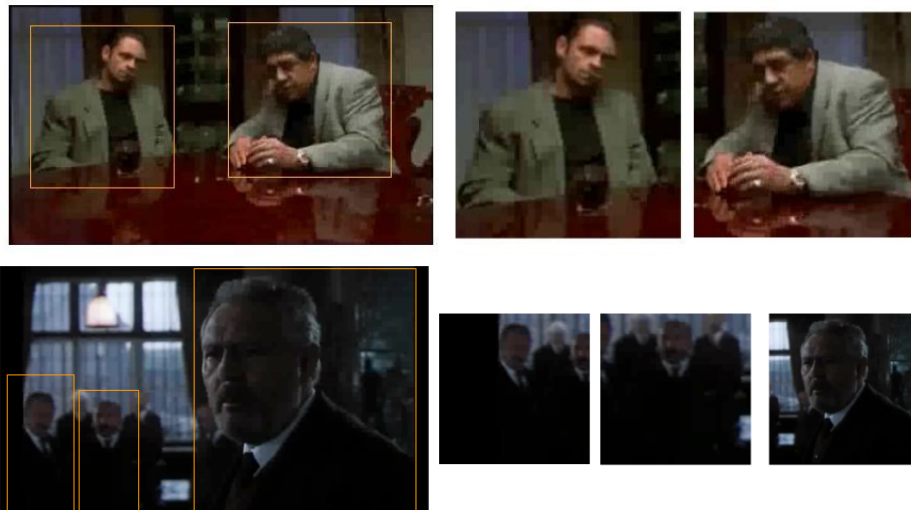


Figure 1: Example of original ground-truth bounding boxes (left) in the AVA dataset, with the cropped actors on the right.

11 1.2 Generating Multiple Predictions and Confidences

12 As previous methods for ZSAR tend to be designed for single-label action classification, we adjust
13 these methods to generate multiple predictions along with prediction confidences. For PS-ZSAR

14 prediction confidences are obtained from the softmax probabilities output by our pair-wise similarity
 15 function. To obtain confidence scores from the method in Brattoli *et al.* [4], we apply a softmax oper-
 16 ation on the inverse cosine distances between the video model’s output and the semantic embeddings:
 17

$$P(y|x) = \frac{\exp(-d(f_\theta(x), \psi(y))/\gamma)}{\sum_{y' \in \mathcal{U}} \exp(-d(f_\theta(x), \psi(y'))/\gamma)}, \quad (1)$$

18 where d is the cosine distance. As the distances between embeddings tend to be small, we use a
 19 temperature parameter $\gamma \leq 1$ increase distances before being passed through the softmax. We find
 20 that selecting $\gamma = 0.1$ leads to best results.

21 To obtain multiple predictions from a given method there are several approaches. One trivial approach
 22 is to select the top-k predictions for a given sample. The main issue with this approach is that it may
 23 over-predict classes when k is too large or under-predict when k is too small. Another approach is to
 24 predict all classes, in which case the mAP evaluation would ignore most low-confidence predictions.
 25 This alleviates the issue of under-predicting, but will always over-predict. Finally, we can predict
 26 classes based on a confidence threshold, in the manner described in equation 3 of the main paper.

Table 1: mAP Results on AVA Dataset

	Top-1	Top-3	Top-5	No threshold	Threshold
Brattoli <i>et al.</i> [4] ($\gamma = 1$)	1.6	2.1	2.4	3.1	6.2
Brattoli <i>et al.</i> [4] ($\gamma = 0.1$)	1.6	2.1	2.3	3.3	6.4
Ours (word2vec)	1.6	3.0	3.4	6.4	6.5
Ours (sent2vec)	1.5	3.0	3.5	5.7	7.0

27 We present results for all approaches in Table 1. It shows that the use of thresholding on predicted
 28 probabilities leads to best results. Interestingly, only the top-1 predictions for both methods achieve
 29 similar performance, but when it is increased to top-5, the gap between mAP scores increases. This
 30 poor performance is due to the nearest neighbor classification which does not allow semantically
 31 dissimilar classes to be predicted confidently. On the other hand, our approach can have multiple
 32 dissimilar classes in the top-5 predictions.

33 2 RareAct Evaluation

34 RareAct is a dataset compiled from rarely co-occurring nouns and verbs such as "microwave show"
 35 or "blend phone". It is meant to be "an evaluation dataset notably meant to be used to evaluate
 36 models trained on the HowTo100M dataset" [39]. We use RareAct in our work to evaluate how
 37 well zero-shot methods can deal with action classes which are extremely different from those seen
 38 during training. In the RareAct work [39], the authors propose different metrics (mWAP and mSAP).
 39 However, we evaluate our method using the top-1 and top-5 accuracy since the purpose of this work is
 40 to create a strong zero-shot classifier rather than learn a joint visual-textual model from a large-scale
 41 instructional dataset (i.e. HowTo100M).

42 3 Evaluation on UCF-101 and HMDB datasets using Random seeds

43 In Brattoli *et al.* [4], one of the primary methods of evaluation involves generating 10 different
 44 testing sets from UCF101 and HMDB by randomly choosing half of the classes. This is the standard
 45 evaluation in all works prior to [2], since lack of access to the Kinetics-700 dataset made testing on
 46 the full UCF-101 or HMDB dataset infeasible. However, we find that this metric is problematic as the
 47 results are dependant almost entirely on the random seed (implemented with numpy’s rand package)
 48 used to choose which classes to test with. To illustrate this issue, we use 10k random seeds, and report
 49 the results in Table 2. The results for Brattoli *et al.* are obtained from the publicly available code and
 50 model weights. When results are averaged over all 10k seeds PS-ZSAR outperforms Brattoli *et al.*.
 51 Furthermore, we find that our method achieves higher accuracy on 58.3% of the seeds on UCF-101
 52 and 76.5% of the seeds on HMDB .

Table 2: Evaluation on 50% of the UCF-101 and HMDB classes over 10k random seeds. Reported are the mean and standard deviation ($\mu \pm \sigma$).

	UCF101	HMDB
Brattoli <i>et al.</i> [4]	39.3 \pm 4.3	25.1 \pm 4.4
PS-ZSAR (ours)	40.1 \pm 3.8	27.3 \pm 4.0



Dataset	UCF101 Class	MEVA Class
Class Name	BaseballPitch	person_opens_car_door
Encoder Input	"Baseball" "Pitch"	"A person opening the door to a vehicle. The only necessary track in this event is the vehicle. The vehicle door is not independently annotated from the vehicle. This event often overlaps with entering/exiting; however, can be independent or absent from these events."
Example		

Figure 2: Example of the natural language descriptions of MEVA classes versus simple class names of UCF101 classes. Note also for MEVA videos are captured through surveillance camera, and thus actions are lower resolution, as well as less visually apparent.

53 As our results show, Bratolli *et al.* [4] scores an average of 39.3 on UCF101 and 25.1 on HMDB.
 54 However, their reported results are 48.0 and 32.7 respectively, nearly two standard deviations above
 55 the mean. In the interest of reporting the most comparable results despite the drawbacks of this
 56 evaluation method, we searched for a seed that resulted in their method achieving as close to their
 57 reported scores as possible. In the main paper, we then reported our accuracy on that same seed: 49.2
 58 and 33.8 for UCF-101 and HMDB respectively. As this evaluation protocol (i.e. selecting only 10
 59 splits with 50% of the classes) can lead to noisy results, we argue future ZSAR should be evaluate on
 60 the entirety of UCF-101 and HMDB.

61 4 MEVA Dataset Activity Descriptions

62 Contrary to conventional video datasets which use class names to generate semantic embeddings, the
 63 MEVA dataset contain natural language descriptions of the action classes. For example, the action
 64 *carrying* has the description "A person carrying an object up to half the size of the person, where
 65 the person's gait has not been substantially modified. The object may be carried in either hand, with
 66 both hands, or on one's back" and the action *falling* has the description "A person falling by either
 67 (1) losing one's balance and possibly collapsing, or (2) moving downward from a higher to a lower
 68 level." These lengthy descriptions allow the ZSAR method to learn a richer semantic embedding
 69 which is useful for classifying surprise activities.

70 5 Method Limitations

71 We analyse how PS-ZSAR performs on the UCF-101 dataset to understand the limitations of the
 72 approach. We find that ZSAR methods achieve strong performance on certain classes, while many
 73 classes tend to be ignored and not predicted. We present 10 classes on which our method achieves
 74 0% in Table 3. PS-ZSAR tends to predict classes which are visually similar to the target class. For

75 instance, videos with the "Jump Rope" and "Jumping Jack" actions tend to be predicted as "Handstand
76 Pushups" since all three actions involve similar motions (i.e. repetitive up and down motions). This is
77 a limitation for not only our approach, but most ZSAR approaches. For example, Bratolli *et al.* [4]
78 achieve 0% accuracy on 36 classes and PS-ZSAR achieves 0% accuracy on 22 classes. We believe
79 solving this problem would be an interesting avenue for future work.

Table 3: Ten classes which PS-ZSAR performs worst on in the UCF-101 dataset. We include the class name, the accuracy, and the class predicted for most videos of the given class.

Class Name	Most Predicted
Jump Rope	Handstand Pushups
Jumping Jack	Handstand Pushups
Hula Hoop	Tai Chi
YoYo	SalsaSpin
Front Crawl	Breast Stroke
Bowling	Basketball
Parallel Bars	Trampoline Jumping
Playing Daf	Head Massage
Playing Violin	Playing Flute
Pole Vault	Trampoline Jumping

80 References

- 81 [1] Valter Estevam, Helio Pedrini, and David Menotti. Zero-shot action recognition in videos: A
82 survey. *Neurocomputing*, 439:159–175, 2021.
- 83 [2] Yi Zhu, Yang Long, Yu Guan, Shawn Newsam, and Ling Shao. Towards universal representation
84 for unseen action recognition. In *Proceedings of the IEEE conference on computer vision and
85 pattern recognition*, pages 9436–9445, 2018.
- 86 [3] Meera Hahn, Andrew Silva, and James M Rehg. Action2vec: A crossmodal embedding
87 approach to action learning. *arXiv preprint arXiv:1901.00484*, 2019.
- 88 [4] Biagio Brattoli, Joseph Tighe, Fedor Zhdanov, Pietro Perona, and Krzysztof Chalupka. Rethink-
89 ing zero-shot video classification: End-to-end training for realistic applications. In *Proceedings
90 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4613–4623,
91 2020.
- 92 [5] Ioannis Alexiou, Tao Xiang, and Shaogang Gong. Exploring synonyms as context in zero-shot
93 action recognition. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages
94 4190–4194. IEEE, 2016.
- 95 [6] Xun Xu, Timothy M Hospedales, and Shaogang Gong. Multi-task zero-shot action recognition
96 with prioritised data augmentation. In *European Conference on Computer Vision*, pages 343–359.
97 Springer, 2016.
- 98 [7] Xun Xu, Timothy Hospedales, and Shaogang Gong. Transductive zero-shot action recognition
99 by word-vector embedding. *International Journal of Computer Vision*, 123(3):309–333, 2017.
- 100 [8] Qian Wang and Ke Chen. Zero-shot visual recognition via bidirectional latent embedding.
101 *International Journal of Computer Vision*, 124(3):356–383, 2017.
- 102 [9] Yongqin Xian, Bernt Schiele, and Zeynep Akata. Zero-shot learning-the good, the bad and the
103 ugly. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
104 pages 4582–4591, 2017.
- 105 [10] Devraj Mandal, Sanath Narayan, Sai Kumar Dwivedi, Vikram Gupta, Shuaib Ahmed, Fa-
106 had Shahbaz Khan, and Ling Shao. Out-of-distribution detection for generalized zero-shot
107 action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
108 Pattern Recognition*, pages 9985–9993, 2019.

- 109 [11] Xingyu Chen, Xuguang Lan, Fuchun Sun, and Nanning Zheng. A boundary based out-of-
110 distribution classifier for generalized zero-shot learning. In *European Conference on Computer*
111 *Vision*, pages 572–588. Springer, 2020.
- 112 [12] William Thong and Cees GM Snoek. Bias-awareness for zero-shot learning the seen and unseen.
113 *arXiv preprint arXiv:2008.11185*, 2020.
- 114 [13] Jiamin Wu, Tianzhu Zhang, Zheng-Jun Zha, Jiebo Luo, Yongdong Zhang, and Feng Wu. Self-
115 supervised domain-aware generative network for generalized zero-shot learning. In *Proceedings*
116 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12767–12776,
117 2020.
- 118 [14] Shaobo Min, Hantao Yao, Hongtao Xie, Chaoqun Wang, Zheng-Jun Zha, and Yongdong Zhang.
119 Domain-aware visual bias eliminating for generalized zero-shot learning. In *Proceedings of*
120 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12664–12673,
121 2020.
- 122 [15] Mina Bishay, Georgios Zoumpourlis, and Ioannis Patras. Tarn: Temporal attentive relation
123 network for few-shot and zero-shot action recognition. *arXiv preprint arXiv:1907.09021*, 2019.
- 124 [16] Chuang Gan, Ming Lin, Yi Yang, Gerard Melo, and Alexander G Hauptmann. Concepts not
125 alone: Exploring pairwise relationships for zero-shot video activity recognition. In *Proceedings*
126 *of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- 127 [17] AJ Piergiovanni and Michael S Ryoo. Learning shared multimodal embeddings with unpaired
128 data. *CoRR*, 2018.
- 129 [18] Bowen Zhang, Hexiang Hu, and Fei Sha. Cross-modal and hierarchical modeling of video and
130 text. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 374–390,
131 2018.
- 132 [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
133 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
134 pages 770–778, 2016.
- 135 [20] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spa-
136 tiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international*
137 *conference on computer vision*, pages 4489–4497, 2015.
- 138 [21] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word
139 representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 140 [22] Qiang Qiu, Zhuolin Jiang, and Rama Chellappa. Sparse dictionary-based representation and
141 recognition of action attributes. In *2011 International Conference on Computer Vision*, pages
142 707–714. IEEE, 2011.
- 143 [23] Jingen Liu, Benjamin Kuipers, and Silvio Savarese. Recognizing human actions by attributes.
144 In *CVPR 2011*, pages 3337–3344. IEEE, 2011.
- 145 [24] Chuang Gan, Ming Lin, Yi Yang, Yueting Zhuang, and Alexander G Hauptmann. Exploring
146 semantic inter-class relationships (sir) for zero-shot action recognition. In *Proceedings of the*
147 *AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- 148 [25] Chung-Wei Lee, Wei Fang, Chih-Kuan Yeh, and Yu-Chiang Frank Wang. Multi-label zero-shot
149 learning with structured knowledge graphs. In *Proceedings of the IEEE conference on computer*
150 *vision and pattern recognition*, pages 1576–1585, 2018.
- 151 [26] Dat Huynh and Ehsan Elhamifar. A shared multi-attention framework for multi-label zero-
152 shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
153 *Recognition*, pages 8776–8786, 2020.
- 154 [27] He Huang, Yuanwei Chen, Wei Tang, Wenhao Zheng, Qing-Guo Chen, Yao Hu, and Philip
155 Yu. Multi-label zero-shot classification by learning to transfer from external knowledge. *arXiv*
156 *preprint arXiv:2007.15610*, 2020.

- 157 [28] Qian Wang and Ke Chen. Multi-label zero-shot human action recognition via joint latent ranking
158 embedding. *Neural Networks*, 122:1–23, 2020.
- 159 [29] Mahdi Naser Moghadasi and Yu Zhuang. Sent2vec: A new sentence embedding representation
160 with sentimental semantic. In *2020 IEEE International Conference on Big Data (Big Data)*,
161 pages 4672–4680. IEEE, 2020.
- 162 [30] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-
163 networks. *arXiv preprint arXiv:1908.10084*, 2019.
- 164 [31] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A
165 closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE
166 conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018.
- 167 [32] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
168 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
169 pytorch. 2017.
- 170 [33] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool.
171 Temporal segment networks: Towards good practices for deep action recognition. In *European
172 conference on computer vision*, pages 20–36. Springer, 2016.
- 173 [34] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
174 arXiv:1412.6980*, 2014.
- 175 [35] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-
176 700 human action dataset. *arXiv preprint arXiv:1907.06987*, 2019.
- 177 [36] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
178 Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference
179 on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010.
- 180 [37] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human
181 actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 182 [38] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre.
183 Hmdb: a large video database for human motion recognition. In *2011 International conference
184 on computer vision*, pages 2556–2563. IEEE, 2011.
- 185 [39] Antoine Miech, Jean-Baptiste Alayrac, Ivan Laptev, Josef Sivic, and Andrew Zisserman. Rareact:
186 A video dataset of unusual interactions. *arXiv preprint arXiv:2008.01018*, 2020.
- 187 [40] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sud-
188 heendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A
189 video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE
190 Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018.
- 191 [41] Kellie Corona, Katie Osterdahl, Roderic Collins, and Anthony Hoogs. Meva: A large-scale
192 multiview, multimodal video dataset for activity detection. In *Proceedings of the IEEE/CVF
193 Winter Conference on Applications of Computer Vision (WACV)*, pages 1060–1068, January
194 2021.
- 195 [42] Kitware inc, the multiview extended video with activities (meva) dataset.
- 196 [43] Hugo Larochelle, Dumitru Erhan, and Yoshua Bengio. Zero-data learning of new tasks. In
197 *AAAI*, volume 1, page 3, 2008.
- 198 [44] Alina Roitberg, Manuel Martinez, Monica Haurilet, and Rainer Stiefelhagen. Towards a fair
199 evaluation of zero-shot action recognition using external data. In *Proceedings of the European
200 Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.
- 201 [45] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer
202 network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
203 Recognition*, pages 244–253, 2019.
- 204 [46] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine
205 learning research*, 9(11), 2008.