

# Multi-Expert Human Action Recognition with Hierarchical Super-Class Learning

Hojat Asgarian Dehkordi, Ali Soltani Nezhad, Hossein Kashiani, Shahriar Baradaran Shokouhi, Ahmad Ayatollahi

*School of Electrical Engineering, Iran University of Science and Technology, Tehran, Iran*

---

## Abstract

In still image human action recognition, existing studies have mainly leveraged extra bounding box information along with class labels to mitigate the lack of temporal information in still images; however, preparing extra data with manual annotation is time-consuming and also prone to human errors. Moreover, the existing studies have not addressed action recognition with long-tailed distribution. In this paper, we propose a two-phase multi-expert classification method for human action recognition to cope with long-tailed distribution by means of super-class learning and without any extra information. To choose the best configuration for each super-class and characterize inter-class dependency between different action classes, we propose a novel Graph-Based Class Selection (GCS) algorithm. In the proposed approach, a coarse-grained phase selects the most relevant fine-grained experts. Then, the fine-grained experts encode the intricate details within each super-class so that the inter-class variation increases. Extensive experimental evaluations are conducted on various public human action recognition datasets, including Stanford40, Pascal VOC 2012 Action, BU101+, and IHAR datasets. The experimental results demonstrate

---

*Email addresses:* [h\\_asgariandehkordi@elec.iust.ac.ir](mailto:h_asgariandehkordi@elec.iust.ac.ir) (Hojat Asgarian Dehkordi), [ali\\_soltaninezhad@elec.iust.ac.ir](mailto:ali_soltaninezhad@elec.iust.ac.ir) (Ali Soltani Nezhad), [hossein\\_kashiyani@alumni.iust.ac.ir](mailto:hossein_kashiyani@alumni.iust.ac.ir) (Hossein Kashiani), [bshokouhi@iust.ac.ir](mailto:bshokouhi@iust.ac.ir) (Shahriar Baradaran Shokouhi), [ayatollahi@iust.ac.ir](mailto:ayatollahi@iust.ac.ir) (Ahmad Ayatollahi)

*Preprint submitted to Knowledge-Based Systems*

that the proposed method yields promising improvements. To be more specific, in IHAR, Sanford40, Pascal VOC 2012 Action, and BU101+ benchmarks, the proposed approach outperforms the state-of-the-art studies by 8.92%, 0.41%, 0.66%, and 2.11 % with much less computational cost and without any auxiliary annotation information. Besides, it is proven that in addressing action recognition with long-tailed distribution, the proposed method outperforms its counterparts by a significant margin.

*Keywords:* Action Recognition, Still Images, Super-Class Learning, Long-Tailed Classification.

---

## 1. Introduction

Recently, human action recognition has attracted significant attention in computer vision due to its real-world applications. Action recognition in computer vision is generally categorized into action recognition in videos [1–4] and in still images [5–7]. Action recognition in still images aims to identify the type of human activity in a static image without any temporal information. Since a wide range of human activities such as *running* and *smoking* can be identified with a single input image and without extra motion cues, image-based action recognition has gained considerable attention. However, the lack of temporal information makes image-based action recognition more challenging rather than video-based action recognition [8]. Human action recognition (HAR) has numerous applications such as sports analysis, abnormal behavior recognition, wildlife observation, image tagging, image retrieval, and human-computer interaction [4, 9]. As still images lack temporal information, typically, action recognition studies extract spatial information from images. In this respect, initial studies in HAR have mostly employed low-level feature extraction techniques to capture low-level structures; however, they fail to achieve reliable and satisfactory results [10]. Another category of approaches explores to leverage object detector

[11] or pose estimator [12] developments to detect the keypoint joints [13], which are the most discriminative regions in the foreground area. Such detected areas favorably contribute to the overall action recognition accuracy. For instance, the detected areas in relation to the bicycle instance and the lower part of the human body could provide the methods with discriminative features in action recognition for *riding a bike* class.

In recent years, Convolutional Neural Networks (CNNs) have demonstrated their superior capabilities in several computer vision tasks and have emerged as a promising tool for HAR [4, 14–16]. CNNs have significantly advanced the keypoint detection and estimation for action recognition. Nevertheless, the standard CNN-based studies mainly have three constraints. First, for the training phase, these studies need a large amount of annotations regarding the human bounding boxes or body parts as well as action labels. Second, though CNNs could extract rich feature hierarchies, they generally struggle to extract the structural features for modeling the relationship among human keypoints in action recognition. This is ascribed to the fact that all activities are treated uniformly without considering any correlation and similarity context. With this learning process, the misclassification between relevant activities such as *Reading a book* and *Writing a book* is penalized same as other irrelevant action activities such as *Fixing a bike*. Third, most studies in HAR conjecture a balanced class distribution in their training phase over the existing well-organized balanced datasets such as Stanford40 [17], Pascal VOC 2012 Action [18] and BU101+ [19] datasets. Nevertheless, human instances exist at various frequencies in different class activities naturally, and this makes the underlying class distribution of the real-world dataset severely imbalanced. Despite the fact that a large number of research activities have been carried out with respect to action recognition [5, 12, 20–22], to the best of our knowledge, none of them have

taken into account the class imbalance issue in nature for action recognition in still images. This challenge impairs their performance when employed in large-scale real-world datasets. To address this challenge, we introduce a new large-scale Imbalanced Human Action Recognition dataset (IHAR) with a long-tailed distribution so that we can fairly assess the generalization of our proposed approach in comparison to the state-of-the-art studies.

In this paper, we propose a super-class learning approach for action recognition, named SCLAR, in still images. The merits of SCLAR are threefold: **1)** Our hierarchical SCLAR can detect visually distinct and subtle differences between various classes for action recognition. When it comes to the super-class learning, researchers attempt to promote inter-class variation between different classes. To reach this objective, we decompose the action recognition task into two-phase classification problem. First, a bucket of separate lightweight CNN classifiers (also called fine-grained classifiers) are pre-trained for the subsets of classes (i.e., super-classes). Then, a coarse-grained classifier is adopted to determine the relevant fine-grained classifiers in the first phase. The fine-grained classifiers ultimately would output the final classification label. With such a methodology, different activity classes are routed downstream to different fine-grained classifiers such that the inter-class variance among different classes increases. Thanks to the discriminative two-phase framework, specialized features are tuned to discriminate subtle visual differences in similar and different action classes from each other. **2)** Our action recognition framework surpasses the state-of-the-art approaches in HAR, while demanding much less computational cost and memory in that the coarse-grained and fine-grained classifiers in our framework are both compact models. As the action recognition task is a fine-grained classification, recent state-of-the-art studies require deeper CNN models with high capacity and computational complexity to detect subtle dif-

ferences between similar classes. However, in a wide range of applications, we require low computation load for deployment on edge devices. SCLAR addresses this bottleneck observed in other studies while outperforming its competitors.

**3)** Our framework can also address action recognition with long-tailed distribution. Since distinctive classes in super-class share knowledge in SCLAR, the under-represented classes manage to generalize better. Class imbalance occurs when there is considerable discrepancy among the cardinality of various classes [23, 24]. In the wild, we mostly encounter an inherent imbalance issue concerning human activity classes in action recognition. This causes the training loss to be biased toward the over-represented classes while simultaneously rendering the under-represented classes unable to reflect intra-class variation. When the knowledge in the coarse-grained phase is transferred to the fine-grained phase through a hierarchical structure, different super-classes value specialized subsets of features in relation to the under-represented classes. Therefore, the learned representation can focus better on the specific classes in a super-class with low cardinality and hard samples. After all, our framework would ideally relieve the data imbalance issue in the action recognition task.

The rest of the paper is structured as follows: Section 2 provides a literature review of the most recent HAR research. Section 3 presents an elaborate description of the proposed two-phase multi-expert classification method. This section also provides a detailed explanation for the Graph-based Class Selection Method. Section 4 introduces our large-scale dataset for action recognition in still images with the long-tailed distribution. Section 5 offers exhaustive assessments of the proposed method on several datasets to gauge the contributions of different ingredients in our framework. Finally, section 6 concludes our research work.

## 2. Related Works

In this section, first, the studies in relation to CNN-based human action recognition are explained. Then, the methods equipped with an object or pose detector are addressed. Ultimately, the state-of-the-art studies regarding data imbalance issue are described in more detail.

### 2.1. Still Image Action Recognition

In recent years, numerous methods in visual object classification [25], object detection [26], object tracking [27], and action recognition [3, 5, 6] have been proposed based on CNNs. In image-based HAR, the state-of-the-art studies exploit CNNs along with object detection module to locate human agents and mitigate irrelevant noisy context. Authors in [28] first detect different objects and their attributes and then compute a weighted sum of such detected instances for image-based HAR. Yan et al. [29] attempt to represent still images as a bag of image patches which are extracted by means of region proposal approaches. To do so, the FV encoding is applied to CNN features of image patches, and the spatial pyramid representation is employed for spatial feature extraction. In [21], the scene-level and region-level contexts are captured at the same time through a well-designed multi-branch attention networks. Two context branches for scene and region attention and a branch for target human region classification are incorporated into the proposed network for HAR. The requirement of human bounding boxes in still images is relaxed for action recognition in [30]. For this relaxation, the detected object proposals via selective search are decomposed into fine-grained components, and the final action predictions are calculated using an efficient product quantization to take into account the human-object interaction areas.

## 2.2. Object/Pose-Based Action Recognition

To comprehend the visual world, we should detect individual object instances in a scene as well as their interactions. In this regard, studies in this category initially recognize the instances in still images and then employ the visual relationship between human and objects to detect different actions. Authors in [31] make use of an inferred target location to find the correct object concerning the specific action. To be more specific, they estimate an action-type specific density around the portion of target objects by means of a human-centric recognition branch in the Faster R-CNN model. Equipped with the keypoint detection network, Wang et al. [32] compute an interaction vector based on the human and object center points to directly determine interactions between human-object pairs. Ma et al. [33] aim to promote features with a human-object relation module to calculate pair-wise interaction context between human and objects for action classification.

As human bodies are structural objects, modeling human activities by means of motion context is feasible. The motion context concerning entire body structure or different body parts can provide valuable geometric interactions in various human activities, which are specifically applicable for fine-grained activity recognition. In this vein, authors in [34] recognize semantic regions in the human bounding box and arrange features of detected regions in a top-down spatial order to strengthen inter-class variance for higher discriminative representation. Mottaghi et al. [35] propose a novel feature descriptor called Histogram of Graph Nodes (HGN), whereby the skeletons of the silhouettes can be converted into a graph to model the articulated human body skeleton. In [36], human body is partitioned into seven parts such as a head with a few semantic part actions to classify human action category. For this objective, a CNN model with two subnetworks is utilized for part localization and action prediction. Li

et al. [37] propose a Hierarchical Activity Graph to encode human instances and their body part to reason out the activity classes using part-level semantics. A unified CNN model is adopted in [22] to capture structural details of body instances and integrate several body structure cues, namely Structural Body Parts and Limb Angle Descriptor cues, for HAR.

While significant progress has been made in action recognition studies [22, 31–35, 35–37], most of which require high computational resources beyond the capabilities of edge devices. However, our two-phase framework is constructed from lightweight CNNs, which is specifically tailored for resource-constrained platforms. As such, under severe constraints on computing power and memory resource, the proposed SCLAR is much faster than its counterparts.

### *2.3. Data Imbalance Issue*

The imbalance class distribution would lead to poor generalization of CNNs, and the learning process with such a distribution would be biased toward over-represented classes. This eventually results in a biased CNN, which fails to detect the subtle discriminant features needed to classify the under-represented samples. Thus, the over-fitting in favor of the over-represented classes is inevitable for most approaches in action recognition. Though learning from imbalanced data has not been addressed in action recognition in still images, it has been well studied in other computer vision tasks such as object detection [38, 39] and image classification [23, 24, 40]. To mitigate this deep-rooted issue, several conventional and also recent studies have been conducted, such as over-sampling the under-represented samples and under-sampling the under-represented samples. Also, recent studies include Prime Sample Attention [41], AP Loss [42], DR Loss [43], pRoI Generator [44], and IoU-based Sampling [45]. While other studies in HAR have not considered the imbalanced data issue, in this paper, different classes are selected and categorized into super-classes to handle the

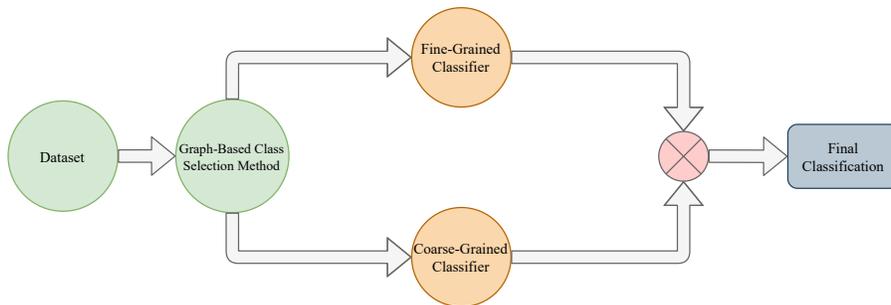


Figure 1: Flowchart of the proposed hierarchical action recognition.

data imbalance issue. With this perspective, the knowledge of different classes would be shared during the training procedure in order that the overall gradient in the training phase would not be dominated by the over-represented classes.

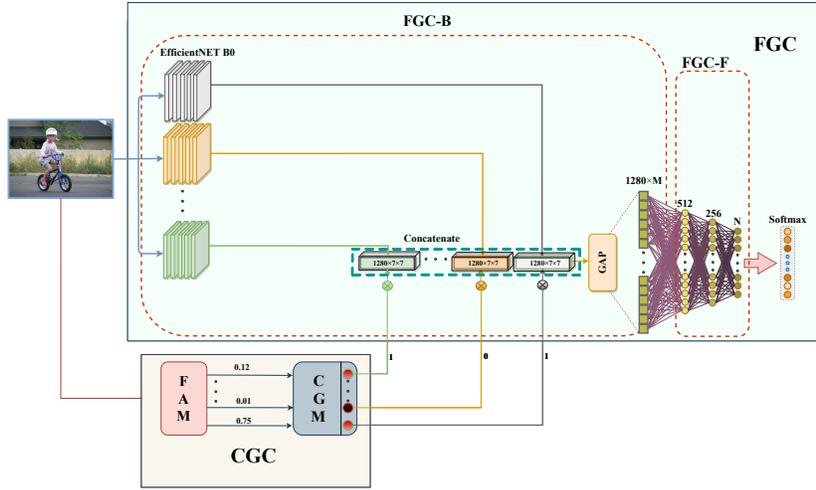
### 3. The Proposed Method

Recall that in this work, a bucket of separate lightweight CNN classifiers is required to be pre-trained for different super-classes. To this end, the entire dataset is first divided into  $M$  super-classes in such a way that all super-classes contain a relatively balanced distribution. Thanks to this dataset division, the previously condensed inter-class boundary would be smoothed. The EfficientNet-B0 model [46] is adopted as the lightweight CNN classifier in the fine-grained classification phase. The EfficientNet models are all pre-trained by means of input data in various classes of their relevant super-class. As such, they would become specialized in capturing domain-specific features in relation to the classes of their corresponding super-classes. A combination of these lightweight expert backbones constructs the fine-grained stage of our architecture configuration for the final action recognition prediction. Finally, the most relevant backbones for Fine-Grained Classifier (FGC) are activated by means of the Coarse-Grained Classifier (CGC). It should be noted that the performance of EfficientNet models would drop with a small number of super-classes (i.e.,

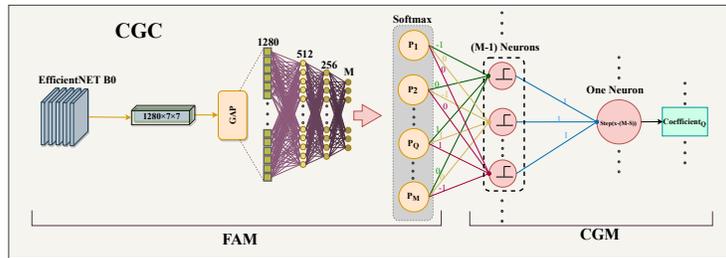
$M$ ). On the other side, a large number of super-classes would incur computation overhead. Thus, to achieve the optimal trade-off between classification performance and computational cost,  $M = 3, 4, 5$  super-classes are ablated in our assessments. Moreover, it should be stressed that the configuration of each super-class is one of the major factors in the performance of our two-phase architecture. Therefore, the GCS algorithm is proposed in this study to probe the best configuration for each super-class whereby the input embedding space is partitioned into several super-classes, and the inter-class variation would be increased. Figure 1 illustrates the flowchart of the hierarchical action recognizer in our method. First, the GCS algorithm pinpoints the best configuration of each super-class. Then, CGC selects the best pre-trained backbone expert and activates FGC. Finally, FGC determines human action recognition classes.

### 3.1. Fine-grained Phase

Figure 2 (a) depicts the entire architecture of FGC in the fine-grained classification phase. The FGC is composed of two components, namely FGC-B and FGC-F (B and F denote backbone and fully connected layer components). The FGC-B is constructed from  $M$  EfficientNet-B0 models, which can readily scale up conventional CNNs in a more principled manner to any resource limitations with a desired efficiency. For more details about the architecture of EfficientNet-B0, such as kernel size and channel size in convolution filters, we refer the readers to [46]. In the fine-grained classification phase, first, all EfficientNets are applied to the input image and each of which outputs a  $7 \times 7 \times 1280$  feature map. The specialized output maps corresponding to different EfficientNets are then concatenated channel-wise and fed into the CGC to be elected for the remaining section of the FGC. Then, the compact subnetwork (FGC-F) runs on the specialized feature map. The global averaging pooling (GAP) is also employed after the concatenation step to regularize the FGC and mitigate



(a)



(b)

Figure 2: The proposed architecture of our two-phase hierarchy. (a) Network architecture. (b) The coarse-grained classifier.

the over-fitting issue. Lastly, to train the FGC-F, the softmax loss is utilized.

### 3.2. Coarse-grained Phase

As represented in Figure 2 (a), to learn discriminative features for each super-class, CGC determines the pre-trained expert backbones for the relevant FGC-F in the fine-grained classification phase. The CGC empowers FGC to differentiate between subtle variations between different action classes. That is to say, it would help the FGC to strengthen inter-class discrepancies. This is due to the fact that the pre-trained backbones in the FGC-B have been sep-

arately tailored for different classes. For selecting the expert backbones, the output probabilities are produced by the Feature Attention Module (FAM), which is made up of EfficientNet-B0 and has been previously trained with softmax loss. The output probabilities of different experts over the corresponding super-classes are then converted to Top-S scores and S-hot masks that are pointwisely multiplied by the FGC-B. We opt for the Top-2 score rather than the Top-1 score since it hedges more relevant fine-grained backbones into a strong one and aggregates the most relevant subset-specific features. Also, the ablation study verifies the effectiveness of Top-2 score. The Coefficient Generation Module (CGM) addresses this task by assigning ones and zeros to the relevant and irrelevant fine-grained backbones in the FGC-B, respectively. To demonstrate this procedure, the FAM is applied to  $M$  super-classes. This results in the output probabilities  $P = \{P_1, P_2, \dots, P_M\}$ , where  $P_Q$  corresponds to the probability of super-class  $Q$  ( $1 \leq Q \leq M$ ). Given that FAM generates the probabilities for the selection process in the FGC-B and  $M$  denotes the number of all super-classes, the CGM would set output  $Q$  to one when  $P_Q$  is greater than at least  $M - S$  number of output probabilities. The number of probabilities less than  $P_Q$  can be computed by step function as follows:

$$u_Q = \sum_{\substack{v=1 \\ v \neq Q}}^M \text{step}(P_Q - P_v), \quad (1)$$

where  $u_Q$  denotes the number of probabilities which are surpassed by  $P_Q$ . This equation compares the probability of  $P_Q$  with other output probabilities  $P_v$  and accumulates all the outputs of step functions. Then, the resulting coefficient corresponding to  $P_Q$  (i.e.,  $c_Q$ ) would be set to one when  $u_Q$  exceeds  $M - S$  as

follows:

$$c_Q = \text{step}(u_Q - (M - S)) = \begin{cases} 1 & M - S < u_Q \\ 0 & M - S \geq u_Q \end{cases} \quad (2)$$

To formulate the general form of Equation 1, a weight  $(M - 1) \times 1$  matrix is defined as follows:

$$\text{step}(P_Q - P_v) = \begin{bmatrix} \text{step}(P_Q - P_1) \\ \text{step}(P_Q - P_2) \\ \vdots \\ \vdots \\ \text{step}(P_Q - P_M) \end{bmatrix}_{(M-1) \times 1}, \quad (3)$$

where each element draws a comparison between  $P_Q$  and other output probabilities.

Equation 3 can be also reformulated as follows:

$$\begin{aligned}
& \text{step}(P_Q - P_v) \\
= & \text{step} \left( \begin{matrix} \left[ \begin{array}{ccccccc} -1 & 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & -1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 1 & 0 & \dots & -1 \end{array} \right]_{(M-1) \times (M)} \times \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_Q \\ \vdots \\ P_M \end{bmatrix}_{(M) \times 1} \end{matrix} \right) \\
& = \text{step}(W_1^T \times P),
\end{aligned} \tag{4}$$

where  $W_1^T$  is a constant matrix which is multiplied by  $P$  and ultimately results in  $P_Q - P_v$ . To implement Equation 4, we use an MLP layer with  $(M - 1)$  neurons. The step function is adopted for the activation function of each neuron. Then, the weights of different neurons can be expressed as  $W_1^T$ , in which each row corresponds to a single neuron. Recall that the input of Equation 3 is the probabilities that have been produced by FAM. After comparing  $P_Q$  with all other output probabilities and applying the step function to them, the computed output can be assessed in comparison with  $M - S$  in Equation 2. Thus, by unifying Equations 2 and 4, we can finally obtain the final discrete coefficients as follows:

$$\begin{aligned}
c_Q &= \text{step} \left( \sum_{\substack{v=1 \\ v \neq Q}}^M \text{step}(P_Q - P_v) - (M - S) \right) \\
&= \text{step} \left( \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}_{1 \times (M-1)} \times \begin{bmatrix} \text{step}(P_Q - P_1) \\ \text{step}(P_Q - P_2) \\ \vdots \\ \text{step}(P_Q - P_M) \end{bmatrix}_{(M-1) \times 1} - (M - S) \right) \\
&= \text{step} \left( (W_2^T \times \text{step}(P_Q - P_v)) - (M - S) \right), \tag{5}
\end{aligned}$$

where  $c_Q$  indicates the final coefficient for  $P_Q$ , and  $W_2^T$  denotes an all-ones matrix. Same as Equation 4, we can implement Equation 5 by means of single-layer MLP with one neuron and  $\text{step}(U - (M - S))$  activation function. Figure 2 (b) illustrates the two-layer MLP network that performs the prementioned operations in CGM and outputs the discrete coefficients to select the best features in the FGC-B. Note that the two-layer MLP network is reused for individual M outputs.

### 3.3. Graph-based Class Selection Method

In action recognition in still images, there exist compact inter-classes boundaries among various activity classes. The proposed GCS method aims to partition the input dataset into M super-classes with relatively balanced distribution so that the inter-class variance would be strengthened desirably. To obtain the optimal configuration for each super-class in the GCS approach,  $N$  classes are divided into M super-classes according to their dependencies. To evaluate the inter-class dependency in a dataset and categorize correlated classes into distinct

super-classes, a baseline network (coined BN) similar to the fine-grained expert classifiers (EfficientNets) is first pre-trained over the input dataset. Then, we can calculate the required dependencies between various classes in our dataset using the BN. More specifically, the pre-trained BN is first applied to all training images in class  $j$ ; then, the output scores are sorted in ascending order. It is evident that the generated Top-1 scores correspond to class  $j$ . Then, the remaining  $i$ -th top scores for the majority of images in class  $j$  would be related to class  $c_i$ , where  $1 \leq c_i \leq N$ . The resulting  $i$ -th top scores ( $i > 1$ ) would determine the desired dependencies, which are required for partitioning different classes in the GCS algorithm. The second-order, third-order, and fourth-order dependencies for each class can be expressed as follows:

$$D_2 = \{d_{21}, d_{22}, d_{23}, \dots, d_{2N}\}, \quad (6)$$

$$D_3 = \{d_{31}, d_{32}, d_{33}, \dots, d_{3N}\}, \quad (7)$$

$$D_4 = \{d_{41}, d_{42}, d_{43}, \dots, d_{4N}\}, \quad (8)$$

where  $D_i$  denotes  $i$ -th dependency set in the input dataset, and  $d_{ij}$  indicates the  $i$ -th top score ( $i > 1$ ) corresponding to class  $j$ . Algorithm 1 summarizes the pseudocode for  $D_i$  computation. For each dependency set  $D_{i \in \{2,3,4\}}$  in the proposed GCS algorithm, the classes are represented as different nodes in a graph configuration, wherein the connections are based on the class dependencies and correlations. To model three dependency sets  $D_{i \in \{2,3,4\}}$  in the GCS algorithm, two types of directed connections, namely one-way or two-way connections, are employed. In the GCS configuration, the second-order dependency is modeled either with the one-way or two-way connection (as depicted in Figure 3). It

---

**Algorithm 1** GCS Algorithm

---

**Input:** Images, EfficientNet, Number of Classes  $N$ .

**Output:** Dependencies order  $D_o$ ,  $o \in \{2, 3, 4\}$

```
1: Pre-train EfficientNet with Input Images.
2:  $D_o = []$ 
3: for  $i = 1, 2, \dots, N$  do
4:    $C_j = [], j \in \{2, 3, 4\}$ 
5:    $ClassImgs = Images[i]$ .
6:   for  $img$  in  $ClassImgs$  do
7:      $score = EfficientNet(img)$ 
8:      $score = sort(score)$ 
9:      $C_r.add(index(score[r])), r \in \{2, 3, 4\}$ 
10:   $d_{m,i} = MaxIter(C_m), m \in \{2, 3, 4\}$ 
11:   $D_n.add(d_{n,i}), n \in \{2, 3, 4\}$ 
```

---

```
1: Function  $MaxIter(A)$ :
2:    $F = zeros(length(A))$ 
3:   for  $k$  in  $range(1, length(A))$ 
4:      $F[A[k]] = F[A[k]] + 1$ 
5:    $F = sort(F)$ 
6: Return  $index(F[1])$ 
```

---

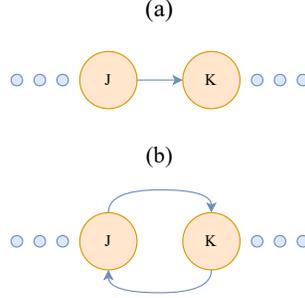


Figure 3: Two types of connections between different nodes in GCS algorithm. (a) one-way connection. (b) two-way connection.

is worth noting that the third-order and fourth-order dependencies can only be represented by the one-way connection (Figure 3). Once class  $j$  is dependant on class  $k$ , similarities between classes  $j$  and  $k$  can be formulated into four categories contingent upon the connection configuration, i.e., the order of dependencies, as below:

$$d_{2j} = k, d_{2k} = j \Rightarrow Type - 1 \text{ similarity} \Rightarrow S_1, \quad (9)$$

$$d_{2j} = k, d_{2k} \neq j \Rightarrow \text{Type} - 2 \text{ similarity} \Rightarrow S_2, \quad (10)$$

$$d_{3j} = k, d_{3k} \neq j \Rightarrow \text{Type} - 3 \text{ similarity} \Rightarrow S_3, \quad (11)$$

$$d_{4j} = k, d_{4k} \neq j \Rightarrow \text{Type} - 4 \text{ similarity} \Rightarrow S_4, \quad (12)$$

where  $d_{ij}$  represents the  $i$ -th top score corresponding to class  $j$ , which is dependant on class  $k$ .  $S_i$  is the status of similarities among different nodes. Clearly,  $S_{t \in \{1,2,3,4\}}$  influences the performance of BN with different scales. To gauge the impact of  $S_t$  on the final classification error, first, all class pairs with similarity  $S_t$  are extracted for the input dataset. Then, to calculate the classification error, we separately train a new EffitientNet for the selected class pairs with  $S_t$ . Finally, the average error for  $S_t$  would be computed. These operations are formulated as follows:

$$E_t = \frac{1}{|N_t|} \sum_{\ell=1}^{|N_t|} E_\ell, \quad (13)$$

$$E_\ell = \frac{1}{|K_\ell| + |J_\ell|} \sum_{\gamma=1}^{|K_\ell| + |J_\ell|} y_\gamma \times \log(\hat{y}_\gamma), \quad (14)$$

where  $E_t$  is the average error for newly trained backbones on the selected class pairs with  $S_{t \in \{1,2,3,4\}}$ .  $N_t$  denotes the number of selected class pairs at each similarity level  $S_t$ .  $\hat{y}_\gamma$  and  $y_\gamma$  indicate ground-truth labels and the predicted label generated by the newly trained backbones.  $|K_\ell|$  and  $|J_\ell|$  correspond to the number of images in the selected pairs  $\ell$ . Now, using Equations 13 and 14, we can compute the average error  $E_t$  for all similarities. The following inequality

summarizes the obtained results in our evaluations as:

$$E_1 > E_2 > E_3 > E_4 > E_0, \quad (15)$$

where  $E_0$  denotes the average classification error for the class pairs without any similarities. To partition the input space such that similar classes with negligible inter-class variation would be placed into different super-classes, the assigning process is carried out under the guidance of Equation 15. Thus, the number of class pairs with  $S_1$ ,  $S_2$ ,  $S_3$ , and  $S_4$  similarities would be minimized in the super-classes, respectively. For the partitioning phase in the GCS algorithm, since  $S_1$ , and  $S_2$  are concerned with  $D_2$  graphs, different classes in  $D_2$  are first assigned to  $M$  super-classes based on the structure of their connections. In this assignment stage, the number of  $S_1$  and  $S_2$  in each sub-class would be minimized with the first and second priorities. Afterward, concerning the neighboring nodes in the  $D_3$  graphs, the configuration of subclasses are reformed to minimize the number of  $S_3$  without constructing new  $S_1$ , and  $S_2$ . In the end, the same reconfiguration is conducted for the neighboring nodes in  $D_4$  graphs to minimize the number of  $S_4$  without constructing new  $S_1$ ,  $S_2$ , and  $S_3$ . For a better understanding of the GCS algorithm, we provide an example in Appendix A.

#### 4. IHAR dataset

Real-world data often contain a long-tailed and open-ended distribution, and the data required for action recognition systems is no exception. This is due to the fact that human instances are present at varying rates in nature. An action recognition system requires to classify various action types among under-represented and over-represented classes and also generalize from a few known instances in under-represented classes. Previous well-known datasets in this field have not paid much attention to this requirement and consequently could

easily struggle to keep working in long-tailed distribution. In this research, we introduce an imbalanced action recognition dataset with long-tailed distribution to fairly demonstrate the performance of our model in long-tailed distribution. Figure 4 illustrates the long-tailed distribution of the IHAR dataset. The IHAR dataset consists of 23854 images for 46 human action classes with a different number of samples. Different classes in the IHAR dataset are displayed in Figure 5. To construct the IHAR dataset, we have incorporated different small-scale and medium-scale datasets, including Stanford40 [17], Pascal VOC 2012 Action[18], BU101+ [19], PPMI [47], Sports [48], and ImSitu [49] datasets. Similar classes in the selected datasets are merged all together to form new classes in the IHAR dataset. For instance, the images relevant to the athletic movements in the selected datasets are integrated into the Sport class in the IHAR dataset. The number of classes chosen from each dataset is reported in Table 1.

## 5. Experiments

In this section, we initially describe the datasets, which are employed in our evaluations. Then, the implementation details and the exhaustive evaluation are addressed to prove the effectiveness of the proposed method. More specifically, subsection 5.1 introduces different well-organized balanced datasets, which are employed in our assessments along with the IHAR dataset. Subsection 5.2 provides the implementation details of the proposed action recognition approach. Subsection 5.3 investigates the contribution of different proposed components in our action recognition approach. The performance of the proposed SCLAR on ablated versions of our method is addressed as well. Finally, quantitative and qualitative experiments are carried out in subsection 5.4, and the comparison with the state-of-the-art studies is also provided in this section.

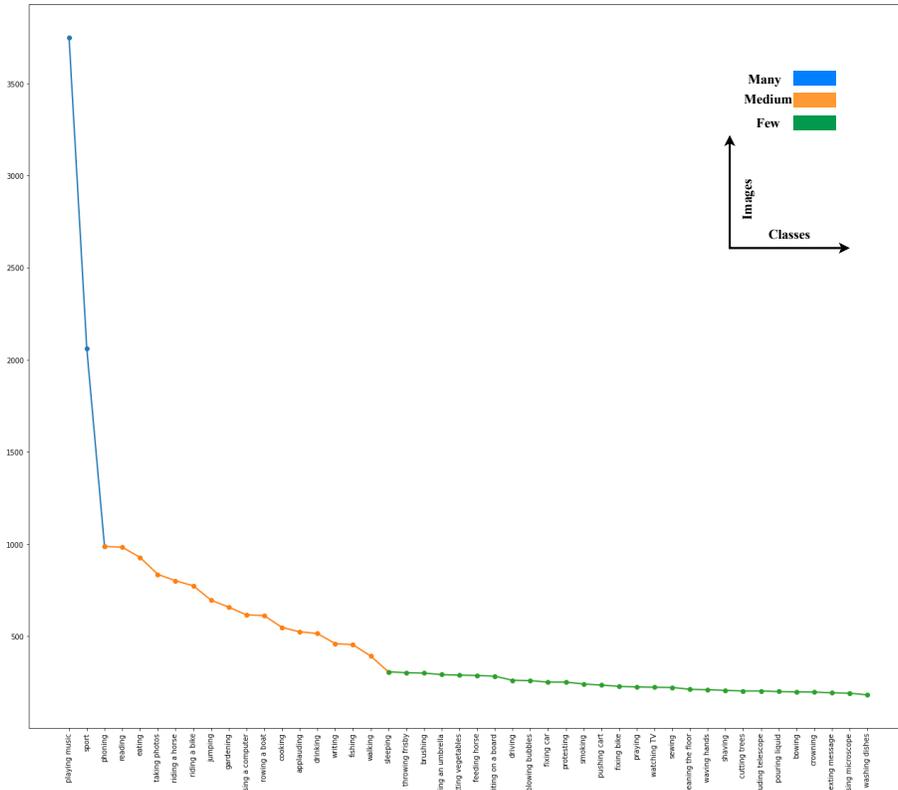


Figure 4: Distribution of different classes in the introduced IHAR dataset.

Table 1: The number of different classes in the selected datasets, which has been utilized for the IHAR dataset construction.

Dataset	Utilized classes
Stanford40 [17]	All classes
Pascal VOC Action [18]	All classes
BU101+ [19]	20 classes <sup>1</sup>
PPMI [47]	All classes
Sports [48]	All classes
ImSitu [49]	23 classes <sup>2</sup>

<sup>1</sup>The classes include baseball pitch, cutting in kitchen, frisbee, playing piano, playing violin, rafting, rowing, soccer penalty, taiChi, writing on board, basketball, biking, brushing teeth, playing cello, playing flute, playing sitar, walking with dog, horse riding, playing daf, playing guitar, and playing tabla.

<sup>2</sup>The classes include sleeping, gardening, reading, praying, eating, shaving, driving, cooking, protesting, crowning, eating, bowling, drinking, sewing, rowing a boat, writing on a book, applauding, taking photos, phoning, eating, gardening, sport, and fishing.



Figure 5: Different action classes in the introduced IHAR dataset.

### 5.1. Datasets

Apart from the introduced IHAR dataset, three other datasets (Stanford 40 [17], Pascal VOC 2012 Action [18], BU101+ [19]) are also employed in our exhaustive evaluations. In what follows, we briefly review the datasets separately. In addition, PPMI, Sports, and ImSitu datasets are explained as follows:

*Stanford 40.* This dataset [17] includes 9532 images in total corresponding to 40 classes of actions, which covers various real-world action classes such as riding a horse, waving hands, running, and shooting an arrow. Eleven action classes are relevant to human body motion, and the others are based on non-body human motion. In each action class, 100 images are selected for the training phase, and the others are utilized for the testing phase.

*Pascal VOC 2012 Action.* Pascal VOC [18] composes of 4588 images in 10 classes, which is divided into 2296 training images, 2292 test images, and 4569 validation images. The images in Pascal VOC 2012 Action cover a large number of activities, including working with a computer, running, riding a horse, taking photos, and playing an instrument. The number of images in each action class for training and testing operations ranges from 140 to 220.

*BU101+.* This dataset [19] is the improved version of the BU101 dataset. To be more specific, it improves the cardinality of more than 90 classes of the original BU101 dataset. In addition, different challenges like dealing with background clutter and confusing objects are also taken into account in the BU101+ dataset. This dataset is constructed from 10100 images with 101 different action classes.

*PPMI.* The PPMI dataset [47] is constructed from different images of humans that interact with various musical instruments such as violin, clarinet, guitar, harp, recorder, and flute.

*Sports.* The Sports dataset [48] contains a limited number of activity classes, including tennis forehand, croquet, tennis serve, and cricket batting. The classes are divided into 30 and 20 images for training and test phases. All images in the test and training sets are cropped and centered such that the persons of interest occupy a large proportion of the input images.

*ImSitu.* This dataset [49] includes 125000 images with 200000 distinctive situations. Each situation corresponds to one of 500 activities, 11000 objects, and 1700 roles. There are lots of images in this dataset that are suitable for human action recognition. We extracted several classes in this dataset to build the IHAR dataset.

### 5.2. Implementation details

The proposed SCLAR contains two-phase training. In the first stage, the FGC-B and FAM components are required to be trained separately. The architecture of the FGC-B and FAM components follows the EfficientNET-B0 network, which has been pre-trained on the ImageNet [50]. Each branch of FGC-B is then trained for 100 epochs on an individual super-class which has been previously formed. In addition, the FAM component is trained for 100 epochs by means of all super-classes to distinguish different super-classes from each other at the coarse-grained phase. To train FGC-B, and FAM, we use Stochastic Gradient Descent (SGD) optimization with a momentum of 0.9, a weight decay of 20, and a batch size of 40. To train the FGC-B and FAM components, we set the learning rate to 0.5 and 0.1 in the Pascal VOC 2012 Action dataset and the remaining datasets, respectively. In the second phase, different branches of FGC-B are weighted by CGC and fed to the FGC-F component. Then, the FGC-F are trained from scratch same as FGC-B to determine the final output for action recognition. To train FGC-B, FAM and FGC-F components from scratch, the input images are first resized to  $224 \times 224$  pixels. Then, the overall architecture is fine-tuned with the input images with  $448 \times 448$  pixels. Also, two types of data augmentations are employed, including random horizontal flipping and random cropping. The optimum number of super-classes  $M$  is set to 4 in our experiments. All the experiments are conducted on a single NVIDIA Geforce GTX 2080 TI GPU with 11 GB memory and the PyTorch toolbox.

### 5.3. Ablation study

In this section, we conduct extensive ablation experiments to demonstrate the effectiveness of the key modules proposed in our SCLAR. All experiments

in this part have been carried out on the Stanford40 dataset, which is the most widely benchmarked dataset in still image action recognition.

*Effect of GCS method for Super-Class Division.* To investigate the impact of the generated super-classes on the GCS method, we set the number of super-classes to  $M = 4$  and assess the performance of each expert backbone in the FGC-B separately in comparison with the baseline model. Note that the baseline model is a single EfficientNet-B0 network which has been pre-trained over all classes. The results of this evaluation are reported in Table 3. The results demonstrate that each expert backbone yields better performance over the dedicated super-classes compared to the baseline model over all classes. This is largely ascribed to the fact that each expert backbone can extract more discriminative features from each super-class compared with the baseline model. In addition, to visually investigate this contribution, Figure 6 illustrates the tSNE visualization for different variations of expert numbers. It is observed that there exists low inter-class variation and no explicit boundaries for the baseline classification model within many classes; yet, the pre-trained expert backbones enjoy more discriminative representations and consequently can easily distinguish different action classes from each other. To explore the performance of different expert backbones for different inputs, we visualize the gradients of the top-class predictions in Figure 7 with respect to various input images by means of Smooth Grad-CAM++ [51]. The green checkmarks denote the expert backbones which are selected by CGC. The red cross marks represent the other irrelevant expert backbones. As shown in Figure 7, for the *drinking* sample in the first row, the top selected expert (#4) highlights the most informative image regions crucial to the proper classification. The other experts, on the other hand, concentrate on the areas uncorrelated with scenes and action-related objects; thus, the proposed CGC attempts to alleviate their impacts. In addition, in the third row,

Table 2: The impact of GCS method for finding the optimum configuration for different super-classes in the Stanford40 dataset.

Super-Class Division Method	mAP (%)
Random Division	88.7
GCS Division	92.86

although the top expert (#2) focuses on the bottles as well as the cell phone, the integration of all experts (the main proposed architecture) captures more discriminative areas with tighter bounds.

We also launch a study to benchmark the defectiveness of Equation 15 for finding the optimum configuration for all super-classes. In this study, we form the configuration of different super-classes randomly multiple times through which we assess the performance of action classification and report the average accuracy for all experiments. Then, we draw a comparison in Table 2 between the results obtained with randomly configured super-classes and the proposed GCS method. Table 2 validates that the proposed GCS method could find the optimum super-classes arrangement, thereby boosting the action recognition performance.

*Effect of CGC Method.* To demonstrate the impact of CGC in our main architecture, we deactivate this component and concatenate the output of expert backbones for the final classification. In this experiment, in the absence of CGC, different number of expert backbones is also ablated, as represented in Table 4. The results verify that CGC manages to assign the best expert to the relevant FGC in the fine classification phase. Interestingly, we observe that without CGC, multi-expert architecture lags behind the baseline model, and even the larger number of experts degrades the classification performance.

*Effect of  $M$  and  $S$ .* In this experiment, the number of super-classes  $M$  along with the threshold of the decision  $S$  in CGM (Equation 2) are investigated. Various experiments are performed with  $M = 3, 4, 5$  and  $S = 1, 2, 3$ . As reported in

Table 3: The impact of GCS method on different expert backbones for generating specialized super-classes in the Stanford40 dataset.

Model	Dataset	mAP (%)
Baseline	All classes	83.5
Expert 1	Super-class 1	93.5
Expert 2	Super-class 2	93.9
Expert 3	Super-class 3	94.2
Expert 3	Super-class 4	93.2

Table 4: The impact of CGC method in the main architecture on action classification in the Stanford40 dataset.

Model	Number of Experts	CGC	mAP (%)
Baseline	✗	✗	83.5
Multi-Expert	3	✗	82.4
Multi-Expert	4	✗	81.5
Multi-Expert	5	✗	79.63
Multi-Expert	4	✓	92.86

Table 5: The impact of  $M$  (the number of super-classes) and  $S$  (the threshold in Equation 2 in CGM) for action classification. The reported numbers in the table are the mAP result for action recognition in the Stanford40 dataset

	S = 1	S = 2	S = 3
M = 3	81.36	89.6	82.4
M = 4	82.4	92.86	85.6
M = 5	81.95	90.4	83.5

Table 5,  $M = 4$  and  $S = 2$  provide the best results in term of mAP metric. This study indicates that the larger values of  $M$  and  $S$  would not necessarily promote the performance of the proposed model.

#### 5.4. Comparisons with state-of-the-art methods

In this section, quantitative and qualitative experiments are performed to assess the contributions of SCLAR and prove its superiority compared with other studies. To this end, different benchmark datasets are adopted, including Stanford40, IHAR, BU101+, and Pascal VOC 2012 Action datasets.

Table 6: Comprehensive results of the proposed method compared to the state-of-the-art action recognition works. The results are in terms of mean average precision (mAP) on the Stanford40 dataset.

Method	mAP (%)	Bounding Box	Pose Estimator	Object Detector	Number of Parameters (M)
Minimum Annotation [30]	82.64				24 >
Attention-Joints Graph [52]	84.6	✓			24 >
Deep VLAD Spatial Pyramids [29]	88.5			✓	24 >
Pose-Guided [12]	89.53		✓	✓	43 >
BSE [22]	90.4		✓		43 >
Multi-branch Attention [21]	90.7		✓		30 >
Ensemble [10]	90.83				80 >
VIT [53]	91.24				85 >
Semantic Body Part [36]	91.2		✓		34 >
Loss Guided [20]	91.2		✓		51 >
Body-Part-aware [54]	91.91		✓		64 >
Multi-Attention Guided [19]	92.45	✓			44 >
SCLAR	92.86				20

#### 5.4.1. Stanford40

Table 6 reports the mean average precision results of our action recognition method compared to its competitors on the Stanford40 dataset. To make an apple-to-apple comparison in still image action recognition, we also report different types of auxiliary components utilized in different studies in Table 6. The state-of-the-art studies in still image action recognition mainly utilize auxiliary components such as a pose estimator or object detector to provide the action classifiers with detailed information about different agents in the environment. Though the auxiliary components promote an accurate action recognition, they detrimentally make their approaches restricted to domain-specific data and also induce a heavy computation burden. As reported in Table 6, SCLAR yields the best performance in terms of accuracy and efficiency in comparison with state-of-the-art works. More specifically, as a runner-up, SCLAR yields gains of 0.41% and 0.66% on mAP compared to the second [19] and third [54] best methods. While the obtained gains in respect to the second-best method [19] is not significantly pronounced, the efficiency gain is dramatically considerable. Thus, SCLAR demonstrates a better accuracy/speed trade-off generally. Moreover, SCLAR relaxes the requirement of human bounding boxes in still images.

Yet, [19] and [54] require additional supervision, such as human bounding boxes or predefined body parts, which confine their practical applications to specific domains with available bounding boxes. Besides, it is proven that SCLAR outperforms other studies [12, 20–22] to a great degree. The primary reason behind such noticeable superiority is that SCLAR could discriminate subtle visual inter-class variations from intra-class ones for action recognition task thanks to the CGC and FGC components. Note that although other studies [12, 12, 20–22, 29] enjoy an additional object detector or pose estimator to localize body keypoints and filter out noisy context, they considerably lag behind SCLAR in both accuracy and efficiency metrics.

To qualitatively assess the proposed method compared to its counterpart, the gradients of top-class predictions are illustrated for different samples in Figure 8. It can be observed that the proposed method extracts tighter and more relevant regions of human agents in comparison with other studies. For instance, for *phoning*, and *reading* samples, the proposed SCLAR captures the action-related interactive objects (book and cell phone) accurately and attends less to the background regions. Moreover, not all areas in an image are helpful for action recognition. Irrelevant details captured by feature extraction phase in Loss Guided [20] and Multi-Branch [21] hurt the performance of action recognition in their works. In challenging samples such as *applauding*, SCLAR pays less attention to distractor objects. We highlight that distractor objects belong to irrelevant action classes in input images, thereby inducing misclassification.

In addition, to evaluate the impact of SCLAR on different challenging classes with low inter-class variation, we compute average precision for all classes in the Stanford40 dataset. Figure 9 depicts the average precision (AP) scores of SCLAR along with single EffitientNet and the state-of-the-art study over all classes in the Stanford40 dataset. We analyze the least AP scores for different

classes as well as their second-dependency classes. This analysis reveals that the proposed SCLAR gains substantial improvements over the baseline EffitientNet on the most challenging classes, which contain low inter-class variation with their second-class dependency classes. For instance, *purring liquid*, *phoning*, *texting message*, *taking photo*, and *waving hand* classes possess low inter-class variations with their second dependency classes, which include *washing dishes*, *taking photo*, *smoking*, *phoning*, and *smoking*, respectively. The proposed SCLAR outperforms the baseline EffitientNet (and also the state-of-the-art study) over these classes more significantly. In short, taking these results into account, we can substantiate the second contribution of this study which claims that SCLAR surpasses its competitors in HAR while requiring much less computational cost.

#### 5.4.2. PASCAL VOC 2012 Action

Table 7 presents the performance of SCLAR on the PASCAL VOC 2012 Action dataset. As reported in Table 7, the superiority of SCLAR is kept with much less computational cost even compared to various state-of-the-art studies [19, 22], which leverage from auxiliary supervision such as annotated bounding boxes or auxiliary modules such as pose estimator. To be more specific, the Multi-Attention Guided [19] approach adopts a teacher-student knowledge distillation which requires annotated bounding boxes and huge computation cost. In addition, the BSE [22] extracts body structure cues such as structural body parts and limb angle descriptors from local and global perspectives by means of a complicated three-branch classifier.

#### 5.4.3. BU101+

To gauge the generalization of SCLAR, we also evaluate SCLAR on the BU101+ dataset. This dataset offers new challenges such as background clutter and confusing human agents for the action recognition evaluation in compar-

Table 7: Comprehensive results of the proposed method compared to the state-of-the-art action recognition works. The results are in terms of mean average precision (mAP) on the PASCAL VOC 2012 Action dataset.

Method	mAP (%)	Bounding Box	Pose Estimator	Number of Parameters (M)
Minimum Annotation [30]	83.23			24 >
SAAM-Nets [55]	84.8			25 >
Multi-branch Attention [21]	87.1		✓	30 >
R*CNN [56]	90.4			24 >
Multi-Attention Guided [19]	91.51	✓		44 >
BSE [22]	91.8		✓	43 >
SCLAR	92.46			20

Table 8: Comprehensive results of the proposed method compared to the state-of-the-art action recognition works. The results are in terms of mean average precision (mAP) on the BU101+ dataset.

Method	mAP (%)	Number of Parameters (M)
ResNET_18 [57]	83.06	11.5
ResNET_34 [57]	87.07	21.6
ResNET_50 [57]	87.44	24
Multi-Attention Guided [19]	90.16	44 >
SCLAR	92.27	20

iosn with the standard Stanford40 dataset. Since BU101+ has been introduced recently, limited research works have been benchmarked on this dataset. As such, we suffice to pursue our comparisons the same as [19]. Based on Table 8, SCLAR obtains better results than Multi-Attention Guided research work [19] even when [19] accesses to annotated bounding boxes. Considering reported results in Table 8, we can demonstrate the effectiveness of the proposed GCS method in SCLAR to strengthen the inter-class variation.

#### 5.4.4. IHAR

In the last part of our evaluations, we seek to study whether the proposed SCLAR can handle the data imbalance issue for action recognition with long-tailed distribution in still images. Table 9 depicts our comparisons on the IHAR dataset. The proposed lightweight SCLAR achieves 92.27 mAP accuracy and remarkably outperforms the state-of-the-art models with a gain of 8.92% com-

Table 9: Comprehensive results of the proposed method compared to the state-of-the-art action recognition works. The results are in terms of mean average precision (mAP) on the introduced IHAR dataset.

Method	mAP (%)	Number of Parameters (M)
ResNET_50 [57]	87.44	24
Multi-branch Attention [21]	81.9	44 >
R*CNN [56]	83.02	24 >
Ensemble [10]	83.35	80
SCLAR	92.27	20

pared to the second-best study [10]. The reason behind such a pronounced gain is that the GCS would divide the input dataset into different super-classes with relatively balanced distribution such that different backbone experts can learn better from under-represented classes in super-classes. A comparison between our evaluations on the IHAR dataset (Table 9) with the other balanced datasets in Table 7 and Table 6 reveals that the superiority of SCLAR over the state-of-the-art studies such as Ensemble [10] is more noticeable on long-tailed distribution than conventional balance datasets. All in all, we can deduce that the contributions of GCS and CGC in dealing with data imbalance issue in still image action recognition with long-tailed distribution are fully corroborated by the reported results in Table 8.

## 6. Conclusion

In this paper, we propose a two-phase multi-expert architecture for still image action recognition, which contains fine-grained and coarse-grained phases. The fine-grained stage includes a combination of pre-trained experts, and the coarse-grained stage contains FAM to select the relevant fine-grained features from the fine-grained phase. To the best of our knowledge, the proposed two-phase multi-expert architecture is the first action recognition study that adaptively hedges features from pre-trained experts. Moreover, the GCS method is

introduced to divide the input dataset into M super-classes for different experts in the fine-grained phase such that the inter-class variance increases. We also introduce a new still image action recognition dataset with long-tailed distribution to assess the robustness of our method against the data imbalance issue. Our experiments verify the effectiveness of different proposed components in our action recognition method and demonstrate its superiority compared with the state-of-the-art studies.

### Appendix A. Example of GCS

In this appendix, we present an example to illustrate GCS algorithm for partitioning various classes in our dataset. Let our dataset include different classes as  $Classes = \{A, B, C, D, E, F, G, H, I, J, K, L, M, N, O\}$ . Then, the second-order, third-order, and fourth-order dependency sets for our dataset can be expressed as follows:

$$D_2 = \{B, C, E, C, F, J, F, F, F, K, J, K, L, O, N\}, \quad (A.1)$$

$$D_3 = \{J, N, G, B, A, H, L, N, D, M, A, E, A, K, J\}, \quad (A.2)$$

$$D_4 = \{M, C, E, A, L, J, B, N, G, K, D, N, E, O, E\}, \quad (A.3)$$

Figure A.10 illustrates the constructed derivative graphs. To find the optimum configuration of different super-classes, the following steps are taken:

- Finding the center node in second order graphs.
- Numbering all routes in the second order graphs (Figure A.11 shows the outcome of the first and second steps).

- In the second order graphs ( $D_2$ ), according to the route score, the connected nodes are placed into the sequential super-classes.
- Optimizing the arrangement of classes in the super-classes according to  $D_2$  such that the number of class pairs with  $S_1$ , and  $S_2$  similarities would be minimized in each super-class.
- Optimizing the arrangement of classes in the super-classes according to  $D_3$  such that the number of class pairs with  $S_3$  similarity would be minimized, and no new  $S_1$ , and  $S_2$  similarities would be created.
- Optimizing the arrangement of classes in the super-classes according to  $D_4$  such that the number of class pairs with  $S_4$  similarity would be minimized, and no new  $S_1$ ,  $S_2$ , and  $S_3$  similarities would be created.

For this example, once there exist three super-classes, i.e.,  $M = 3$ , the arrangement of classes for the super-classes would be as Figure A.12. In the classes configuration, classes  $I$  and  $F$  contain  $S_2$  similarity. Since rearranging each of  $I$  and  $F$  classes with other classes would lead to a new  $S_2$  similarity with other classes, we would not change the configuration of these classes. However, pair classes such as  $(N, B)$ ,  $(D, K)$  and  $(B, G)$  in the second and third super-classes contain  $S_3$  and  $S_4$  similarities, respectively. Rearranging these classes with others would properly eliminate their similarities.

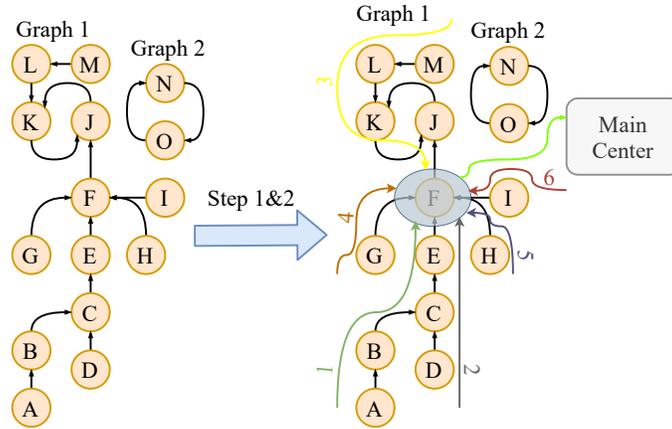


Figure A.11: Illustration of the first and second steps in the GCS algorithm. The main center node has the largest number of connections. Node F is the main center node in this figure. In the second order graph ( $D_2$ ), a set of connected nodes which end with the center node and begin from a boundary node is defined as a route.  $ABCEF$ ,  $GF$ , and  $MLKJ$  exemplify different routes which end with the main center node  $F$ .

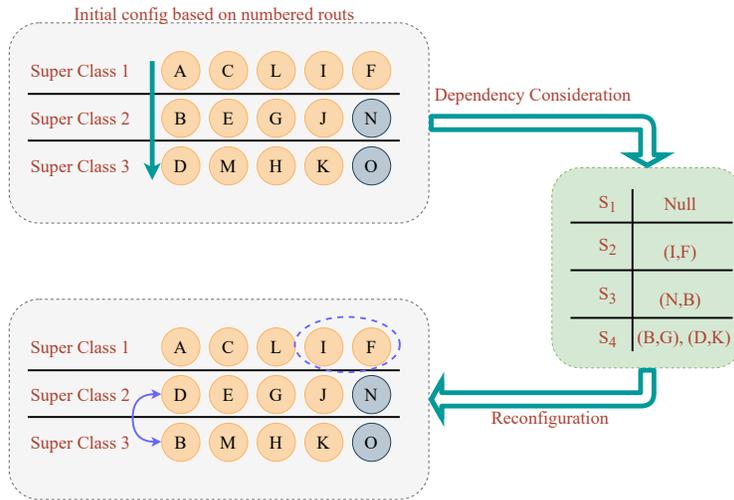


Figure A.12: The arrangement of classes for the super-classes with  $M = 3$ .

## References

- [1] C. Si, Y. Jing, W. Wang, L. Wang, T. Tan, Skeleton-based action recognition with hierarchical spatial reasoning and temporal stack learning net-

- work, *Pattern Recognition* 107 (2020) 107511.
- [2] L. Wu, Z. Yang, M. Jian, J. Shen, Y. Yang, X. Lang, Global motion estimation with iterative optimization-based independent univariate model for action recognition, *Pattern Recognition* 116 (2021) 107925.
- [3] H. Wang, B. Yu, J. Li, L. Zhang, D. Chen, Multi-stream interaction networks for human action recognition, *IEEE Transactions on Circuits and Systems for Video Technology* (2021).
- [4] T. Özyer, D. S. Ak, R. Alhajj, Human action recognition approaches with video datasets—a survey, *Knowledge-Based Systems* 222 (2021) 106995.
- [5] J. Dong, W. Yang, Y. Yao, F. Porikli, Knowledge memorization and generation for action recognition in still images, *Pattern Recognition* 120 (2021) 108188.
- [6] Y. Ji, Y. Zhan, Y. Yang, X. Xu, F. Shen, H. T. Shen, A context knowledge map guided coarse-to-fine action recognition, *IEEE Transactions on Image Processing* 29 (2019) 2742–2752.
- [7] S. Herath, B. Fernando, M. Harandi, Using temporal information for recognizing actions from still images, *Pattern Recognition* 96 (2019) 106989.
- [8] Z. Zheng, G. An, D. Wu, Q. Ruan, Spatial-temporal pyramid based convolutional neural network for action recognition, *Neurocomputing* 358 (2019) 446–455.
- [9] S. K. Yadav, K. Tiwari, H. M. Pandey, S. A. Akbar, A review of multi-modal human activity recognition with special emphasis on classification, applications, challenges and future directions, *Knowledge-Based Systems* (2021) 106970.

- [10] S. Mohammadi, S. G. Majelan, S. B. Shokouhi, Ensembles of deep neural networks for action recognition in still images, in: 2019 9th International Conference on Computer and Knowledge Engineering (ICCKE), IEEE, 2019, pp. 315–318.
- [11] D.-J. Kim, X. Sun, J. Choi, S. Lin, I. S. Kweon, Detecting human-object interactions with action co-occurrence priors, in: Proceedings of the European Conference on Computer Vision, Springer, 2020, pp. 718–736.
- [12] S. Mi, Y. Zhang, Pose-guided action recognition in static images using lie-group, *Applied Intelligence* (2021) 1–9.
- [13] M. Safaei, P. Balouchian, H. Foroosh, Ucf-star: A large scale still image dataset for understanding human actions, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, 2020, pp. 2677–2684.
- [14] Z. Xiao, X. Xu, H. Xing, F. Song, X. Wang, B. Zhao, A federated learning system with enhanced feature extraction for human activity recognition, *Knowledge-Based Systems* 229 (2021) 107338.
- [15] C. Plizzari, M. Cannici, M. Matteucci, Skeleton-based action recognition via spatial and temporal transformer networks, *Computer Vision and Image Understanding* 208 (2021) 103219.
- [16] Y. Yoshikawa, Y. Shigeto, A. Takeuchi, Metavd: A meta video dataset for enhancing human action recognition datasets, *Computer Vision and Image Understanding* (2021) 103276.
- [17] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, L. Fei-Fei, Human action recognition by learning bases of action attributes and parts, in: Proceedings of the IEEE International Conference on Computer Vision, IEEE, 2011, pp. 1331–1338.

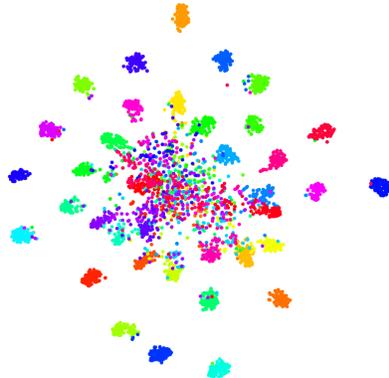
- [18] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, *International Journal of Computer Vision* 111 (2015) 98–136.
- [19] S. S. Ashrafi, S. B. Shokouhi, A. Ayatollahi, Action recognition in still images using a multi-attention guided network with weakly supervised saliency detection, *Multimedia Tools and Applications* (2021) 1–27.
- [20] L. Liu, R. T. Tan, S. You, Loss guided activation for action recognition in still images, in: *Asian Conference on Computer Vision*, Springer, 2018, pp. 152–167.
- [21] S. Yan, J. S. Smith, W. Lu, B. Zhang, Multibranch attention networks for action recognition in still images, *IEEE Transactions on Cognitive and Developmental Systems* 10 (2018) 1116–1125.
- [22] Y. Li, K. Li, X. Wang, Recognizing actions in images by fusing multiple body structure cues, *Pattern Recognition* 104 (2020) 107341.
- [23] S. Suh, H. Lee, P. Lukowicz, Y. O. Lee, Cegan: Classification enhancement generative adversarial networks for unraveling data imbalance problems, *Neural Networks* 133 (2021) 69–86.
- [24] B. Kim, Y. Ko, J. Seo, Novel regularization method for the class imbalance problem, *Expert Systems with Applications* (2021) 115974.
- [25] Z. Yang, Z. Wang, L. Luo, H. Gan, T. Zhang, Sws-dan: Subtler ws-dan for fine-grained image classification, *Journal of Visual Communication and Image Representation* (2021) 103245.
- [26] Y. Li, B. Fan, W. Zhang, W. Ding, J. Yin, Deep active learning for object detection, *Information Sciences* (2021) 418–433.

- [27] Y. Wang, X. Wei, L. Luo, W. Wen, Y. Wang, Robust rgb-d tracking via compact cnn features, *Engineering Applications of Artificial Intelligence* 96 (2020) 103974.
- [28] A. Rosenfeld, S. Ullman, Action classification via concepts and attributes, in: 2018 24th International Conference on Pattern Recognition (ICPR), IEEE, 2018, pp. 1499–1505.
- [29] S. Yan, J. S. Smith, B. Zhang, Action recognition from still images based on deep vlad spatial pyramids, *Signal Processing: Image Communication* 54 (2017) 118–129.
- [30] Y. Zhang, L. Cheng, J. Wu, J. Cai, M. N. Do, J. Lu, Action recognition in still images with minimum annotation efforts, *IEEE Transactions on Image Processing* 25 (2016) 5479–5490.
- [31] G. Gkioxari, R. Girshick, P. Dollár, K. He, Detecting and recognizing human-object interactions, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8359–8367.
- [32] T. Wang, T. Yang, M. Danelljan, F. S. Khan, X. Zhang, J. Sun, Learning human-object interaction detection using interaction points, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4116–4125.
- [33] W. Ma, S. Liang, Human-object relation network for action recognition in still images, in: 2020 IEEE International Conference on Multimedia and Expo (ICME), IEEE, 2020, pp. 1–6.
- [34] Z. Zhao, H. Ma, X. Chen, Semantic parts based top-down pyramid for action recognition, *Pattern Recognition Letters* 84 (2016) 134–141.

- [35] A. Mottaghi, M. Soryani, H. Seifi, Action recognition in freestyle wrestling using silhouette-skeleton features, *Engineering Science and Technology, an International Journal* 23 (2020) 921–930.
- [36] Z. Zhao, H. Ma, S. You, Single image action recognition using semantic body part actions, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3391–3399.
- [37] Y.-L. Li, L. Xu, X. Liu, X. Huang, Y. Xu, S. Wang, H.-S. Fang, Z. Ma, M. Chen, C. Lu, Pastanet: Toward human activity knowledge engine, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 382–391.
- [38] A. Bria, C. Marrocco, F. Tortorella, Addressing class imbalance in deep learning for small lesion detection on medical images, *Computers in biology and medicine* 120 (2020) 103735.
- [39] Y. Li, T. Wang, B. Kang, S. Tang, C. Wang, J. Li, J. Feng, Overcoming classifier imbalance for long-tail object detection with balanced group softmax, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10991–11000.
- [40] S. Suh, P. Lukowicz, Y. O. Lee, Discriminative feature generation for classification of imbalanced data, *Pattern Recognition* (2021) 108302.
- [41] Y. Cao, K. Chen, C. C. Loy, D. Lin, Prime sample attention in object detection, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11583–11591.
- [42] K. Chen, J. Li, W. Lin, J. See, J. Wang, L. Duan, Z. Chen, C. He, J. Zou, Towards accurate one-stage object detection with ap-loss, in: *Proceedings*

- of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 5119–5127.
- [43] Q. Qian, L. Chen, H. Li, R. Jin, Dr loss: Improving object detection by distributional ranking, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2020, pp. 12164–12172.
- [44] K. Oksuz, B. C. Cam, E. Akbas, S. Kalkan, Generating positive bounding boxes for balanced training of object detectors, in: Proceedings of the IEEE Winter Conference on Applications of Computer Vision, 2020, pp. 894–903.
- [45] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, D. Lin, Libra r-cnn: Towards balanced learning for object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 821–830.
- [46] M. Tan, Q. Le, Efficientnet: Rethinking model scaling for convolutional neural networks, in: International Conference on Machine Learning, PMLR, 2019, pp. 6105–6114.
- [47] B. Yao, L. Fei-Fei, Grouplet: A structured image representation for recognizing human and object interactions, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2010, pp. 9–16.
- [48] A. Gupta, A. Kembhavi, L. S. Davis, Observing human-object interactions: Using spatial and functional compatibility for recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 31 (2009) 1775–1789.
- [49] M. Yatskar, L. Zettlemoyer, A. Farhadi, Situation recognition: Visual semantic role labeling for image understanding, in: Proceedings of the

- IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 5534–5542.
- [50] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Ieee, 2009, pp. 248–255.
- [51] D. Omeiza, S. Speakman, C. Cintas, K. Weldermariam, Smooth grad-cam++: An enhanced inference level visualization technique for deep convolutional neural network models, 2019.
- [52] T. Ahmad, H. Mao, L. Lin, G. Tang, Action recognition using attention-joints graph convolutional neural networks, *IEEE Access* 8 (2019) 305–313.
- [53] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv:2010.11929 (2020).
- [54] B. Bhandari, G. Lee, J. Cho, Body-part-aware and multitask-aware single-image-based action recognition, *Applied Sciences* 10 (2020) 1531.
- [55] Y. Zheng, X. Zheng, X. Lu, S. Wu, Spatial attention based visual semantic learning for action recognition in still images, *Neurocomputing* 413 (2020) 383–396.
- [56] G. Gkioxari, R. Girshick, J. Malik, Contextual action recognition with r\* cnn, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 1080–1088.
- [57] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.



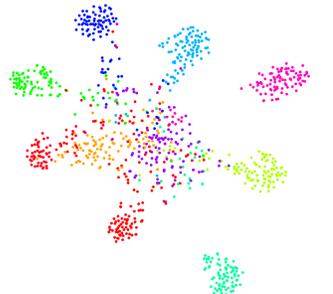
(a) tSNE visualisation for the baseline model on the entire dataset.



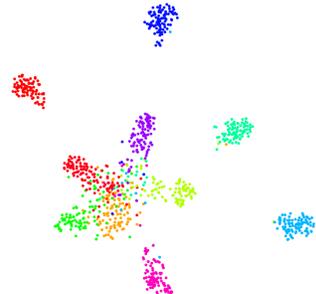
(b) tSNE visualisation for the first expert on the first super-class.



(c) tSNE visualisation for the second expert on the second super-class.



(d) tSNE visualisation for the third expert on the third super-class.



(e) tSNE visualisation for the fourth expert on the fourth super-class.

Figure 6: tSNE visualisation for the features learned by the baseline model and different experts in the main architecture.

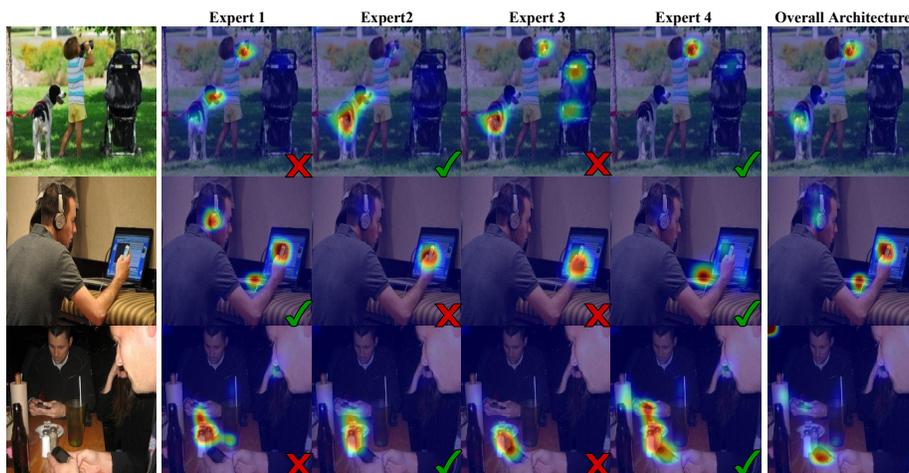


Figure 7: The heatmaps computed from different expert backbones and the main architecture with CGC.

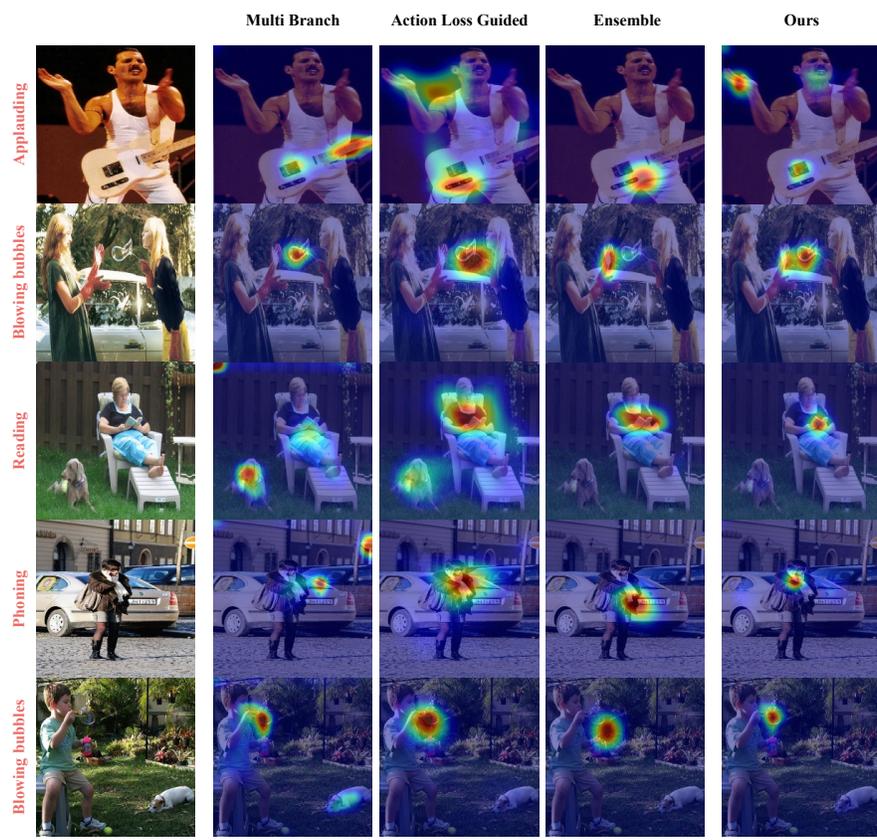


Figure 8: The heatmaps computed from different state-of-the-art studies and the proposed SCLAR.

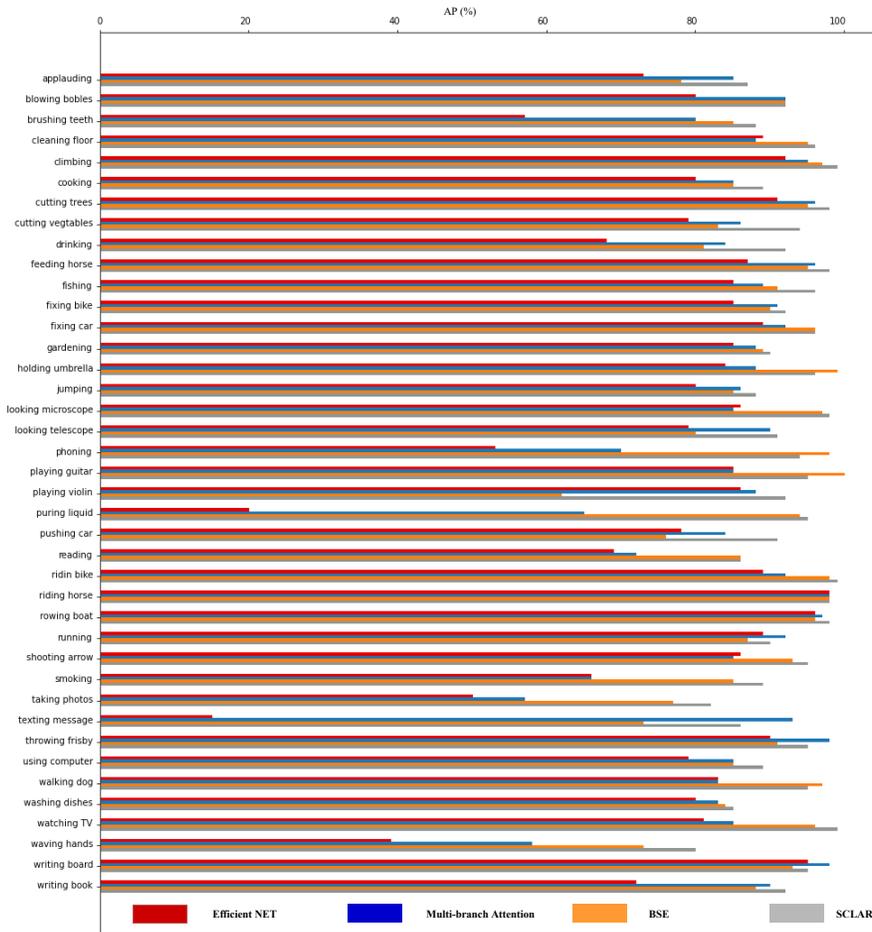


Figure 9: Results of the proposed SCLAR in comparison with the baseline EfficientNet and Multi-branch Attention [21] on different classes of the Stanford40 dataset.

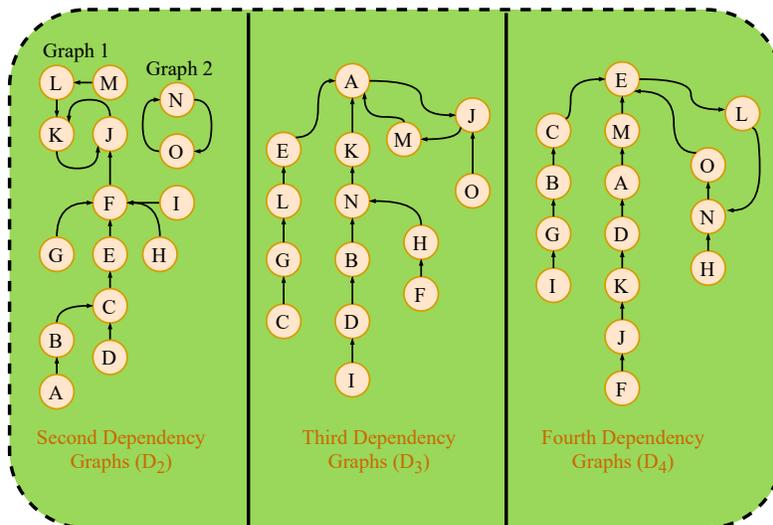


Figure A.10: Different dependency graphs in the proposed GCS algorithm.