

---

# Enabling Detailed Action Recognition Evaluation Through Video Dataset Augmentation

---

**Jihoon Chung**  
Princeton University  
jc5933@princeton.edu

**Yu Wu**  
Princeton University  
yuwu@princeton.edu

**Olga Russakovsky**  
Princeton University  
olgarus@princeton.edu

## Abstract

1 It is well-known in the video understanding community that human action recogni-  
2 tion models suffer from background bias, over-relying on scene cues in making  
3 their prediction. However, it is difficult to *quantify* this effect given the existing  
4 evaluation frameworks. We introduce the Human-centric Analysis Toolkit (HAT),  
5 which enable evaluation of the learned background bias without the need for new  
6 manual video annotation. It does so by automatically generating synthetically ma-  
7 nipulated videos, leveraging the recent advances in image segmentation and video  
8 inpainting. Using HAT we perform an extensive analysis of 74 action recognition  
9 models trained on the Kinetics dataset. We confirm that all these models focus  
10 more on the scene background than on the human motion; further, we demonstrate  
11 that certain model design decisions (such as training with fewer frames per video  
12 or using dense as opposed to uniform temporal sampling) appear to worsen the  
13 background bias. We open-source HAT to enable the community to leverage its  
14 metrics to design more robust and generalizable human action recognition models.<sup>1</sup>

## 15 1 Introduction

16 Human action recognition is about understanding what the *human* in the video is doing; however,  
17 human action recognition models frequently rely on background cues to make their predictions.  
18 Works such as [6, 33, 59, 60, 70] have leveraged visualization tools like GradCam [46] to demonstrate  
19 that the video background significantly influences the prediction of human action recognition models.  
20 This occurs due to representation bias in the dataset, where particular actions (e.g., eating) tend to  
21 occur in particular environments (e.g., kitchens). Such concerns limit the practical usability and  
22 generalizability of models despite the impressive overall progress in the field [36, 61, 67].

23 While it is known that this background bias phenomenon is occurring, *quantifying* the degree to which  
24 it is occurring is still necessary. Being able to accurately assess how much human action recognition  
25 models rely on human features rather than background scene cues would allow researchers to compare  
26 different model designs and select the ones that would be robust to their unique test domains. Efforts  
27 such as [34, 63] have introduced datasets specifically curated for quantifying background bias;

---

<sup>1</sup><https://github.com/princetonvisualai/HAT>

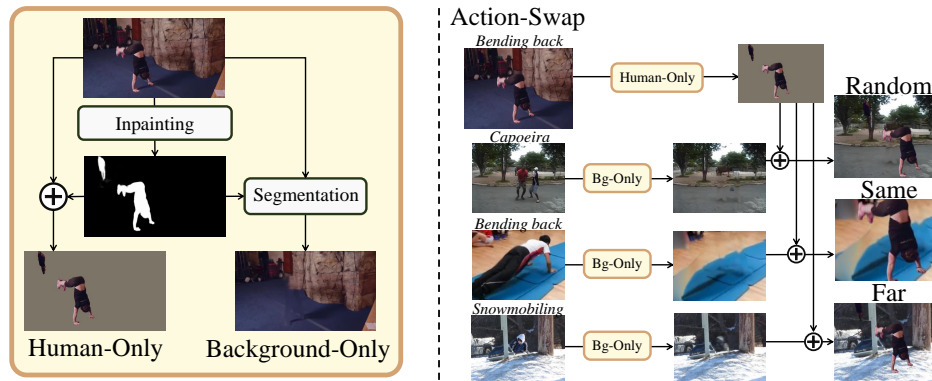


Figure 1: The pipeline of our Human-centric Analysis Toolkit (HAT). **Left:** HAT takes a video, segments the spatio-temporal human figure, and generates the Human-Only and Background-Only videos. **Right:** HAT generates Action-Swap videos by pasting the same human figure onto the Background-Only video from the same, a randomly-selected, and a far (dissimilar) action class.

28 however, scaling up or generalizing their approaches may be prohibitively expensive due to the  
 29 reliance on manual annotation.

30 In this work, we introduce the Human-centric Analysis Toolkit (HAT) to measure background bias in  
 31 human action recognition models without the need for costly human annotation. We leverage recent  
 32 improvements in image segmentation [24, 32, 38] and video inpainting [29, 35, 37] to automatically  
 33 synthesize counterfactual videos containing Human-Only (a spatio-temporal segmentation of the  
 34 human figure against a gray background), Background-Only (the video with the human removed via  
 35 inpainting) or Action-Swap (human figure against an unusual background). Examples are shown in  
 36 Figure 1. This process is efficient and scalable, requiring no manual annotation. HAT thus enables  
 37 us to evaluate the sensitivity of human action understanding models to the different visual cues by  
 38 comparing the accuracy on the original and synthetically manipulated videos.

39 We demonstrate the capabilities of HAT by running extensive analysis of human action recogni-  
 40 tion models trained on the Kinetics-400 [28] dataset. Concretely, we evaluate 74 trained models,  
 41 corresponding to 14 different model designs (TSN [61], I3D [5], Non-local Neural Networks [62],  
 42 R(2+1)D [54], TSM [36], SlowFast/SlowOnly [15], CSN [58], TIN [49], TPN [68], X3D [14],  
 43 OmniSource [12], TANet [39], and TimeSformer [4]) with varying hyperparameters and backbone  
 44 architectures provided by the MMAAction2 [8] implementation. We find, for example:

- 45 • All 74 models exhibit strong background bias. When evaluated on the Action-Swap videos,  
 46 the 74 models predicted the action class of the human 16.8% of the time on average – but  
 47 predicted the action class of the randomly-selected background 29.5% of the time!
- 48 • Models trained with fewer frames per video appear to be more prone to background bias. For  
 49 example, the TSN-based models [61] trained with 8, 5 and 3 frames per video retain 0.679,  
 50 0.683 and 0.694 of their original accuracy respectively when evaluated on the Background-  
 51 Only videos, demonstrating consistently high and somewhat *increasing* background bias.
- 52 • Models trained with dense temporal sampling around a single timestep appear to be more  
 53 prone to background bias compared to models trained with uniform sampling throughout the  
 54 video. For example, when evaluated on the Background-Only videos as above, TSM-based  
 55 models [36] with dense sampling exhibit strong background bias by retaining 0.703 of the  
 56 original accuracy compared to only 0.675 with uniform sampling.

57 Overall, we make three contributions. First, we develop and open-source the Human-centric Analysis  
 58 Toolkit (HAT), which generates synthetic videos to evaluate the background bias learned by human  
 59 action recognition models. Second, we demonstrate its capabilities through extensive evaluation of 74  
 60 released models. Finally, we show that HAT can identify the design choices that appear to influence  
 61 the amount of background bias learned by the model, helping inform future model design.

## 62 2 Related Work

63 **Human Action Recognition Models.** Currently, human action recognition is largely dominated by  
64 deep learning methods. With strong success in image-based tasks [10, 20, 30, 51], CNN-based deep  
65 learning models [12, 14, 15, 36, 39, 49, 54, 58, 61, 62, 68] were the go-to method for human action  
66 recognition, with gradual improvements in the model structure going from 2D-CNN [20, 44, 61]  
67 to 3D-CNN [5, 14, 15, 54] to CNN models with specific temporal modeling [36, 39, 49]. A recent  
68 trend [2, 4] in human action recognition is to use a transformer module [57] as it has shown good  
69 performance [11] in image-based tasks. Another trend [4, 12, 58] is to incorporate large-scale  
70 datasets [16, 26, 65] into the training. In this work, we evaluate multiple action recognition models [4,  
71 5, 12, 14, 15, 36, 39, 49, 54, 58, 61, 62, 68] in an effort to identify design decisions which appear to  
72 correlated with learned background bias.

73 **Human Action Recognition Dataset.** The early datasets [45, 66] offered a handful of human action  
74 classes that were collected in a controlled environment. UCF101 [52] and HMDB51 [31] were one  
75 of the first few datasets that were suitable for machine learning tasks. Although there are many  
76 different human action datasets [7, 9, 19, 47, 48, 72] these days, the most popular dataset must  
77 be Kinetics-400 [28], due to its large size and the variety of actions. However, due to the cost of  
78 collecting video datasets, the size of the dataset is still smaller compared to image datasets. Synthetic  
79 datasets [13, 17, 27, 41, 53], often used mixed with the real dataset, are popular methods of collecting  
80 data in an affordable manner. In human action recognition, the synthetic datasets are often used for  
81 training dataset [55, 56], and the model is tested on real videos. In this work, we generate synthetic  
82 counterfactual videos to enable detailed model evaluation without the need for costly annotation.

83 **Human-centric Analysis.** As the models have grown more complex, there has been an increased  
84 need for frameworks that provide insights into the model behavior beyond just a single accuracy  
85 number. Efforts have included model interpretability techniques [3, 46], detailed error analysis  
86 using additional manual annotations [1, 22, 43, 50], and (recently) stress-testing using automatically  
87 generated text or image data [25, 42]. There are a number of works studying specifically the  
88 impact of the human figure on human action recognition models. A common strategy employed  
89 by [6, 33, 59, 60, 70] is to use the GradCam [46] visualization to qualitatively demonstrate that the  
90 model’s attention is on the background cues rather than on the human in the video. Several of these  
91 papers [6, 59, 60, 70] propose methods to mitigate the effects of background bias during training;  
92 they evaluate its success both qualitatively through GradCam and quantitatively via accuracy on a  
93 downstream action recognition task (after fine-tuning the model trained with their new background-  
94 debiasing method). While this successfully demonstrates that their innovation is effective for model  
95 pre-training, it does not directly measure the learned background bias. The most natural analysis is  
96 to collect specific datasets [7, 18, 34, 48, 63], such that the trained models can have high accuracy  
97 on the dataset if and only if they can understand the human body movement. One such example is  
98 Mimetics [63] with 713 hand-collected videos of 50 human action classes from Kinetics-400 [28]  
99 happening against irrelevant backgrounds. However, scaling up or generalizing this effort would be  
100 extremely costly due to the need for manual annotation. In contrast, our toolkit provides quantitative  
101 metrics for directly measuring the effect of background bias without the need for manual annotation.

## 102 3 Human-centric Analysis Toolkit

103 Our Human-centric Analysis Toolkit (HAT) is a general framework that can be used to measure the  
104 amount of background bias learned by a human action recognition model. HAT takes two inputs: (1)  
105 a trained human action recognition model and (2) a set of validation videos each annotated with the  
106 human action class. HAT then proceeds in three steps. First, it leverages human segmentation models  
107 to separate the human visual cues from the background visual cues in the validation videos. Second,  
108 it generates six sets of counterfactual validation videos, including Human-Only, Background-Only,  
109 and four sets of Action-Swap videos (see Figure 1 for examples). Finally, it evaluates the trained  
110 model on these counterfactual videos and returns a set of ten metrics which quantify the different  
111 effects of background bias. This methodology can expand the dataset without any need to manually  
112 collect new data, allowing deeper analysis of human action recognition in an affordable manner.

### 113 3.1 Separating human from background

114 The first step of HAT is separating the visual cues corresponding directly to the *human* from the rest of  
115 the cues in the video. This can be done using a pre-trained human segmentation model. Interestingly,  
116 in our internal experiments we find that modern image-based segmentation models [24, 71, 69]  
117 tend to have better results than video-based segmentation models [40]. (We hypothesize that this  
118 might be due to the differences in training set size.) While older CNN-based image segmentation  
119 models [71, 69] suffer from low temporal consistency, missing human segments in some of the  
120 frames, the modern transformer-based SeMask [24] actually appears to have overcome this limitation.  
121 We use SeMask trained on ADE20K [74] in our implementation and include sample videos in the  
122 supplementary material.

123 One thing to note is that in the current instantiation of HAT we consider any *objects* that the human is  
124 interacting to be part of the background. Thus, for example, a person performing the “drinking coffee”  
125 action would be expected to be segmented separately from the coffee mug that they are holding  
126 (which becomes part of the background). One way of partially avoiding this would be to use a human  
127 bounding box instead of a segmentation mask – however, undesirable background cues would then  
128 also be included. Different tradeoffs can be considered in future instantiations of HAT.

### 129 3.2 Generating counterfactual validation videos

130 The core of our toolkit is generating synthetic validation videos with different visual cues, which  
131 allows us to investigate the effect of the different cues on human action recognition models.

132 The first two sets of videos are **Background-Only** (where only the background is shown and all human  
133 cues are removed) and **Human-Only** (where only the human cues are shown). For Background-Only,  
134 we leverage the video inpainting model [29] to remove all human pixels segmented by the model of  
135 Section 3.1. In contrast to prior works [6, 21] which fill the human pixels with a frame average color  
136 value (e.g., grey), we use inpainting to generate a more realistic-looking video. For Human-Only, we  
137 instead keep only the segmented human pixels and fill in the rest with an average color. We use the  
138 *dataset’s* average color rather than the *frame* average, since that can reveal a lot about the background,  
139 e.g., green for a sports field or blue for a body of water.

140 The other four sets of videos are more complex **Action-Swap** videos, which combine different visual  
141 cues to investigate their additive effects. We synthesize these videos by combining the segmented  
142 human figure with the background from a different video, similar to [64]. While the Background-Only  
143 and Human-Only video sets are both decidedly outside the model’s training data distribution, these  
144 Action-Swap videos are arguably somewhat more realistic since they do contain a human figure  
145 against a viable background – although in an unexpected combination. Example frames are in Figure 1  
146 and videos in supplementary material; more details on Action-Swap generation below.

#### 147 3.2.1 Details of generating Action-Swap videos

148 HAT includes four different types of Action-Swap videos:

- 149 • **Random**: The background is swapped with a video from a different class.
- 150 • **Close**: The background is swapped with a video from a class with a similar background.
- 151 • **Far**: The background is swapped with a video from a class with a very different background.
- 152 • **Same**: The background is swapped with a video from the same class. This can be used as a  
153 theoretical upper bound of Action-Swap Accuracy.

154 To determine the appropriate classes for **Close** and **Far** Action-Swap videos, we need to determine  
155 how similar the backgrounds are across different classes. To do so, we first feed the frames from the  
156 original validation videos into a Places365 [73] trained scene classification model. For each action  
157 class, we then compute the average scene prediction vector by averaging the prediction probabilities  
158 from all frames of all videos of this class. We can then rank all the other classes according to the L1  
159 distance in their average scene prediction vector. We consider the class to be “close” if it’s among the  
160 5 classes with the smallest L1 distance and “far” if it’s among the 200 largest (of 399 classes total).

161 For generating an Action-Swap counterfactual video, we thus:

- 162 (1) segment the human figure from the video using [24] as if creating a Human-Only video,
- 163 (2) randomly sample a background action class, depending on the particular Action-Swap set,
- 164 (3) randomly sample a video of the class from (2),
- 165 (4) generate the Background-Only version of the video from (3),
- 166 (5) paste in the human figure from (1) onto the video from (4)

167 One additional challenge is that we want to ensure that sufficient human *and* background cues are  
168 present in every generated Action-Swap video. Thus, we only consider videos where all frames have  
169 human masks taking up 5-50% of the pixels; when sampling background videos in step (3) we relax  
170 the lower bound to allow videos with few human pixels.<sup>2</sup> Therefore, unlike the Background-Only  
171 and Human-Only sets, the Action-Swap sets have fewer video samples than the original dataset. In  
172 Kinetics-400, we end up with 5,631 videos, whereas the original validation set has 19,877 videos. To  
173 compensate for this, we run steps (2-5) three times for each video to generate three different videos.

### 174 3.3 Metrics

175 We use the generated counterfactual validation videos from Section 3.2 to evaluate the trained human  
176 action recognition models. We measure how much of the original recognition accuracy comes from  
177 the different cues:

$$\text{Background-Only Ratio (BOR)} = \frac{\text{Background-Only Accuracy}}{\text{Original Accuracy}} \quad (1)$$

$$\text{Human-Only Ratio (HOR)} = \frac{\text{Human-Only Accuracy}}{\text{Original Accuracy}} \quad (2)$$

178 If a model shows high BOR, i.e., a model can get close to the original accuracy with just the  
179 background cues, we see this as “right for the wrong reason.” In contrast, ideally models would have  
180 high HOR since they should be able to recognize the human action even without the background  
181 cues.

182 Finally, for Action-Swap videos recall that each counterfactual video is generated by combining  
183 the human figure foreground from class A with the background from a different class B. We then  
184 measure the **Swap Human Accuracy (SHAcc)** as the fraction of counterfactual videos the model  
185 predicts correctly as class A, and **Swap Background Error (SBErr)** as the fraction of times the  
186 model incorrectly predicts the video as the background class B. Human action recognition models  
187 that successfully rely on human motion cues would be expected to have high SHAcc; those that are  
188 driven primarily by background cues would be expected to have high SBErr.

## 189 4 Analyzing Action Recognition Models

190 We now demonstrate the capabilities of HAT by evaluating human action recognition models trained  
191 on the popular Kinetics-400 [28] dataset. We present the results on the different types of counterfactual  
192 videos in order (Background-Only in Section 4.2, Human-Only in Section 4.3, and Action-Swap  
193 in Section 4.4), along with discussing our findings and drawing conclusions about different model  
194 design decisions that appear to have contributed to the learned background bias. HAT is not limited  
195 to Kinetics-400, and can be used on other human action recognition datasets [19, 23, 52]. Please refer  
196 to the supplementary material for the experiments on UCF101.

### 197 4.1 Experimental Details

198 We test a number of different model designs, including TSN [61], I3D [5], Non-local Neural  
199 Networks [62], R(2+1)D [54], TSM [36], SlowFast, and SlowOnly [15], CSN [58], TIN [49],

---

<sup>2</sup>Please see visualization examples here [https://github.com/princetonvisualai/HAT/blob/main/doc/review\\_discussion.md#percentage-of-synthetic-pixels](https://github.com/princetonvisualai/HAT/blob/main/doc/review_discussion.md#percentage-of-synthetic-pixels)

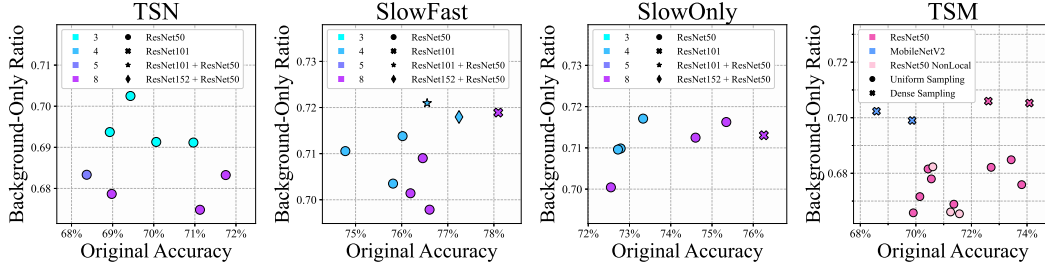


Figure 2: **1-3**: We visualize how the number of frames used for training can worsen the Background-Only Ratio. **4**: Background-Only Ratio on sampling strategy. A single dot represents a single trained weight under a specific training configuration.

200 TPN [68], X3D [14], OmniSource [12], TANet [39], and TimeSformer [4]. In total, we test 74  
 201 different trained models offered by the MMAAction2 [8] implementation.

202 We extract the videos in 30 FPS with original resolution. For other pre-processing, such as resizing  
 203 and temporal sampling, we follow the configuration that each model specified. We list the details  
 204 of the tested models and their configuration in the supplementary material. Within the scope of the  
 205 paper, we chose not to retrain any models and rely on publicly released model weights. In drawing  
 206 conclusions we try to do an apples-to-apples comparison whenever possible; however, we are not able  
 207 to guarantee that all hyperparameter settings are directly comparable between the different models.

208 For image segmentation and video inpainting, we used 20 Nvidia RTX 3090 GPUs with 20 hours of  
 209 forward pass to generate synthetic videos of the full Kinetics-400 validation set. See supplementary  
 210 material for examples of the synthetic videos on Kinetics-400.

## 211 4.2 Background-Only Video

212 **Accuracy and Background-Only Ratio.** Table 1 tabulates the model accuracies in Background-  
 213 Only Videos. For a fair comparison, we have separated the weights that use additional large-scale  
 214 datasets [10, 12, 16, 26, 65]. Despite removing a human body from the video, thus removing any  
 215 human action, all the models still show a strong tendency to predict the removed action. This hints at  
 216 the possibility that the performance of the human action recognition models is highly dependent on  
 217 the background, rather than the action itself.

218 Table 1 tabulates the Background-Only Ratio. It shows that on all the tested models, we see around  
 219 70% of the accuracy is coming from the non-human regions, revealing the problematic behavior,  
 220 “right for the wrong reason”, is common in human action recognition. Next, we show examples of  
 221 using Background-Only Ratio to analyze and improve the model design.

222 **Number of Frames used to Train.** The first three plots of Figure 2 visualize how the number of  
 223 video frames used during the training can worsen the Background-Only Ratio. This shows that the  
 224 models trained with fewer temporal frames tend to suffer more, with a lot of their accuracy coming  
 225 from the background. A possible explanation is that when fewer frames are given, the model is  
 226 not able to learn to understand temporal information, thus given a video with or without the human  
 227 movement, the model will perform similarly, as they never learned to understand such complex human  
 228 movement during training. Thus the accuracy would come from the temporally static background.  
 229 While exact behavior can be different per model structure, we see this to be most severe on TSN [61]  
 230 which lacks any sophisticated temporal modeling.

231 **Sampling Strategy.** We check if the frame sampling strategy can affect the Background-Only Ratio.  
 232 The results are visualized on the last plot of Figure 2. Unlike uniform sampling, i.e., getting uniformly  
 233 distributed frames, dense sampling strategy, i.e., sampling frames with a specified stride, shows higher  
 234 BOR in general. We believe this is due to the dense sampling strategy having a smaller temporal

Table 1: Accuracies over Background-Only Videos. OAcc and BAcc denote original accuracy and accuracy on Background-Only Videos, respectively. Models using additional large-scale data are tabulated separately below. We only tabulate the setting with the highest OAcc per backbone. Check supplementary material for the experiment results of all 74 weights.

Model	Backbone	Pre-trained	OAcc (%)	BAcc (%)	$\frac{\text{BAcc}}{\text{OAcc}}$
<i>Normal-scale dataset</i>					
TSM [36]	MNetV2 [44]	ImageNet	69.87	48.84	0.6990
R(2+1)D [54]	ResNet34	-	74.22	52.99	0.7140
TSN [61]	ResNet50	ImageNet	71.75	49.02	0.6833
TIN [49]	ResNet50	TSM-Kinetics400	70.89	48.32	0.6816
TSM [36]	ResNet50	ImageNet	74.09	52.25	0.7053
I3D [5]	ResNet50	ImageNet	73.57	52.26	0.7104
NL-TSM [62]	ResNet50	ImageNet	71.57	47.62	0.6654
NL-I3D [62]	ResNet50	ImageNet	74.91	52.84	0.7054
NL-SlowOnly [62]	ResNet50	ImageNet	75.78	53.51	0.7062
CSN [58]	ResNet50	-	73.22	51.97	0.7098
TPN [68]	ResNet50	ImageNet	76.16	54.40	0.7143
SlowOnly [15]	ResNet50	ImageNet	75.35	53.97	0.7163
SlowFast [15]	ResNet50	-	76.61	53.46	0.6978
SlowOnly [15]	ResNet101	-	76.26	54.38	0.7131
SlowFast [15]	ResNet101+50	-	76.55	55.19	0.7210
SlowFast [15]	ResNet101	-	<b>78.10</b>	56.14	0.7189
CSN [58]	ResNet152	-	77.62	54.33	0.6999
SlowFast [15]	ResNet152+50	-	77.24	55.46	0.7179
X3D [14]	X3D_S	-	72.67	50.61	0.6964
X3D [14]	X3D_M	-	75.55	52.47	0.6944
TANet [39]	TANet	ImageNet	76.10	53.71	0.7059
<i>Large-scale dataset</i>					
TSN [61]	ResNet50	IG-1B [65]	70.96	49.05	0.6912
Omni-TSN [12]	ResNet50	IG-1B [65]	74.70	52.09	0.6973
Omni-SlowOnly [12]	ResNet50	-	76.49	55.00	0.7190
CSN [58]	ResNet50	IG65M [16]	79.09	55.83	0.7059
Omni-SlowOnly [12]	ResNet101	-	80.00	58.05	0.7255
CSN [58]	ResNet152	IG65M [16]	<b>82.38</b>	58.97	0.7159
TimeSFormer [4]	TimeSformer	ImageNet-21K [10]	77.97	53.88	0.6910

235 window so that the model was not able to learn the body movement sufficiently. Surprisingly, the  
 236 effect of the sampling strategy would have not been clear if we only used original accuracy alone (see  
 237 x-axis), showing a clear benefit of using BOR for model training analysis.

### 238 4.3 Human-Only Video

239 **Accuracy and Human-Only Ratio.** We plot HOR in left of Figure 3. We tabulated the evaluation  
 240 results in the supplementary material. Given only the human action, all the models suffer significantly  
 241 with an accuracy of around 20%. Despite Human-Only modification keeping the human action intact,  
 242 the ratio is far lower than Background-Only Accuracy. By comparing BOR (with around 0.7) and  
 243 HOR (with around 0.3), we quantitatively measure the well-believed problem of the current state of  
 244 human action recognition, that most existing methods are all highly influenced by the background,  
 245 more than the foreground human action.

246 Thankfully, we see a strong correlation between Human-Only Ratio and the original accuracy. This  
 247 could hint that the performance improvement of the action recognition model is benefited from a  
 248 better understanding of the human body, showing the important direction of where the human action  
 249 recognition field needs to focus. Next, we show one example case where HOR can be used to evaluate  
 250 different model structures.

251 **TSN vs TSM.** While the original paper on TSM [36] claims +4% accuracy improvements over  
 252 TSN [61] on Kinetics-400, using different training and testing conditions, MMAAction2 [61] shows

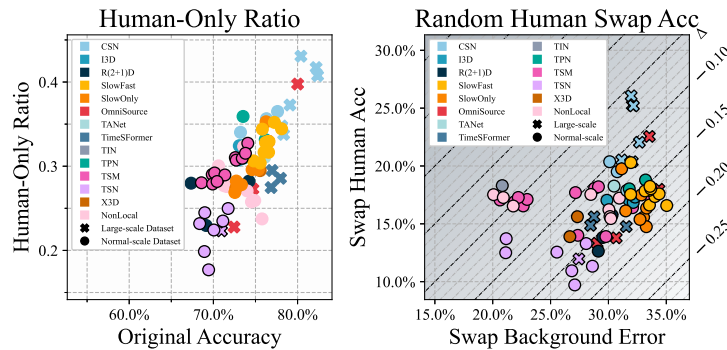


Figure 3: **Left:** Human-Only Ratio. Dots from TSM and TSN are bordered. **Right:** Action-Swap results in random swap.

253 that the accuracy of TSN can be achieved on par with TSM, as shown in the x-axis of left of Figure 3.  
 254 However, using Human-Only Ratio as a metric, we show that TSM does indeed show superior  
 255 performance over TSN when a non-human region is removed. One possible explanation is that,  
 256 as TSM design makes use of temporal difference, e.g., human body movement, it can capture the  
 257 information of the human body better, as TSN cannot distinguish between human and background  
 258 using its basic temporal modeling design.

#### 259 4.4 Action-Swap Videos

260 **Accuracy on Action-Swap.** Table 2 and right of Figure 3 details the performance of different  
 261 models over the Action-Swap Videos. It shows that when we randomly swap background with other  
 262 videos, all the models lean towards predicting the class of the background, rather than the foreground  
 263 human action. Swapping between classes that are similar/different shows a gain/drop in SHAcc,  
 264 showing that the output of a human action recognition model is largely dependent on the background.

265 **Original Accuracy vs. Action-Swap Accuracy.** Among models using normal-scale datasets,  
 266 SlowFast-Res101 [15] shows the best accuracy when the background is relevant to the foreground  
 267 action, on both original Kinetics accuracy (See Tab. 1) and Same Swap (See Tab. 2). However, given  
 268 counterfactual videos that have irrelevant backgrounds, their performance drops to 18%, while the  
 269 model falsely predicts 34 percent of the validation videos as their background class, one of the highest  
 270 among the models we have tested. Such low performance on human action could be due to its reliance  
 271 on the background, as models with better Random Swap SHAcc (CSN-Res152, SlowFast-Res152+50,  
 272 etc.) show fewer background errors. Such experiment shows that models showing good accuracy in  
 273 original Kinetics-400, might not be a good human action recognition model, due to their reliance on  
 274 the background.

275 **Use of Non-local Module.** To demonstrate the evaluation of a model design using Action-Swap,  
 276 we select Non-local [62] module as an example. Table 3 tabulates the evaluation results on Random  
 277 Swap. We see that the Non-local module not only improves the original accuracy, but also drops the  
 278 background error on all the tested models, showing reduced background bias. However, Non-local  
 279 module do not always improve the focus on the human body, as for I3D [5] models, we see that  
 280 SHAcc tends to drop.

281 **Use of Large-scale Dataset for Pre-training.** Table 4 tabulates the performance of models where  
 282 we compare trained weight with/without additional large-scale pre-training. It shows that in all the  
 283 cases, using a large-scale dataset improves the original accuracy and Random Swap Human Accuracy.  
 284 However, as CSN shows an increase in the Background Error, this does not necessarily mean that  
 285 the model is being better at recognizing the human. We expect the model is recognizing the image  
 286 feature better when pre-trained with large-scale dataset, regardless of the scene or the person.



Table 2: Action-Swap experiment results. We average the numbers from 3 random runs. We show standard deviation as well. See supplementary material for the full experiments.

Model	Backbone	Pre-trained	Same			Random Swap		Close		Far	
			SHAcc $\uparrow$	SHAcc $\uparrow$	SBErr $\downarrow$	SHAcc $\uparrow$	SBErr $\downarrow$	SHAcc $\uparrow$	SBErr $\downarrow$	SHAcc $\uparrow$	SBErr $\downarrow$
<i>Normal-scale dataset</i>											
TSM [36]	MNetV2	ImgNet	62.2 $\pm$ .3	13.9 $\pm$ .1	29.8 $\pm$ .2	24.4 $\pm$ .3	26.4 $\pm$ .4	11.2 $\pm$ .1	35.5 $\pm$ .2		
R(2+1)D [54]	Res34	-	64.5 $\pm$ .3	15.8 $\pm$ .3	30.3 $\pm$ .5	26.6 $\pm$ .3	27.1 $\pm$ .4	13.0 $\pm$ .1	35.6 $\pm$ .1		
TSN [61]	Res50	ImgNet	60.2 $\pm$ .2	13.3 $\pm$ .1	28.1 $\pm$ .2	23.4 $\pm$ .2	26.7 $\pm$ .3	11.9 $\pm$ .1	32.7 $\pm$ .2		
TIN [49]	Res50	Kin400	58.6 $\pm$ .1	18.3 $\pm$ .2	<b>20.8</b> $\pm$ .1	27.1 $\pm$ .2	<b>21.0</b> $\pm$ .2	16.6 $\pm$ .1	<b>23.5</b> $\pm$ .3		
TSM [36]	Res50	ImgNet	66.6 $\pm$ .4	17.2 $\pm$ .5	33.7 $\pm$ .5	27.8 $\pm$ .1	29.2 $\pm$ .2	14.3 $\pm$ .3	40.4 $\pm$ .2		
I3D [5]	Res50	ImgNet	64.9 $\pm$ .4	17.0 $\pm$ .2	29.9 $\pm$ .1	27.4 $\pm$ .3	26.6 $\pm$ .5	14.8 $\pm$ .5	34.8 $\pm$ .5		
NL-TSM [62]	Res50	ImgNet	58.6 $\pm$ .4	16.5 $\pm$ .2	21.8 $\pm$ .2	25.9 $\pm$ .6	21.7 $\pm$ .2	15.0 $\pm$ .3	25.0 $\pm$ .1		
NL-I3D [62]	Res50	ImgNet	64.9 $\pm$ .4	16.2 $\pm$ .2	30.0 $\pm$ .4	27.0 $\pm$ .2	26.6 $\pm$ .1	13.4 $\pm$ .3	35.6 $\pm$ .4		
NL-SlowOnly [62]	Res50	ImgNet	63.8 $\pm$ .1	17.5 $\pm$ .2	28.5 $\pm$ .5	27.0 $\pm$ .2	25.6 $\pm$ .3	14.8 $\pm$ .5	34.0 $\pm$ .4		
CSN [58]	Res50	-	65.9 $\pm$ .3	17.9 $\pm$ .2	31.6 $\pm$ .2	28.2 $\pm$ .2	27.6 $\pm$ .5	15.2 $\pm$ .2	37.1 $\pm$ .5		
TPN [68]	Res50	ImgNet	69.3 $\pm$ .2	18.8 $\pm$ .2	33.2 $\pm$ .5	29.0 $\pm$ .3	29.0 $\pm$ .6	15.8 $\pm$ .2	38.9 $\pm$ .4		
SlowOnly [15]	Res50	ImgNet	68.2 $\pm$ .1	17.5 $\pm$ .2	32.8 $\pm$ .5	28.1 $\pm$ .4	28.7 $\pm$ .2	14.8 $\pm$ .3	38.8 $\pm$ .4		
SlowFast [15]	Res50	-	68.4 $\pm$ .3	18.0 $\pm$ .3	33.7 $\pm$ .6	28.8 $\pm$ .2	29.7 $\pm$ .5	15.0 $\pm$ .2	40.0 $\pm$ .2		
SlowOnly [15]	Res101	-	69.4 $\pm$ .4	19.8 $\pm$ .3	31.1 $\pm$ .6	<b>31.0</b> $\pm$ .1	28.1 $\pm$ .3	17.0 $\pm$ .2	37.0 $\pm$ .4		
SlowFast [15]	Res101+50	-	67.9 $\pm$ .2	17.5 $\pm$ .3	31.9 $\pm$ .4	28.4 $\pm$ .3	29.0 $\pm$ .2	15.1 $\pm$ .1	37.7 $\pm$ .6		
SlowFast [15]	Res101	-	<b>69.6</b> $\pm$ .3	18.2 $\pm$ .2	33.6 $\pm$ .6	29.2 $\pm$ .3	29.4 $\pm$ .3	15.4 $\pm$ .1	40.0 $\pm$ .5		
CSN [58]	Res152	-	67.8 $\pm$ .4	<b>20.4</b> $\pm$ .5	30.1 $\pm$ .3	30.8 $\pm$ .2	26.3 $\pm$ .3	<b>17.6</b> $\pm$ .3	35.2 $\pm$ .0		
SlowFast [15]	Res152+50	-	69.3 $\pm$ .5	20.3 $\pm$ .6	31.9 $\pm$ .7	31.0 $\pm$ .1	28.5 $\pm$ .3	17.5 $\pm$ .2	36.9 $\pm$ .2		
X3D [14]	X3D_S	-	60.8 $\pm$ .3	13.9 $\pm$ .3	26.7 $\pm$ .7	24.2 $\pm$ .2	24.7 $\pm$ .3	11.0 $\pm$ .1	32.0 $\pm$ .3		
X3D [14]	X3D_M	-	64.3 $\pm$ .3	15.6 $\pm$ .2	27.3 $\pm$ .1	26.5 $\pm$ .4	25.5 $\pm$ .1	12.8 $\pm$ .0	32.8 $\pm$ .6		
TANet [39]	TANet	ImgNet	67.1 $\pm$ .3	18.3 $\pm$ .3	30.5 $\pm$ .4	28.5 $\pm$ .2	27.0 $\pm$ .3	15.5 $\pm$ .1	36.6 $\pm$ .4		
<i>Large-scale dataset</i>											
TSN [61]	Res50	IG-1B [65]	57.7 $\pm$ .5	12.0 $\pm$ .3	<b>27.4</b> $\pm$ .3	21.4 $\pm$ .3	<b>25.7</b> $\pm$ .1	10.1 $\pm$ .3	<b>32.1</b> $\pm$ .2		
Omni-TSN [12]	Res50	IG-1B [65]	63.9 $\pm$ .6	13.8 $\pm$ .4	30.7 $\pm$ .1	24.4 $\pm$ .1	27.9 $\pm$ .2	11.8 $\pm$ .2	36.8 $\pm$ .6		
Omni-Slow [12]	Res50	-	69.5 $\pm$ .3	18.0 $\pm$ .6	34.4 $\pm$ .5	29.1 $\pm$ .2	29.8 $\pm$ .2	15.0 $\pm$ .2	40.8 $\pm$ .2		
CSN [58]	Res50	IG65M [16]	70.4 $\pm$ .3	22.1 $\pm$ .5	32.7 $\pm$ .2	32.4 $\pm$ .4	28.9 $\pm$ .4	18.8 $\pm$ .1	37.7 $\pm$ .2		
TSFormer [4]	TSformer	Img21K [10]	65.3 $\pm$ .3	15.6 $\pm$ .3	28.8 $\pm$ .1	25.8 $\pm$ .1	27.4 $\pm$ .3	13.0 $\pm$ .3	33.2 $\pm$ .5		
Omni-Slow [12]	Res101	-	<b>73.3</b> $\pm$ .4	22.6 $\pm$ .2	33.5 $\pm$ .4	33.4 $\pm$ .5	30.1 $\pm$ .5	19.4 $\pm$ .2	39.2 $\pm$ .3		
CSN [58]	Res152	IG65M [16]	72.9 $\pm$ .1	<b>25.2</b> $\pm$ .4	32.2 $\pm$ .5	<b>35.6</b> $\pm$ .3	28.4 $\pm$ .6	<b>22.1</b> $\pm$ .3	38.0 $\pm$ .3		

## 287 5 Conclusion

288 We introduce a general framework for human-centric analysis for human action recognition models.  
 289 We test Human-centric Analysis Toolkit on the Kinetics-400 dataset and evaluate the generated  
 290 dataset on a number of existing action recognition models.

291 Through extensive experiments over 74 trained models, we find that all the models we tested have  
 292 stronger background bias. However, we found that the background bias can be mitigated when  
 293 more frames are fed during the training, the temporal stride between frames is increased, and  
 294 temporal/spacial modeling is improved using Non-local module. Moreover, we see that the original  
 295 accuracy do not fully represent the human understanding as the accuracy cannot differentiate TSN  
 296 and TSM, large-scale dataset and Non-local module improves original accuracy but not necessarily  
 297 SHAcc.

298 From our findings, we suggest the future researchers to (1) not rely on the accuracy as the only metric,  
 299 as original accuracy do not fully represent the performance of the model based on the human action;  
 300 (2) carefully select the temporal hyper-parameters, as temporal parameters can improve/worsen the  
 301 background bias of human action recognition models; and (3) use HAT toolkit to see if the model  
 302 design (e.g., as Non-local) can improve your model on accuracy and reduce the background bias. We  
 303 hope that this tool can be adopted by future researchers for a better human-centric analysis of human  
 304 action recognition models.

Table 3: Performance comparison when using a Non-local module [62]. NL-EG, NL-G, and NL-Dot denote Non-local method using embedded Gaussian, Gaussian, and dot product, respectively. Numbers are bolded when the Non-local module improves the metric.

Model	frames	OAcc $\uparrow$	SHAcc $\uparrow$	SBErr $\downarrow$
TSM	8	72.89	16.55	22.64
TSM + NL-EG	8	<b>74.06</b> <sub>(+1.18)</sub>	16.54 <sub>(-0.01)</sub>	<b>21.77</b> <sub>(-0.86)</sub>
TSM + NL-G	8	72.61 <sub>(-0.27)</sub>	<b>17.52</b> <sub>(+0.98)</sub>	<b>20.07</b> <sub>(-2.56)</sub>
TSM + NL-Dot	8	<b>73.52</b> <sub>(+0.63)</sub>	<b>17.27</b> <sub>(+0.73)</sub>	<b>20.91</b> <sub>(-1.72)</sub>
I3D	32	75.33	17.05	31.78
I3D + NL-EG	32	<b>76.90</b> <sub>(+1.58)</sub>	16.23 <sub>(-0.83)</sub>	<b>30.04</b> <sub>(-1.73)</sub>
I3D + NL-G	32	<b>75.96</b> <sub>(+0.63)</sub>	<b>17.22</b> <sub>(+0.17)</sub>	<b>30.94</b> <sub>(-0.83)</sub>
I3D + NL-Dot	32	<b>76.17</b> <sub>(+0.84)</sub>	15.63 <sub>(-1.43)</sub>	<b>30.14</b> <sub>(-1.64)</sub>
SlowOnly	4	75.28	14.75	33.27
SlowOnly + NL-EG	4	<b>76.10</b> <sub>(+0.82)</sub>	<b>15.46</b> <sub>(+0.70)</sub>	<b>30.21</b> <sub>(-3.07)</sub>
SlowOnly	8	75.18	16.12	31.49
SlowOnly + NL-EG	8	<b>77.74</b> <sub>(+2.56)</sub>	<b>17.54</b> <sub>(+1.42)</sub>	<b>28.48</b> <sub>(-3.01)</sub>

Table 4: Performance when using a large-scale dataset. We compare the same settings except for the initial weight. Numbers are bolded when the large-scale dataset improves the metric.

Model	Backbone	Pre-trained	OAcc $\uparrow$	SHAcc $\uparrow$	SBErr $\downarrow$
TSN [61]	ResNet50	ImageNet	72.55	11.34 $\pm$ 0.15	28.62 $\pm$ 0.16
TSN [61]	ResNet50	IG-1B	<b>73.39</b>	<b>11.96</b> $\pm$ 0.35	<b>27.45</b> $\pm$ 0.28
ir-CSN [58]	ResNet50	None	75.51	17.88 $\pm$ 0.18	31.58 $\pm$ 0.17
ir-CSN [58]	ResNet50	IG65M	<b>81.46</b>	<b>22.05</b> $\pm$ 0.49	32.68 $\pm$ 0.17
ir-CSN [58]	ResNet152	None	78.08	19.51 $\pm$ 0.11	30.76 $\pm$ 0.23
ir-CSN [58]	ResNet152	Sports1M	<b>78.98</b>	<b>20.52</b> $\pm$ 0.21	31.14 $\pm$ 0.51
ir-CSN [58]	ResNet152	IG65M	<b>83.17</b>	<b>25.25</b> $\pm$ 0.36	32.07 $\pm$ 0.39
ip-CSN [58]	ResNet152	None	79.26	20.37 $\pm$ 0.50	30.11 $\pm$ 0.34
ip-CSN [58]	ResNet152	Sports1M	<b>79.38</b>	20.37 $\pm$ 0.36	32.06 $\pm$ 0.31
ip-CSN [58]	ResNet152	IG65M	<b>83.92</b>	<b>25.19</b> $\pm$ 0.41	32.16 $\pm$ 0.46

## 305 6 Discussion

306 **Limitation** As we use an off-the-shelf image semantic segmentation model and a video inpainting  
 307 model, the quality of the synthetic dataset is limited by the performance of the aforementioned models.

308 **Ethical Concerns** Our tool requires the use of image segmentation and inpainting tool to generate a  
 309 dataset, requiring computation cost for the initial setup. However, as human-centric analysis using  
 310 our tool does not require any new training, we believe our tool is more environmentally friendly than  
 311 the existing methods. Moreover, as our tool is automated, human labor for data collection is not  
 312 required. Also, as we generate a dataset from an existing dataset, we show fewer concerns about  
 313 privacy issues when a new video dataset is generated.

314 **License** MMAAction2 [8] and SeMask [24] follow Apache License 2.0. We used author-released code  
 315 for Deep Video Inpainting [29] which did not specify any license. Kinetics-400 annotation data is  
 316 licensed under a Creative Commons Attribution 4.0 International License, but some of the video  
 317 sources do not specify any license. Please refer to the individual licenses when using our released  
 318 code.

319 **Acknowledgements.** We are grateful for the support from the National Science Foundation under  
 320 Grant No. 2112562, Microsoft, Princeton SEAS Project X Innovation Fund, and Princeton First Year  
 321 Ph.D. Fellowship to JC.

## 322 References

- 323 [1] Humam Alwassel, Fabian Caba, Victor Escorcia, and Bernard Ghanem. Diagnosing error in temporal  
324 action detectors. In *ECCV*, 2018.
- 325 [2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A  
326 video vision transformer. In *ICCV*, 2021.
- 327 [3] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantify-  
328 ing interpretability of deep visual representations. In *CVPR*, 2017.
- 329 [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video  
330 understanding? In *ICML*, 2021.
- 331 [5] J. Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset.  
332 In *CVPR*, 2017.
- 333 [6] Jinwoo Choi, Chen Gao, C. E. Joseph Messou, and Jia-Bin Huang. Why can't i dance in the mall? learning  
334 to mitigate scene bias in action recognition. In *NeurIPS*, 2019.
- 335 [7] Jihoon Chung, Cheng hsin Wu, Hsuan ru Yang, Yu-Wing Tai, and Chi-Keung Tang. Haa500: Human-  
336 centric atomic action dataset with curated videos. In *ICCV*, 2021.
- 337 [8] MMAction2 Contributors. Openmmlab's next generation video understanding toolbox and benchmark.  
338 <https://github.com/open-mmlab/mmdetection>, 2020.
- 339 [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos  
340 Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision:  
341 The epic-kitchens dataset. In *ECCV*, 2018.
- 342 [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
343 image database. In *CVPR*, 2009.
- 344 [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
345 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth  
346 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- 347 [12] Haodong Duan, Yue Zhao, Yuanjun Xiong, Wentao Liu, and Dahua Lin. Omni-sourced webly-supervised  
348 learning for video recognition. In *ECCV*, 2020.
- 349 [13] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for  
350 instance detection. In *ICCV*, 2017.
- 351 [14] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *CVPR*, 2020.
- 352 [15] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video  
353 recognition. In *ICCV*, 2019.
- 354 [16] Deepti Ghadiyaram, Du Tran, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video  
355 action recognition. In *CVPR*, 2019.
- 356 [17] Rohit Girdhar and Deva Ramanan. CATER: A diagnostic dataset for Compositional Actions and TEmporal  
357 Reasoning. In *ICLR*, 2020.
- 358 [18] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal,  
359 Heuna Kim, Valentin Haefliger, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something  
360 something" video database for learning and evaluating visual common sense. In *ICCV*, 2017.
- 361 [19] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra  
362 Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of  
363 spatio-temporally localized atomic visual actions. In *CVPR*, 2018.
- 364 [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
365 In *CVPR*, 2016.
- 366 [21] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also  
367 snowboard: Overcoming bias in captioning models. In *ECCV*, 2018.
- 368 [22] Derek Hoiem, Yodsawalai Chodpathumwan, and Qieyun Dai. Diagnosing error in object detectors. In  
369 *ECCV*, 2012.
- 370 [23] Oana Ignat, Laura Burdick, Jia Deng, and Rada Mihalcea. Identifying visible actions in lifestyle vlogs. In  
371 *ACL*, 2019.
- 372 [24] Jitesh Jain, Anukriti Singh, Nikita Orlov, Zilong Huang, Jiachen Li, Steven Walton, and Humphrey Shi.  
373 Semask: Semantically masking transformer backbones for effective semantic segmentation. *arXiv*, 2021.
- 374 [25] Carlos E. Jimenez, Olga Russakovsky, and Karthik Narasimhan. Carets: A consistency and robustness  
375 evaluative test suite for vqa. In *ACL*, 2022.
- 376 [26] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei.  
377 Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- 378 [27] Kevin Karsch, Varsha Hedau, David Forsyth, and Derek Hoiem. Rendering synthetic objects into legacy  
379 photographs. *ACMTOG*, 2011.
- 380 [28] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,  
381 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The  
382 kinetics human action video dataset, 2017.

- 383 [29] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Deep video inpainting. In *CVPR*, 2019.
- 384 [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional  
385 neural networks. In *NeurIPS*, 2012.
- 386 [31] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human  
387 motion recognition. In *ICCV*, 2011.
- 388 [32] Feng Li, Hao Zhang, Huaizhe xu, Shilong Liu, Lei Zhang, Lionel M. Ni, and Heung-Yeung Shum. Mask  
389 dino: Towards a unified transformer-based framework for object detection and segmentation. In *CVPR*,  
390 2022.
- 391 [33] Rui Li, Yiheng Zhang, Zhaofan Qiu, Ting Yao, Dong Liu, and Tao Mei. Motion-focused contrastive  
392 learning of video representations. In *ICCV*, 2021.
- 393 [34] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation  
394 bias. In *ECCV*, 2018.
- 395 [35] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng. Towards an end-to-end  
396 framework for flow-guided video inpainting. In *CVPR*, 2022.
- 397 [36] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In  
398 *ICCV*, 2019.
- 399 [37] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai,  
400 and Hongsheng Li. Fuseformer: Fusing fine-grained information in transformers for video inpainting. In  
401 *ICCV*, 2021.
- 402 [38] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin  
403 transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- 404 [39] Zhaoyang Liu, Limin Wang, Wayne Wu, Chen Qian, and Tong Lu. Tam: Temporal adaptive module for  
405 video recognition. In *ICCV*, 2021.
- 406 [40] Jiaxu Miao, Yunchao Wei, Yu Wu, Chen Liang, Guangrui Li, and Yi Yang. Vspw: A large-scale dataset for  
407 video scene parsing in the wild. In *CVPR*, 2021.
- 408 [41] Yair Movshovitz-Attias, Takeo Kanade, and Yaser Sheikh. How useful is photo-realistic rendering for  
409 visual learning? In *ECCV*, 2016.
- 410 [42] Amir Rosenfeld, Richard Zemel, and John K Tsotsos. The elephant in the room. *arXiv preprint*  
411 *arXiv:1808.03305*, 2018.
- 412 [43] Olga Russakovsky, Jia Deng, Zhiheng Huang, Alexander C Berg, and Li Fei-Fei. Detecting avocados to  
413 zucchinis: what have we done, and where are we going? In *ICCV*, 2013.
- 414 [44] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2:  
415 Inverted residuals and linear bottlenecks. In *CVPR*, 2018.
- 416 [45] Christian Schuldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: a local svm approach. In  
417 *ICCV*, 2004.
- 418 [46] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and  
419 Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*,  
420 2017.
- 421 [47] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+ d: A large scale dataset for 3d  
422 human activity analysis. In *CVPR*, 2016.
- 423 [48] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. Finegym: A hierarchical video dataset for fine-grained  
424 action understanding. In *CVPR*, 2020.
- 425 [49] Hao Shao, Shengju Qian, and Yu Liu. Temporal interlacing network. In *AAAI*, 2020.
- 426 [50] Gunnar Sigurdsson, Olga Russakovsky, and Abhinav Gupta. What actions are needed for understanding  
427 human actions in videos? In *ICCV*, 2017.
- 428 [51] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-  
429 tion. In *ICLR*, 2015.
- 430 [52] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions  
431 classes from videos in the wild. *arXiv*, 2012.
- 432 [53] Hao Su, Charles R Qi, Yangyan Li, and Leonidas J Guibas. Render for cnn: Viewpoint estimation in  
433 images using cnns trained with rendered 3d model views. In *ICCV*, 2015.
- 434 [54] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at  
435 spatiotemporal convolutions for action recognition. In *CVPR*, 2018.
- 436 [55] Gül Varol, Ivan Laptev, Cordelia Schmid, and Andrew Zisserman. Synthetic humans for action recognition  
437 from unseen viewpoints. *IJCV*, 2021.
- 438 [56] Gul Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J Black, Ivan Laptev, and Cordelia  
439 Schmid. Learning from synthetic humans. In *CVPR*, 2017.
- 440 [57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
441 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 442 [58] Heng Wang, Matt Feiszli, and Lorenzo Torresani. Video classification with channel-separated convolutional  
443 networks. In *ICCV*, 2019.

- 444 [59] Jinpeng Wang, Yuting Gao, Ke Li, Jianguo Hu, Xinyang Jiang, Xiaowei Guo, Rongrong Ji, and Xing Sun.  
 445 Enhancing unsupervised video representation learning by decoupling the scene and the motion. In *AAAI*,  
 446 2021.
- 447 [60] Jinpeng Wang, Yuting Gao, Ke Li, Yiqi Lin, Andy J Ma, Hao Cheng, Pai Peng, Feiyue Huang, Rongrong  
 448 Ji, and Xing Sun. Removing the background by adding the background: Towards background robust  
 449 self-supervised video representation learning. In *CVPR*, 2021.
- 450 [61] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal  
 451 segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016.
- 452 [62] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*,  
 453 2018.
- 454 [63] Philippe Weinzaepfel and Grégory Rogez. Mimetics: Towards understanding human actions out of context.  
 455 *arXiv*, 2019.
- 456 [64] Yu Wu and Yi Yang. Exploring heterogeneous clues for weakly-supervised audio-visual video parsing. In  
 457 *CVPR*, 2021.
- 458 [65] I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-supervised  
 459 learning for image classification. *arXiv*, 2019.
- 460 [66] Junji Yamato, Jun Ohya, and Kenichiro Ishii. Recognizing human action in time-sequential images using  
 461 hidden markov model. In *CVPR*, 1992.
- 462 [67] Shen Yan, Xuehan Xiong, Anurag Arnab, Zhichao Lu, Mi Zhang, Chen Sun, and Cordelia Schmid.  
 463 Multiview transformers for video recognition. In *CVPR*, 2022.
- 464 [68] Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network for action  
 465 recognition. In *CVPR*, 2020.
- 466 [69] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation.  
 467 In *ECCV*, 2020.
- 468 [70] Manlin Zhang, Jinpeng Wang, and Andy J Ma. Suppressing static visual cues via normalizing flows for  
 469 self-supervised video representation learning. In *AAAI*, 2022.
- 470 [71] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing  
 471 network. In *CVPR*, 2017.
- 472 [72] Hang Zhao, Zhicheng Yan, Lorenzo Torresani, and Antonio Torralba. Hacs: Human action clips and  
 473 segments dataset for recognition and temporal localization. *arXiv*, 2019.
- 474 [73] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million  
 475 image database for scene recognition. In *TPAMI*, 2017.
- 476 [74] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing  
 477 through ADE20k dataset. In *CVPR*, 2017.

## 478 Checklist

- 479 1. For all authors...
- 480 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 481 contributions and scope? [\[Yes\]](#)
- 482 (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 6
- 483 (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See  
 484 Section 6
- 485 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 486 them? [\[Yes\]](#)
- 487 2. If you are including theoretical results...
- 488 (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
- 489 (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
- 490 3. If you ran experiments (e.g. for benchmarks)...
- 491 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 492 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
- 493 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 494 were chosen)? [\[Yes\]](#) See Section 4.1
- 495 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 496 ments multiple times)? [\[Yes\]](#) See appendix.
- 497 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 498 of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Section 4.1
- 499 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 500 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 6

- 501 (b) Did you mention the license of the assets? [Yes] See Section 6
- 502 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 503 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 504 using/curating? [N/A]
- 505 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 506 information or offensive content? [N/A]
- 507 5. If you used crowdsourcing or conducted research with human subjects...
- 508 (a) Did you include the full text of instructions given to participants and screenshots, if
- 509 applicable? [N/A]
- 510 (b) Did you describe any potential participant risks, with links to Institutional Review
- 511 Board (IRB) approvals, if applicable? [N/A]
- 512 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 513 spent on participant compensation? [N/A]