

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053054
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

Audio-Adaptive Activity Recognition Across Video Domains

Anonymous CVPR submission

Paper ID 62

Abstract

This paper strives for activity recognition under domain shift, for example caused by change of scenery or camera viewpoint. The leading approaches reduce the shift in activity appearance by adversarial training and self-supervised learning. Different from these vision-focused works we leverage activity sounds for domain adaptation as they have less variance across domains and can reliably indicate which activities are not happening. We propose an audio-adaptive encoder and associated learning methods that discriminatively adjust the visual feature representation as well as addressing shifts in the semantic distribution. To further eliminate domain-specific features and include domain-invariant activity sounds for recognition, an audio-infused recognizer is proposed, which effectively models the cross-modal interaction across domains. We also introduce the new task of actor shift, with a corresponding audio-visual dataset, to challenge our method with situations where the activity appearance changes dramatically. Experiments on this dataset, EPIC-Kitchens and CharadesEgo show the effectiveness of our approach. For instance, we achieve a 5% absolute improvement over previous works in EPIC-Kitchens. This paper has been accepted at CVPR2022.

1. Introduction

The goal of this paper is to recognize activities such as *eating*, *sleeping* or *cutting* under domain shift caused by change of scenery, camera viewpoint or actor, as shown in Figure 1. Existing solutions align distribution-shifted domains inside a single visual video network by adversarial training [5, 20, 27, 29] and self-supervised learning [9, 22, 34]. Although successful, projecting the visual features from different source and target domains into a shared space can make the ability of the model to distinguish between classes in the target domain suffer. We observe that activity sounds can act as natural domain-invariant cues, as they carry rich activity information while exhibiting less variance across domains. We thus propose a video model which adapts to video distribution shifts with the aid of sound.

Many have considered sound in addition to visual analysis for activity recognition within a single domain [18, 24, 25,

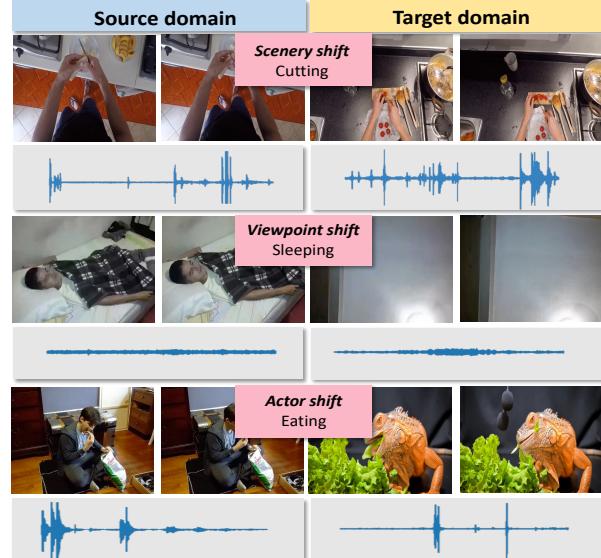


Figure 1. We recognize activities under domain shifts, caused by change of scenery, camera viewpoint or actor, with the aid of sound.

28, 31, 38, 39, 45, 46, 49]. For instance, both Gao *et al.* [18] and Korbar *et al.* [24] reduce the computational cost by previewing the audio track, while Lee *et al.* [25] show that combining visual features with audio can better localize actions. However, the cross-modal correspondences become harder to discover when shifting domains, causing existing cross-modal fusion schemes to degrade in performance. Yang *et al.* [48] and Planamente *et al.* [30] propose to directly fuse visual and audio features or predictions for cross-domain activity classification. However, the effectiveness of these methods is reduced when not all activities make a characteristic sound. Different from previous works, we introduce audio-adaptive learning methods and a cross-modal interaction that utilizes the reliable domain-invariant cues within sound to help the video model adapt to the distribution shift.

We make three contributions in this paper. First, we propose an audio-adaptive encoder which exploits the rich information from sound to adjust the visual feature representation causing the model to learn more discriminative features in the target domain. This is done by preventing the model from

108 overfitting to domain-specific visual content, while simultaneously dealing with imbalanced semantic distributions 109 between domains. Second, we introduce an audio-infused 110 recognizer, which eliminates domain-specific features further 111 and allows effective cross-modal interaction across domains 112 by considering domain-invariant activity information within 113 sound. As a third contribution, we introduce the new task of 114 *actor shift*, and a corresponding audio-visual video dataset 115 *ActorShift*, to challenge our approach when the change in 116 actors results in large variation in activity appearance. Our 117 experiments on EPIC-Kitchens [12], CharadesEgo [33] and 118 *ActorShift*, demonstrate the advantage of our approach under 119 various video distribution shifts for both audible and silent 120 activities.

2. Related Work

123 **Sound for activity recognition.** Many works have utilized 124 sound for within-domain activity recognition in videos, 125 *e.g.*, [18, 21, 24, 25, 38, 39]. Since there is a natural correlation 126 between the visual and auditive elements of a video, Korbar 127 *et al.* [23] and Asano *et al.* [1] learn audio-visual models in a 128 self-supervised manner. As processing audio signals is much 129 faster than video frames, both Gao *et al.* [18] and Korbar *et* 130 *al.* [24] reduce computation by previewing the audio track 131 for video analysis. Cross-modal attention is widely used 132 in activity localization [25, 39, 46] and audiovisual video 133 parsing [38, 45] to guide the visual model to focus on the 134 audible regions. Zhang *et al.* [49] conduct repetitive activity 135 counting by using audio signals to decide the sampling rate 136 and predict the reliability of the visual features. As opposed 137 to most works which rely on sound for within-domain activi- 138 ty recognition, we consider its domain-invariant nature for 139 activity recognition across different domains.

140 **Video domain adaptation by vision.** The field of vision- 141 focused domain adaptation is extensive (see recent surveys 142 [43, 51]). Here, we focus on video domain adaptation for 143 activity recognition. State-of-the-art visual-only solutions 144 learn to reduce the shift in activity appearance by adversarial 145 training [5, 6, 8, 9, 20, 27, 29] and self-supervised learning 146 techniques [9, 22, 27, 34]. While Jamal *et al.* [20] and Munro 147 and Damen [27] directly penalize domain specific features 148 with an adversarial loss at every time stamp, Chen *et al.* [5], 149 Choi *et al.* [9] and Pan *et al.* [29] attend to temporal seg- 150 ments that contain important cues. Self-supervised learning 151 objectives are also incorporated in [27] and [9] to better align 152 the features across domains by utilizing the correspondences 153 between RGB and optical flow or the temporal order of video 154 clips. Song *et al.* [34] and Kim *et al.* [22] obtain remarkable 155 performance by contrastive learning for self-supervised 156 learning to align the feature distributions between video 157 domains. Instead of relying on the vision modality only, which 158 may present large activity appearance variance, we consider 159 the domain-invariant information within sound to help the 160

161 model adapt to the visual distribution shift.

162 **Video domain adaptation by vision and audio.** As audio 163 signals contain valuable domain-invariant cues, some recent 164 works recognize activities across domains with the aid of 165 sound. Yang *et al.* [48] directly fuse the features from visual 166 and audio modalities before classification. However, this 167 can lead to the visual features dominating the classification 168 since many activities are silent and the audio features are 169 less discriminative. As a result, the complementary informa- 170 tion from sound may not be considered. Planamente *et* 171 *al.* [30] instead align the two modalities with an audio-visual 172 loss. Nonetheless, the audio predictions for silent activities 173 remain unreliable and limit their performance improvements. 174 Instead, we propose audio-adaptive learning that exploits the 175 supervisory signals from sound to adjust to the distribution 176 shift and handle both audible and silent activities.

177 Additionally, existing datasets *e.g.*, [3, 12, 33, 35] focus 178 on humans meaning activities are inherently close in appear- 179 ance and share commonalities in hand-object interactions. 180 Inspired by the A2D dataset by Xu *et al.* [47], which contains 181 multiple actor classes for activity recognition, we introduce 182 the challenging domain adaptation setting of *actor shift*, in 183 which the shift between humans and animals performing the 184 action results in large appearance and motion differences 185 across domains, further facilitating video domain adaptation 186 by the use of vision and audio.

3. Approach

187 For activity recognition under domain shift, we consider 188 unsupervised domain adaptation where we have: a set of 189 labeled source videos $\mathcal{S}=\{(X_1^{\mathcal{S}}, y_1^{\mathcal{S}}), \dots, (X_N^{\mathcal{S}}, y_N^{\mathcal{S}})\}$ and a 190 set of unlabeled target videos $\mathcal{T}=\{X_1^{\mathcal{T}}, \dots, X_M^{\mathcal{T}}\}$. In each 191 domain, X and y indicate a video sample and the correspond- 192 ing activity class label, while N and M are the number of 193 samples in the source and target domain. Using all available 194 training data from the source and the target domains, the 195 task is to train an activity recognition model, which performs 196 well on (unseen) videos from the target domain.

197 We train our audio-adaptive model in two stages using 198 videos from source and target domains with accompanying 199 audio. In the first stage we train our audio-adaptive encoder 200 (Section 3.1) that uses audio to adapt a visual encoder to 201 be more robust to distribution shifts. In the second stage 202 we train our audio-infused recognizer (Section 3.2) using 203 pseudo-labels from the audio-adaptive encoder for the target 204 domain and the ground-truth labels for the source domain. 205 The audio-infused recognizer maps the source and target 206 domains into a common space and fuses audio and visual 207 features to produce an activity prediction for either domain.

3.1. Stage 1: Audio-Adaptive Encoder

208 Our audio-adaptive encoder $\mathcal{E}(\cdot)$, detailed in Figure 2, 209 consists of a visual encoder $\mathcal{V}(\cdot)$, an audio encoder $\mathcal{A}(\cdot)$, 210

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230

Target domain audio
Source domain audio
Target domain video
Source domain video

Audio encoder $\mathcal{A}(\cdot)$
Audio-based attention $\Psi(\cdot)$
Transformer
Attention vector
Visual encoder $\mathcal{V}(\cdot)$

Audio prediction
Visual prediction

Absent-activity learning
Target domain
Absent-activity loss
Pseudo-absent label
Visual prediction

Audio-balanced learning
Source domain
Clustering
Assigning weights
Groundtruth
Activity class
Audio-balanced loss

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284

231 **Figure 2. Audio-adaptive encoder for activity recognition under domain shift.** With a pretrained audio encoder, we train the visual
232 encoder and audio-based attention module, which guides the visual encoder to focus on the activity relevant features. We do this with two
233 audio-adaptive learning methods: absent-activity learning and audio-balanced learning. The absent activity learning operates in the target
234 domain and uses the audio predictions to indicate which activities cannot be heard in the video. The visual predictions are then encouraged
235 to have low probabilities for these ‘pseudo-absent’ activities. The audio-balanced learning uses audio in the source domain to cluster samples
236 in each activity class into clusters according to the sounds of the object/environment interacted with. In the audio-balanced loss the rare
237 activities and interactions are weighted higher to handle the semantic shift between domains.

285
286
287
288
289
290
291
292

239 and an audio-based attention module $\psi(\cdot)$. Since the sounds
240 of activities have less variance across domains, $\mathcal{E}(\cdot)$ aims
241 to extract visual features that are invariant but discriminative
242 under domain shift with the aid of $\mathcal{A}(\cdot)$ pretrained for
243 audio-based activity recognition. To this end, we train $\mathcal{V}(\cdot)$
244 and $\psi(\cdot)$ with two audio-adaptive learning methods: absent-
245 activity learning for unlabeled target data and audio-balanced
246 learning for labeled source data. The former aims to remove
247 irrelevant parts of the visual features while the latter helps
248 to handle the differing label distribution between domains.
249 Once trained, for each video, we can extract an audio feature
250 vector from $\mathcal{A}(\cdot)$ and a series of visual features from
251 $\mathcal{V}(\cdot)$ with which to train our audio-infused recognizer (Section
252 3.2) for activity classification.

293
294
295
296
297
298
299
300
301
302

253 **Audio-based attention.** We use an audio-based attention
254 module $\psi(\cdot)$ to adapt the visual encoder to focus on activity-
255 relevant features. For example, the visual model may pre-
256 dict the activity *washing* because of the presence of a sink.
257 However, without the sound of water the attention module
258 suppresses the channels encoding the sink thus increasing
259 the prediction of the correct class. The attention module
260 is based on the transformer encoder [13, 14, 42]. It takes
261 the audio features as input and outputs the channel attention
262 feature vector, which is multiplied with the visual features.

303
304
305
306
307
308
309
310
311
312
313
314
315

263 **Absent-activity learning.** The absent-activity learning uses
264 audio in the target domain to train the attention module and
265 visual encoder. Naively, we could treat the class with the
266 highest probability from the visual encoder as the pseudo
267 label. However, doing so can create biased pseudo-labels
268 as irrelevant objects often appear in a scene. Instead, we

use the audio predictions to guide the visual pseudo-labels.
While we may not be confident which activity is happening
in a video, particularly for silent videos, we can often be
confident that certain activities with distinctive sounds are
not occurring in a video. We call these “absent activities”.
To learn from these absent activities, we generate pseudo-
absent labels for the unlabeled target domain videos, which
indicate the activities with the lowest probabilities from the
audio encoder. The visual encoder is then encouraged to
predict these unlikely classes with low probability.

Specifically, for an unlabeled video X^T in the target
domain, we obtain the audio-based activity probability dis-
tribution $\mathbf{p}_a^T \in \mathbb{R}^K$ (K is the number of classes) from the
audio encoder $\mathcal{A}(\cdot)$ trained on labeled source data. From
this we obtain the set of absent activities \mathcal{Q} by taking the
lowest r predictions in \mathbf{p}_a^T , i.e., the classes with the lowest
probabilities from the audio encoder. We also extend this to
multi-label classification by instead assuming the $(1 - \alpha_k)\gamma$
percent videos with the lowest probabilities do not contain
class k , where $\gamma \in (0, 1]$ and α_k is the percentage of videos
containing each activity class in the labeled source domain.

Our loss for absent-activity learning is formulated as:

$$l_A(\mathbf{p}_v^T, \mathcal{Q}) = - \sum_{q \in \mathcal{Q}} \log(1 - p_{v,q}^T), \quad (1)$$

where $p_{v,q}^T$ is the probability output for the q th class for the
video X^T . With this loss, the visual encoder is able to ignore
confounding visual features and generate less-noisy pseudo-
labels for the target domain. This allows our model to better
capture high-level semantic information between domains

316
317
318
319
320
321
322
323

324 based on both appearance and motion cues.
 325 **Audio-balanced learning.** Besides change in visual appearance,
 326 domain shift can also be caused by change in label
 327 distributions. We address this challenge with our audio-
 328 balanced learning, which not only handles imbalance in
 329 activity classes, but also imbalance in terms of the objects or
 330 the environment being interacted with.

331 To this end, we first use k -means to group the video samples
 332 inside each activity class by their audio feature \mathbf{f}_a^S with
 333 the assumption that each group represents a different types
 334 of objects or environments. We use audio features for clustering
 335 as they can indicate the material of the interacted objects
 336 or the environment the action is performed in, while being
 337 invariant to appearance changes. The number of interaction
 338 clusters per activity class is determined by the Elbow
 339 method [37], which favours a small number while obtaining
 340 a low ratio of dispersion both between and within clusters.
 341

342 We based our *audio-balanced loss* on the class-balanced
 343 loss by Cui *et al.* [11]. When using the original class-
 344 balanced loss on a source domain video X^S with visual
 345 probabilities \mathbf{p}_v^S we can balance over our activity classes:

$$l_{CB}(\mathbf{p}_v^S, y^S) = \frac{1 - \beta}{1 - \beta^{n_y}} \mathcal{L}(\mathbf{p}_v^S, y^S), \quad (2)$$

346 where \mathcal{L} is a classification loss, *e.g.*, Softmax cross-entropy
 347 loss and n_y is the number of training samples of ground-
 348 truth activity class y . $\beta \in [0, 1]$ is a hyper-parameter which
 349 controls the weighting factor $\frac{1-\beta}{1-\beta^{n_y}}$. As $\beta \rightarrow 1$, this weight-
 350 ing factor becomes inversely proportional to the effective
 351 number of samples inside each class so that tail classes in
 352 the source domain are weighted higher in training.
 353

354 With our *audio-balanced loss* we include an additional
 355 weighting factor so the long tail of object interactions are
 356 also accounted for with our interaction clusters:
 357

$$l_B(\mathbf{p}_v^S, y^S) = \frac{1 - \beta}{1 - \beta^{n_{y,j}}} l_{CB}(\mathbf{p}_v^S, y^S). \quad (3)$$

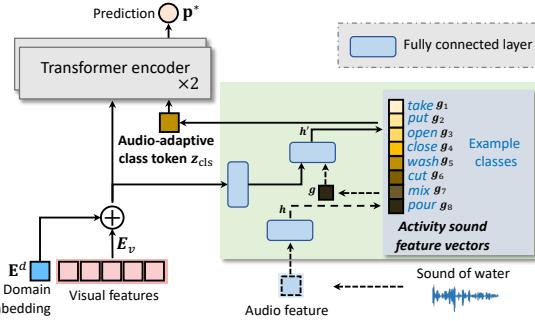
358 $n_{y,j}$ is the number of samples for the j th interaction cluster
 359 within ground-truth activity y^S . By this loss, both rare activi-
 360 ties and rare interactions from frequent activities are given
 361 a high weight during training. This means the classifier can
 362 generalize well to the target domain where the distribution
 363 of activities and interactions may not be the same.
 364

365 **Audio-adaptive encoder loss.** The absent-activity loss and
 366 the audio-balanced loss are combined to obtain the overall
 367 loss for training the visual encoder $\mathcal{V}(\cdot)$ and audio-based
 368 attention $\psi(\cdot)$ inside the audio-adaptive encoder $\mathcal{E}(\cdot)$:

$$l_E = \sum_{(X_i) \in \mathcal{T}} l_A(\mathbf{p}_{i,v}^T, Q_i) + \sum_{(X_j, y_j) \in \mathcal{T}} l_B(\mathbf{p}_{j,v}^S, y_j^S). \quad (4)$$

374 3.2. Stage 2: Audio-Infused Recognizer

375 While audio can help focus on the activity-relevant visual
 376 features, there is still a large difference between the



377 **Figure 3. Audio-infused recognizer.** We add domain embedding
 378 E_d to encourage a common visual representation across domains.
 379 Then, an audio-adaptive class token is obtained from a series of
 380 activity sound feature vectors, considering both audio and visual
 381 features. It is sent into the transformer together with the visual
 382 features. By the transformer’s self attention, this token aggregates
 383 information from visual features with the domain-invariant audio
 384 activity cues for activity classification.
 385

386 appearance of activities in different domains. To further
 387 eliminate domain-specific visual features and fuse the activi-
 388 ty cues from the audio and visual modalities we propose the
 389 audio-infused recognizer $\mathcal{R}(\cdot)$, visualized in Figure 3.
 390

391 **Transformer with domain embedding.** We adopt a trans-
 392 former encoder since its core mechanism, self-attention, can
 393 efficiently encode multi-modal representations [16, 36, 50].
 394 For a vanilla version, we take the input sequence:
 395

$$\mathbf{z}^m = [\mathbf{z}_{cls}^m; \mathbf{f}_{v,1}\mathbf{E}_v; \dots; \mathbf{f}_{v,n}\mathbf{E}_v; \mathbf{f}_{a,1}\mathbf{E}_a; \dots; \mathbf{f}_{a,n}\mathbf{E}_a], \quad (5)$$

396 where \mathbf{z}_{cls}^m is the learnable class token defined as in [14],
 397 and $\{\mathbf{f}_{v,1}, \dots, \mathbf{f}_{v,n} | \mathbf{f}_{v,\cdot} \in \mathbb{R}^{C_v}\}$ and $\{\mathbf{f}_{a,1}, \dots, \mathbf{f}_{a,n} | \mathbf{f}_{a,\cdot} \in \mathbb{R}^{C_a}\}$ are the visual and audio features of n clips from video
 398 X . $\mathbf{E}_v \in \mathbb{R}^{C_v \times D}$ and $\mathbf{E}_a \in \mathbb{R}^{C_a \times D}$ are linear projections
 399 to map the visual and audio features to D dimensions. To
 400 map source and target domains into a common space, we
 401 first learn a domain embedding $\mathbf{E}^d \in \mathbb{R}^D$ ($d \in \{\mathcal{S}, \mathcal{T}\}$),
 402 which is added to suppress domain-specific visual features.
 403 Then, the input sequence for the transformer becomes:
 404

$$\mathbf{z}' = [\mathbf{z}_{cls}^m; \mathbf{f}_{v,1}\mathbf{E}_v + \mathbf{E}^d; \dots; \mathbf{f}_{v,n}\mathbf{E}_v + \mathbf{E}^d; \mathbf{f}_{a,1}\mathbf{E}_a; \dots; \mathbf{f}_{a,n}\mathbf{E}_a]. \quad (6)$$

405 **Audio-adaptive class token.** Ideally, the transformer’s
 406 self attention will aggregate audio and visual features with
 407 the class token to predict the correct activity. However, the
 408 cross-modal correspondences are difficult to find under dis-
 409 tribution shift, meaning the prediction may rely on the more
 410 discriminative, but less domain-invariant, visual features.
 411 To address this, we propose to generate an audio-adaptive
 412 class token, which is initialized from the audio activity class
 413 prediction and gradually aggregates the visual features while
 414 keeping its own audio-based activity information through
 415 the transformer. As shown in Figure 3, the audio-adaptive
 416 class token is obtained from a series of activity sound vec-
 417 tors $\{\mathbf{g}_k \in \mathbb{R}^D\}_{k=1}^K$, with each representing an activity class.
 418

			Source Domain Setting		Target Domain Setting		
	Shift	Video Dataset	Source Domain	Train	Target Domain	Train	Test
432	Scenery	EPIC-Kitchens-55 [12]	Kitchens	7,935	Kitchens	7,935	2,114
433	Viewpoint	CharadesEgo [33]	Third-person view	3,083	Ego-centric view	3,083	825
434	Actor	ActorShift (<i>ours</i>)	Human actors	1,305	Animal actors	66	235

Table 1. **Domain adaptation benchmarks for activity recognition under scenery, viewpoint and actor shift** with the datasets used and number of videos per source and target split. Scenery and viewpoint shift are present in existing datasets. We propose the actor shift setting and dataset to tackle the challenge of a severe change in activity appearance. The dataset will be made available on the project website.

They capture global context information and serve as the representation bottleneck to provide regularization for model learning [2, 32]. For selection, the feature vector from the audio adaptive encoder $\mathcal{A}(X)$ is first processed by a fully connected layer to give the activity probabilities $\mathbf{h} \in \mathbb{R}^K$. Then, an initial vector is obtained by $\mathbf{g} = \sum_{k=1}^K h_k * \mathbf{g}_k$. We include visual features to help silent activities select the representative vector. To avoid the visual features dominating, we project them to a lower dimension with a fully connected layer before concatenating them with the initial vector \mathbf{g} . The concatenated vector is given to another fully connected layer which outputs the probabilities \mathbf{h}' for each type of activity sound. Finally, we obtain the audio representation $\mathbf{z}_{cls} = \sum_{k=1}^K h'_k * \mathbf{g}_k$, which serves as the class token. Consequently, the input sequence for the transformer becomes:

$$\mathbf{z} = [\mathbf{z}_{cls}; \mathbf{f}_{v,1} \mathbf{E}_v + \mathbf{E}^d, ; \dots; \mathbf{f}_{v,n} \mathbf{E}_v + \mathbf{E}^d], \quad (7)$$

where \mathbf{z}_{cls} is the audio-adaptive class token. The class token output state is further sent to a fully connected layer to get the final prediction \mathbf{p}^* . For audible activities, the activity sound vector can be accurately selected and kept discriminative for audiovisual interaction. For silent activities, the vector is obtained from environmental sound, which indicates the presence of multiple possible activities. The vector becomes more discriminative as the transformer progressively enhances it through the visual features.

Audio-infused recognizer loss. We train the audio-infused recognizer on both source and target videos with the loss:

$$l_{\mathcal{R}} = \sum_{(X_i, y_i) \in \{\mathcal{S}, \mathcal{T}\}} \mathcal{L}(\mathbf{p}_i^*, y_i) + \eta \left(\mathcal{L}(\mathbf{h}_i, y_i) + \mathcal{L}(\mathbf{h}'_i, y_i) \right), \quad (8)$$

where hyperparameter η balances the loss terms and y_i is the groundtruth or, in the case of the unlabeled video, the hard pseudo-label. \mathbf{p}_i^* is the final classification prediction, and \mathbf{h}_i and \mathbf{h}'_i are the probabilities for the activity sound vectors outputted by the first and second fully connected layers. The first term $\mathcal{L}(\mathbf{p}_i^*, y_i)$ optimizes the transformer to predict the correct activity class, while the second term $\mathcal{L}(\mathbf{h}_i, y_i) + \mathcal{L}(\mathbf{h}'_i, y_i)$ optimizes the activity sound vectors. We are now ready to validate the effectiveness of our approach on three domain adaptation benchmarks as highlighted in Figure 1, summarized in Table 1 and detailed next.

4. Domain Adaptation Benchmarks

Scenery shift. We study scenery shift in the *EPIC-Kitchens-55* [12] dataset, which contains 1st-person videos of fine-grained kitchen activities. The domain adaptation benchmark proposed by Munro and Damen [27] uses three domain partitions (D1, D2 and D3), where each domain is a different person in different kitchen. The task is to adapt between each pair of domains. This benchmark focuses on eight activity classes (verbs), which occur in combination with different objects, with a severe class imbalance. The kitchens have different appearances and contain different utensils.

Viewpoint shift. We consider viewpoint shift in the *CharadesEgo* dataset by Sigurdsson *et al.* [33]. It contains paired videos of the same activities, recorded from first and third-person perspective. It has 3,083 and 825 videos per viewpoint for training and testing, spanning 157 activity classes. Following [8], we treat the 3rd-person videos as the source domain and the 1st-person videos as the target domain. The changing views make the activities appear visually different, resulting in a large domain gap.

Actor shift. While both EPIC-Kitchens and CharadesEgo contain considerable domain shifts, there are still some inherent similarities between the domains in these datasets. Since all the actors are humans, latent signals describing the way hands and objects interact are shared between domains. Therefore, we introduce an even more challenging domain shift setting to further facilitate video domain adaption research and demonstrate the potential of our method. We introduce *ActorShift*, where the domain shift comes from the change in actor species: we use humans in the source domain and animals in the target domain. This causes large variances in the appearance and motion of activities.

For the corresponding dataset we select 1,305 videos of 7 human activity classes from Kinetics-700 [3] as the source domain: *sleeping*, *watching tv*, *eating*, *drinking*, *swimming*, *running* and *opening a door*. For the target domain we collect 200 videos from YouTube of animals performing the same activities. We divide them into 35 videos for training (5 per class) and 165 for evaluation. The target domain data is scarce, meaning there is the additional challenge of adapting to the target domain with few unlabeled examples.

Evaluation criteria. Following standard practice [27, 33],

540 we report top-1 accuracy on EPIC-Kitchens and ActorShift
 541 for single-label classification, and mAP (mean average pre-
 542 cision) on CharadesEgo for multi-label classification.
 543

544 5. Results

545 We first describe the implementation details before ablating
 546 the components of our method and comparing to prior
 547 works in each type of domain shift.

548 **Implementation details.** For our visual encoder $\mathcal{V}(\cdot)$ we
 549 use SlowFast [15], unless stated otherwise. For the audio
 550 encoder $\mathcal{A}(\cdot)$ we use ResNet-18 [19]. The audio-based at-
 551 tention module $\psi(\cdot)$ consists of eight transformer encoder
 552 layers [14] with a final fully connected layer to obtain the
 553 attention vector for the visual encoder. The inputs are
 554 intermediate audio features from $\mathcal{A}(\cdot)$ (conv3) along with a
 555 learnable class token defined as in [14] (note this is different
 556 from our audio-adaptive class token used in $\mathcal{R}(\cdot)$). The out-
 557 put state of the class token passes through the fully connected
 558 layer to obtain the attention vector for the visual encoder. We
 559 set the parameters of our absent activity loss to $r=3$, $\gamma=0.05$
 560 and $\beta=0.999$. For the audio-infused recognizer $\mathcal{R}(\cdot)$, we
 561 use two transformer encoder layers with the same architec-
 562 ture in [14]. The sequence dimension D is 512 and each
 563 layer has 8 self-attention heads. We provide more details in
 564 the supplementary and share anonymous demo code¹ based
 565 on mmaction2 [10]. The full version will be released.

566 5.1. Ablation Study

567 For ablations we use RGB and audio modalities on both
 568 EPIC-Kitchens and CharadesEgo. During training, all la-
 569 beled source videos are used. With EPIC-Kitchens all target
 570 videos are unlabelled, while for CharadesEgo we use half
 571 labelled and half unlabelled for semi-supervised domain
 572 adaptation as in [8]. Since EPIC-Kitchens contains multiple
 573 adaptation settings, we report the average. Ablations on
 574 component internals are provided in the supplementary.

575 **Stage 1: Audio-adaptive encoder.** We report results in Ta-
 576 ble 2. We first consider the audio-adaptive encoder alone.
 577 Initially, we train only the visual encoder with a standard
 578 softmax cross-entropy loss on the source domain. Simply
 579 generating channel attention for the visual features with our
 580 audio-based attention module already improves performance
 581 by 3.2% top-1 accuracy on EPIC-Kitchens and 0.4% mAP
 582 on CharadesEgo. Since audio contains useful activity in-
 583 formation, this attention helps the visual encoder focus on
 584 relevant features. Adding the absent-activity learning re-
 585 sults in 2.5% and 0.9% improvements, demonstrating that
 586 the pseudo-absent labels increase the discriminative ability
 587 of the model in the target domain. We observe that adopt-
 588 ing the audio-balanced learning and replacing the softmax
 589 cross-entropy with our audio-balanced loss delivers a further
 590

591 ¹[https://anonymous.4open.science/r/Domain-
 592 Adaptation-Demo-B413](https://anonymous.4open.science/r/Domain-Adaptation-Demo-B413)

Model	EPIC-Kitchens	CharadesEgo	594
	Top-1 (%) ↑	mAP (%) ↑	595
Stage 1: Audio-adaptive encoder $\mathcal{E}(\cdot)$			
Visual encoder $\mathcal{V}(\cdot)$	48.0	23.1	596
+ Audio-based attention $\psi(\cdot)$	51.2	23.5	597
+ Absent-activity learning	53.7	24.4	598
+ Audio-balanced learning	55.7	25.0	599
Stage 2: Audio-infused recognizer $\mathcal{R}(\cdot)$			
+ Vanilla multi-modal transformer \mathbf{z}^m	56.1	25.0	600
+ Domain embedding \mathbf{z}'	57.2	25.4	601
+ Audio-adaptive class token \mathbf{z}	59.2	26.3	602

Table 2. **Model components ablation.** All components in the audio-adaptive encoder and the audio-infused recognizer contribute to performance improvement under distribution shift. For both EPIC-Kitchens and CharadesEgo the improvements over a vanilla SlowFast visual encoder are considerable.

Model	Activities		Overall	609
	Silent	Audible	mAP (%) ↑	610
Visual encoder $\mathcal{V}(\cdot)$	23.2	22.7	23.1	611
Full model	26.3	25.9	26.3	612

Table 3. **Benefit over silent and audible activities** on CharadesEgo. Our audio-adaptive model benefits both activity types.

2.0% and 0.6% increase. This highlights the importance of addressing the label distribution shift in domain adaption.

Stage 2: Audio-infused recognizer. For the audio-infused recognizer, we first consider a vanilla transformer. It takes as input \mathbf{z}^m (Eq. 5), *i.e.* the audio and visual features from the audio-adaptive encoder, mapped by \mathbf{E}_v and \mathbf{E}_a into a common space, alongside a learnable class token. This only gives a marginal improvement in results. Adding the domain embedding \mathbf{E}^d to reduce domain-specific visual features in \mathbf{z}' (Eq. 6) gives a benefit of 1.1% on EPIC-Kitchens and 0.4% on CharadesEgo. This is because the cross-modal correspondences become easier to discover. When we replace the plain audio features and single learnable class token with our audio-adaptive class token to get \mathbf{z} (Eq. 7), we observe larger improvements of 2.0% and 0.9%. This is expected, as the audio-adaptive class token better incorporates complementary information from sound for the final activity classification, with a standard learnable class token the visual features will dominate the fusion inside the transformer.

Benefit for silent activities. In Table 3, we demonstrate the effect of our full model on silent and audible activities separately. We focus on CharadesEgo since only 13 out of 157 classes have a characteristic sound (see supplementary). Our model obtains ~3% absolute increase for both silent and audible activities over a visual-only encoder. We conclude that audio is helpful for handling visual distribution shifts even for activities which do not have a characteristic sound.

Benefit for silent videos. We have also tested our approach when the audio track is available for training but unavailable during inference. On EPIC-Kitchens, the audio-adaptive encoder achieves 50.7% top-1 accuracy, still an improvement

	EPIC-Kitchen Activity Recognition Across Domains											
	Method	RGB	Flow	Audio	D2 → D1	D3 → D1	D1 → D2	D3 → D2	D1 → D3	D2 → D3	Mean	
I3D backbone												702
Source-only [27]		✓	✓		42.5	44.3	42.0	56.3	41.2	46.5	45.5	705
Munro and Damen [27]		✓	✓		48.2	50.9	49.5	56.1	44.1	52.7	50.3	706
Planamente <i>et al.</i> [30] [†]		✓	✓	✓	48.5	50.9	49.7	56.3	44.8	52.5	50.5	707
Yang <i>et al.</i> [48] [†]		✓	✓	✓	49.2	51.0	49.8	56.5	45.7	52.3	50.8	708
Kim <i>et al.</i> [22]		✓	✓		49.5	51.5	50.3	56.3	46.3	52.0	51.0	709
Song <i>et al.</i> [34]		✓	✓		49.0	52.6	52.0	55.6	45.5	52.5	51.2	710
<i>This paper</i>		✓	✓	✓	53.6	53.4	54.9	61.2	52.3	61.1	56.1	711
SlowFast backbone												712
<i>This paper</i>		✓	✓	✓	59.3	59.1	59.5	69.1	54.8	64.3	61.0	713

[†] Based on our re-implementation using our features for RGB, flow and audio.

Table 4. **Activity recognition under scenery shift** on EPIC-Kitchens for the unsupervised domain adaptation setting. Our audio-adaptive model achieves state-of-the-art top-1 accuracy, and benefits from audio more than the audio-visual fusion methods used in prior works [30, 48]. Results increase further with a SlowFast backbone. More comparisons and modality-combinations are provided in the supplementary.

over visual encoder only (48.0%). With both the audio-adaptive encoder and audio-infused recognizer, the result improves to 51.2%. This indicates our approach effectively uses audio to help the visual encoder learn a more discriminative feature representation in the target domain, even when audio is absent during inference.

Benefit for the long-tail. In Figure 4, we demonstrate the benefit of audio-balanced learning towards activities that are rare in the source domain but are more frequent in the target domain. We use EPIC-Kitchens since it contains a long-tail of different object interactions (nouns) in each activity class (verb). We treat verb-noun pairs as frequent when they occur more than 10 times in the source domain, else they are considered rare. As the distribution of activities (verbs) changes across domains, the class-balanced loss [11] improves over the standard softmax cross-entropy loss. However, the domain shift also causes imbalance in the distribution of interactions (nouns). Because we balance the loss of each pseudo-interaction by clustering, our audio-balanced loss is especially helpful for the rare interactions (0-1 and 2-10 instances) where it obtains ~3.5% improvement. In comparison to the class-balanced loss we are slightly worse on frequent interactions, as we give higher weight to less common interactions. As interactions have a long-tail, our audio-balanced loss does result in an overall improvement.

5.2. Comparison with State-of-the-Art

Scenery shift. We first demonstrate the effectiveness of our approach for domain adaptation on EPIC-Kitchens, as defined by Munro and Damen [27]. Here, different domains mean a change in scenery. The results are shown in Table 4. We first note that our approach gives ~5% improvement over the best performing prior works with the same I3D backbone. A further ~5% improvement can be gained from using SlowFast as the backbone [15]. There are several reasons for this improvement. First, our model utilizes the domain-invariant nature of audio signals to produce reliable

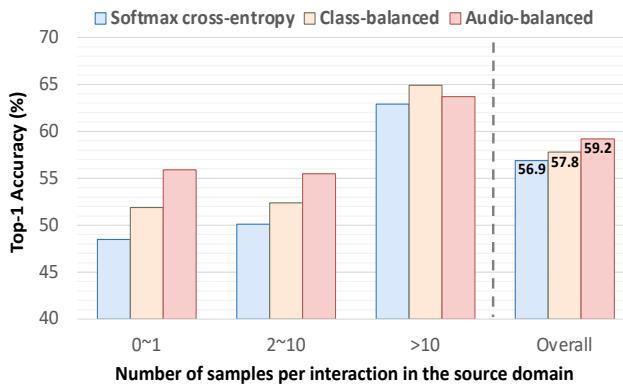


Figure 4. **Benefit for the long-tail** on EPIC-Kitchens. Our audio-balanced loss learns rare activities in the source domain to generalize better to unknown activity distributions in the target domain.

pseudo-absent labels for the target domain video during training. This is particularly helpful for 1st-person videos where the activity may happen out of view. In addition, both RGB and Flow suffer from large appearance variance making it harder to guide domain-adaption through these modalities alone. Second, since the dataset has imbalanced label distributions, treating all the classes and interactions equally, as in prior works, results in inaccurate predictions when the semantic distribution shifts.

We also compare our full model with alternative audio-visual approaches proposed for cross-domain activity recognition [30, 48]. We let both of them use the same inputs, *i.e.*, the features as outputted by the visual and audio encoders. Both of them use an adversarial loss to first align the visual features between domains and fuse visual and audio features or predictions afterwards. This causes the visual features to dominate the classification while the complementary information from sound may not be considered. Planamente *et al.* [30] introduce an audio-visual loss, so the two modalities make a more balanced contribution towards the prediction. However, the audio predictions for silent activities are unreliable.

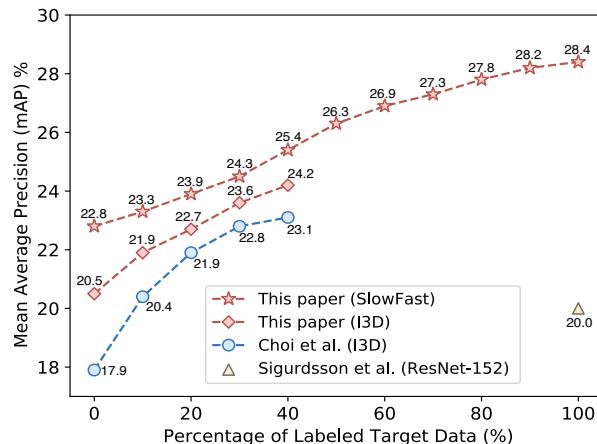


Figure 5. **Activity recognition under viewpoint shift** on CharadesEgo. Using all source data, we compare with [8, 33] under varying amounts of labeled target training data. Our model obtains favourable results under all settings.

able and harm their accuracy. Our model better combines the complementary information in the audio and visual modalities, effectively coping with many activities being silent.

Viewpoint shift. In this comparison we consider viewpoint shift in CharadesEgo [33], following the semi-supervised setting of Choi *et al.* [8]. Meaning we have some labeled target domain videos available during training. The results are shown in Figure 5. Our method achieves better results than Choi *et al.* [8] with the same I3D RGB backbone [3], for all amounts of labeled target videos. When adopting the SlowFast RGB backbone [15], we again further improve performance for all settings. In the supplementary, we also provide a favorable comparison with Li *et al.* [26] under their fully-supervised setting. Since CharadesEgo contains paired 1st person and 3rd person videos, we can test whether our method needs to see the same action instance from different viewpoints as in previous methods [33] or whether it can make use of unpaired videos. When half of the paired videos from both views are used, we achieve a mAP of 29.9. When we use unpaired videos, the performance remains unchanged. We conclude our approach does not require paired training videos to be robust to viewpoint shift.

Actor shift. For this experiment, we use our ActorShift dataset and compare our model with the method by Munro and Damen [27], as their code is available. For fair comparison, we replace their I3D backbone with the same SlowFast backbone used for our model. We also show a baseline of the SlowFast model trained on source domain video only. The results are shown in Figure 6. While the method proposed by Munro and Damen [27] achieves good performance, our audio-adaptive approach better handles the large activity appearance variance caused by the shift in actors. For example, humans and animals sleep in visually different places and positions, while the sound of snoring or breathing is common

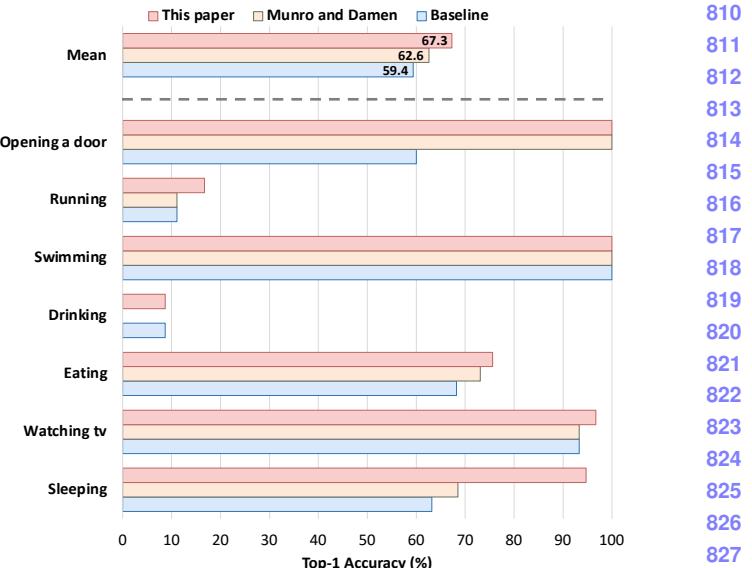


Figure 6. **Activity recognition under actor shift** on our ActorShift dataset. When the visual similarities for the same activity are difficult to discover between domains, our model can use additional cues from sound to improve the recognition accuracy.

to both. All models struggle with silent activities when there is both a large shift in appearance and a significant difference in sounds of activities between the domains, such as *drinking* and *running*. We provide examples in the supplemental material, which are of interest for future work.

6. Discussion

Limitations. During training our approach needs videos from both source and target domains, and all should have an audio track, limiting our approach to multi-modal video training sets. While audio at test-time is not required, it benefits our activity recognition results considerably.

Potential negative impact. When deployed our approach will have to record, store and process video and audio information related to human activities, which will have privacy implications for some application domains.

Conclusions. We propose to recognize activities under domain shift with the aid of sound, using a novel audiovisual model. By leveraging the domain-invariant activity information within sound, our model improves over both silent and audible activities as well as rare activities in the source domain. Experiments on two domain adaptation benchmarks demonstrate that our approach has better adaptation ability than visual-only solutions and benefits from audio more than alternative audiovisual fusion methods used in prior works. We also show that our model better handles large activity appearance variance caused by the shift in actors.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

References

- [1] Yuki M. Asano, Mandela Patrick, Christian Rupprecht, and Andrea Vedaldi. Labelling unlabelled videos from scratch with multi-modal self-supervision. In *NeurIPS*, 2020. 2
- [2] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *ICML*, 2018. 5
- [3] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *CVPR*, 2017. 2, 5, 8
- [4] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. VGGSound: a large-scale audio-visual dataset. In *ICASSP*, 2020. 11
- [5] Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. Temporal attentive alignment for large-scale video domain adaptation. In *ICCV*, 2019. 1, 2
- [6] Min-Hung Chen, Baopu Li, Yingze Bao, Ghassan Al-Regib, and Zsolt Kira. Action segmentation with joint self-supervised temporal domain adaptation. In *CVPR*, 2020. 2, 11
- [7] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *CVPR*, 2018. 11
- [8] Jinwoo Choi, Gaurav Sharma, Manmohan Chandraker, and Jia-Bin Huang. Unsupervised and semi-supervised domain adaptation for action recognition from drones. In *WACV*, 2020. 2, 5, 6, 8
- [9] Jinwoo Choi, Gaurav Sharma, Samuel Schulter, and Jia-Bin Huang. Shuffle and attend: Video domain adaptation. In *ECCV*, 2020. 1, 2, 11
- [10] MMAAction2 Contributors. Openmmlab's next generation video understanding toolbox and benchmark. <https://github.com/open-mmlab/mmaction2>, 2020. 6, 11
- [11] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019. 4, 7
- [12] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, 2018. 2, 5
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3, 11
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 3, 4, 6, 11, 13
- [15] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019. 6, 7, 8, 11
- [16] Valentin Gabeur, Chen Sun, Karteek Alahari, and Cordelia Schmid. Multi-modal transformer for video retrieval. In *ECCV*, 2020. 4, 14
- [17] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 11
- [18] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. In *CVPR*, 2020. 1, 2
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6, 11
- [20] Arshad Jamal, Vinay P Namboodiri, Dipti Deodhare, and KS Venkatesh. Deep domain adaptation in action space. In *BMVC*, 2018. 1, 2, 11
- [21] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *ICCV*, 2019. 2
- [22] Donghyun Kim, Yi-Hsuan Tsai, Bingbing Zhuang, Xiang Yu, Stan Sclaroff, Kate Saenko, and Manmohan Chandraker. Learning cross-modal contrastive features for video domain adaptation. In *ICCV*, 2021. 1, 2, 7, 11, 14
- [23] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, 2018. 2
- [24] Bruno Korbar, Du Tran, and Lorenzo Torresani. SCSampler: Sampling salient clips from video for efficient action recognition. In *ICCV*, 2019. 1, 2
- [25] Jun-Tae Lee, Mihir Jain, Hyoungwoo Park, and Sungrock Yun. Cross-attentional audio-visual fusion for weakly-supervised action localization. In *ICLR*, 2021. 1, 2, 14
- [26] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *CVPR*, 2021. 8, 13
- [27] Jonathan Munro and Dima Damen. Multi-modal domain adaptation for fine-grained action recognition. In *CVPR*, 2020. 1, 2, 5, 7, 8, 11, 14
- [28] Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. In *NeurIPS*, 2021. 1, 14
- [29] Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Niebles. Adversarial cross-domain action recognition with co-attention. In *AAAI*, 2020. 1, 2
- [30] Mirco Planamente, Chiara Plizzari, Emanuele Alberti, and Barbara Caputo. Cross-domain first person audio-visual action recognition through relative norm alignment. *arXiv preprint arXiv:2106.01689*, 2021. 1, 2, 7
- [31] Nishant Rai, Haofeng Chen, Jingwei Ji, Rishi Desai, Kazuki Kozuka, Shun Ishizaka, Ehsan Adeli, and Juan Carlos Niebles. Home action genome: Cooperative compositional action understanding. In *CVPR*, 2021. 1
- [32] Alexandre Rame and Matthieu Cord. Dice: Diversity in deep ensembles via conditional redundancy adversarial estimation. In *ICLR*, 2021. 5
- [33] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Actor and observer: Joint modeling of first and third-person videos. In *CVPR*, 2018. 2, 5, 8

- 972 [34] Xiaolin Song, Sicheng Zhao, Jingyu Yang, Huanjing Yue,
973 Pengfei Xu, Runbo Hu, and Hua Chai. Spatio-temporal con-
974 trastive domain adaptation for action recognition. In *CVPR*,
975 2021. 1, 2, 7, 11, 14 1026
976 [35] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah.
977 Ucf101: A dataset of 101 human actions classes from videos
978 in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 2 1027
979 [36] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and
980 Cordelia Schmid. Videobert: A joint model for video and
981 language representation learning. In *ICCV*, 2019. 4 1028
982 [37] Robert L Thorndike. Who belongs in the family? *Psychome-
983 trika*, 18(4):267–276, 1953. 4, 12 1029
984 [38] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified mul-
985 tisensory perception: weakly-supervised audio-visual video
986 parsing. In *ECCV*, 2020. 1, 2, 14 1030
987 [39] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chen-
988 liang Xu. Audio-visual event localization in unconstrained
989 videos. In *ECCV*, 2018. 1, 2 1031
990 [40] Yi-Hsuan Tsai, Kihyuk Sohn, Samuel Schulter, and Manmo-
991 han Chandraker. Domain adaptation for structured output via
992 discriminative patch representations. In *ICCV*, 2019. 11 1032
993 [41] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell.
994 Adversarial discriminative domain adaptation. In *CVPR*, 2017.
995 11 1033
996 [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,
997 Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia
998 Polosukhin. Attention is all you need. In *NIPS*, 2017. 3, 11 1034
999 [43] Mei Wang and Weihong Deng. Deep visual domain adapta-
1000 tion: A survey. *Neurocomputing*, 312:135–153, 2018. 2 1035
1001 [44] Weiyao Wang, Du Tran, and Matt Feiszli. What makes train-
1002 ing multi-modal classification networks hard? In *CVPR*, 2020.
1003 13 1036
1004 [45] Yu Wu and Yi Yang. Exploring heterogeneous clues for
1005 weakly-supervised audio-visual video parsing. In *CVPR*,
1006 2021. 1, 2 1037
1007 [46] Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention
1008 matching for audio-visual event localization. In *ICCV*, 2019.
1009 1, 2 1038
1010 [47] Chenliang Xu, Shao-Hang Hsieh, Caiming Xiong, and Jason J
1011 Corso. Can humans fly? action understanding with multiple
1012 classes of actors. In *CVPR*, 2015. 2 1039
1013 [48] Lijin Yang, Yifei Huang, Yusuke Sugano, and Yoichi Sato.
1014 Epic-kitchens-100 unsupervised domain adaptation challenge
1015 for action recognition 2021: Team m3em technical report.
1016 *arXiv preprint arXiv:2106.10026*, 2021. 1, 2, 7 1040
1017 [49] Yunhua Zhang, Ling Shao, and Cees GM Snoek. Repetitive
1018 activity counting by sight and sound. In *CVPR*, 2021. 1, 2 1041
1019 [50] Linchao Zhu and Yi Yang. Actbert: Learning global-local
1020 video-text representations. In *CVPR*, 2020. 4 1042
1021 [51] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi,
1022 Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A
1023 comprehensive survey on transfer learning. *Proceedings of
1024 the IEEE*, 109(1):43–76, 2020. 2 1043
1025