

25 algorithm is superior in both accuracy and robustness compared with the state of the
26 arts.

27 **Keywords:** Large-scale image classification; Spatial attention module; Object
28 decision making; Channel attention module; Explainable AI

29 **1 Introduction**

30 In computer vision, image classification is a subject worthy of long-time research,
31 and it is an important foundation in some fields such as object detection (Diba et al.,
32 2017; Ren et al., 2017), face recognition (Girshick et al., 2014), pose estimation
33 (Chen et al., 2019; Chen et al., 2021), population density estimation(Zhang et al.,
34 2012; Chen et al., 2021), image segmentation (Chen et al., 2020 and Huang et al.,
35 2019) etc. Therefore, image classification has always been a hot topic (Lin et al., 2019
36 and Lin et al., 2021). At present, the image classification algorithm mainly relies on
37 the deep neural network model. Besides, a lot of high quality images are needed. After
38 training the deep neural network model, an input image can be correctly recognized.
39 However, due to the large number of image categories and the limitation of computing
40 resources, it is difficult for traditional classification algorithms to achieve satisfactory
41 accuracy. In fact, deep learning is an extension of machine learning and many
42 scholars have done a lot of research on it. Different from the traditional image
43 classification method, it does not require complex feature decomposition of the target
44 image. Deep learning uses deep neural network models and a large number of images
45 to learn features. Thus, the deep learning algorithm is suitable for image classification
46 task. The unexplained black box in AI system makes users can only get decision
47 results, and cannot distinguish the specific reasons, processes and action logic of the
48 system's decision-making (Florian et al., 2015 and Leung et al., 2013). This kind of

49 AI system is difficult to achieve the trust and understanding of different
50 decision-makers, because we fail to know what is controlling the design, operation
51 and decision-making, especially for autonomous decision-making with uncontrollable
52 risks. In this regard, the Explainable AI (XAI) is widely developed in handling these
53 issues (Parkhi et al., 2015). Generally, different decision-makers possess different
54 goals in XAI to ensure the high performance of intelligent decision-making and
55 provide a reasonable explainable model synchronously, so that different stakeholders
56 could effectively understand, trust and manage AI systems (Jamie et al., 2011). In this
57 process, multi-attribute group decision-making has been widely used, which can
58 effectively evaluate different decision schemes through extracting and fusing the
59 cognitive judgment information of multiple decision-makers (Alejandro 2016; Alex et
60 al., 2017), ensuring that multiple decision-makers can achieve group consensus and
61 trust for the Explainable AI.

62 In daily life, when we look at a image, the most basic task is what the image is,
63 whether it is a landscape image or a figure image, whether it describes a building or
64 food. For computer vision, this is an image classification task. The main difficulty of
65 image classification is the step of feature extraction. Once a distinguishable feature is
66 found, the image classification becomes very easy. The so-called feature extraction
67 refers to constructing an algorithm to extract features in the target image, such as the
68 edge feature of the face, the color feature of the skin etc. The extracted features are
69 used to distinguish the target object as much as possible from other objects. For
70 example, what you need to distinguish is a black cat from a white cat, so the color
71 feature is definitely a good feature. However, the difficulties encountered in life are
72 often difficult to extract features, such as detecting pedestrians and vehicles on noisy
73 streets. The task requires high accuracy of the detection algorithm to avoid a car

74 accident.

75 At present, the best approach is to use deep neural network model for image
76 classification tasks (Cires et al., 2012). The experimental result of VGGnet shows that
77 the block with the same shape can obtain better classification accuracy by
78 constructing deeper convolutional neural network model. Inspired by this idea, the
79 deep residual network is constructed by cross-layer connection method and it achieves
80 higher classification accuracy. GoogLeNet increases the adaptability of the network to
81 different scales, showing that adjusting the width of model is also an important
82 method to obtain better classification accuracy. ResNeXt and Xception added
83 cardinality to the network model, proving that the cardinality can not only reduce the
84 overall parameters of the model, but also has a strong ability of representation.

85 The attention mechanism is similar to that of human vision. When we look at the
86 scene around us, we always focus our attention on the main things to get key
87 information (Pentina et al., 2015). The main purpose of the attention mechanism is to
88 make the system focus its attention on key information in the scene. Attention
89 mechanics can be used in a wide range of scenarios. The neural network captures key
90 information with the help of the attention mechanism and we can take use of attention
91 mechanism to observe things in the environment better. In the traditional neural
92 network model, adding convolution channels and operating multiple convolutions on
93 features in the same channel usually bring a certain degree of accuracy improvement.
94 The attention mechanism make neural network model know how to pay attention on
95 channel dimension and spatial dimension. In order to verify the role of attention
96 mechanism in computer vision more clearly, the CAM and the SAM are analyzed
97 from the point of attention domain. The attention domain mainly consists of three
98 types: spatial domain, channel domain and mixed domain. In our experiments, it is

99 found that better performance is obtained with using the CAM before the SAM.

100 Daniel Wolpert believed that the reason for the evolution of the brain is not for
101 thinking and feeling, but for controlling motion, which is the core idea of the Deep
102 Reinforcement Learning (Lin et al., 2017 and Holzinger et al., 2017). The problem
103 model studied by it includes an environment and an agent interacting with the
104 environment. The goal of reinforcement learning is to design a behavioral strategy for
105 the agent that maximizes its benefits in interacting with the environment. Google
106 DeepMind used this strategy in 2016 to get computers to go beyond the level of top
107 professionals. However, the development of object detection and behavior recognition
108 in complex environment is slow, that is mainly because :(1) the efficiency of model
109 optimization algorithm is not high; (2) Lack of interpretability of constructed model;
110 (3) Small sample with label; (4) Complex and unrestricted scenes; (5)
111 High-dimensional video data leads to difficulty in parameter adjustment. Curriculum
112 Learning and self-paced Learning represent the recently proposed Learning strategy.
113 Their core idea is to simulate the cognitive mechanism of human beings, they first
114 learn simple and general knowledge structure, then gradually increase the difficulty
115 degree and transition to more complex and professional knowledge. These two
116 methods have similar conceptual learning paradigms, but differ in their specific
117 learning schemes.

118 This paper proposes a novel image classification framework based on
119 multi-perspective deep transfer learning. Attention mechanism is used to describe
120 varied image characteristics, the self-paced learning strategy is adopted to solve the
121 problems of small number of labeled samples and model interpretation. A nonlinear
122 model based on deep transfer learning is proposed to solve the problem that it is
123 difficult to distinguish the perspective-related features from image-related features,

124 the interpretability of the model is further improved.

125 The main structure of this paper is as follows: Section 2 starts with a brief
126 presentation of necessary related concepts. Section 3 proposes a channel and spatial
127 attention module based on explainable machine learning. Designed experiments and
128 discussion are depicted in Section 4. Finally, some conclusions are presented in
129 Section 5.

130 **2 Related work**

131 Explainable AI (XAI) refers to those Artificial Intelligence techniques aimed at
132 explaining, to a given audience, the details or reasons by which a model produces its
133 output (Arrieta et al., 2020). To this end, XAI borrows concepts from philosophy,
134 cognitive sciences and social psychology to yield a spectrum of methodological
135 approaches that can provide explainable decisions for users without a strong
136 background on Artificial Intelligence. Therefore, XAI targets at bridging the gap
137 between the complexity of the model to be explained, and the cognitive skills of the
138 audience for which explainability is sought. Interdisciplinary XAI methods have so
139 far embraced assorted elements from multiple disciplines, including signal processing,
140 adversarial learning, visual analytics or cognitive modeling, to mention a few.
141 Although reported XAI advances have risen sharply in recent times (Russakovsky et
142 al., 2020 and Deng et al., 2015), there is global consensus around the need for further
143 studies around the explainability of ML models. A major focus has been placed on
144 XAI developments that involve the human in the loop and thereby, become
145 human-centric. This includes interpretable reasoning of models, neurosymbolic
146 reasoning or systems based on fuzzy rules, etc. At present, explanations provided by
147 different algorithms are fragmented and independent, which makes it difficult to
148 determine reasonable decisions and explain model structures. In addition, in the

149 design of interpretable classifier, the selection of optimal training set, correlation
150 selection of heat graph, semantic analysis, model visual interpretation and error
151 analysis cannot be combined compulsively (Wang et al., 2017 and Yang et al., 2020).
152 Moreover, the text interpretation can't match the characteristics of a certain layer of
153 deep learning network, and it lacks of continuous interpretability. In general, the text
154 explanation generated for classification comes from training data based on model
155 annotations. Up to now, data labels are set manually and are very subjective, and it
156 doesn't take into account the differences between the different elements. Therefore, it
157 is not possible to determine the relevant region of the image that is most useful for
158 classification. In most cases, experts are encouraged to use their attribute label data as
159 interpretable evidence. Existing interpretable artificial intelligence models can provide
160 the basis behind the classification. However, in the existing classification model, there
161 is no mechanism for identifying potential misclassification of classifiers. Warning
162 users about misclassification will help prevent errors from entering the system. One of
163 the reasons for misclassification is the reduction of distance between classes. Some
164 outliers or edge elements of a class can share the common characteristics of adjacent
165 classes. However, there is no mechanism to ensure the number of subclasses of a
166 given class and whether it makes sense to merge closely related subclasses of two
167 adjacent classes into a new class and implement the correct classification. Xu (Xu et
168 al., 2017) proposes a solution based on database transaction model interpretation,
169 whose explanation is on the basis of logical structure or reasoning. The static structure
170 makes it unsuitable for the deep network classifier. In the constructed model, the
171 system dynamically give appropriate explanations from stored vocabularies, which in
172 turn are generated based on model learning. It provides a consistent view of models
173 and interpretations beyond the scope of existing technology. Reference (Sachan et al.,

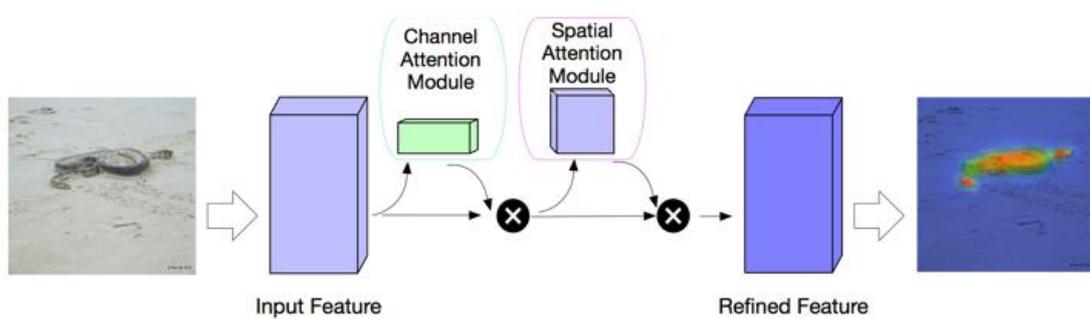
174 2020) proposes an interpretable model in which they are interactively validated
175 through visual features and similarity. Moreover, k-means clustering was used to
176 analyze the similar features, so that the average features obtained had greater
177 robustness and relatively low time complexity. In addition, it does not consider the
178 importance or relevance of model, nor does it cluster with respect to output classes.
179 He et al. (He et al., 2016) conducts the explanatory demonstration of the model
180 mainly from the aspect of correlation calculation. By observing the change of the
181 connection mode of the network, the hidden layer is explained visually. However, the
182 feature learned by the network layer are not described in detail. Simonyan et al.
183 (Simonyan et al., 2020) proposed an image interpretation and generation method
184 based on visual features. Take an image signature with a fixed length of 8000 to
185 generate a caption. In this model, the correlation of features is firstly determined and
186 the signature generation is carried out on this basis. Since eigenvalues can be of any
187 length, strategies that follow highly correlated features are interpretable. As the power
188 of interpretation becomes more important in intelligent decision-making, AI systems
189 are no longer there to serve as black boxes (Xie et al., 2016) . Decision makers of AI
190 services have the right to know the reasons behind their decisions so that they can
191 better play to their strengths.

192 Abhronil et al. (Abhronil et al., 2019) propose a novel model based on the
193 attention mechanism and it is named the residual attention network. As the network
194 layer deepens, attention modules can extract key information from different layers.
195 Finally, it got 4.8% Top-5 error rate on ImageNet. HU J et al. (Hu et al., 2019)
196 proposed the SENet 44 network model. In the training process, the model can
197 distinguish the importance of different channels, then enhances useful features and
198 inhibits useless features according to the importance of feature. Finally, it won the

199 ILSVRC2017 classification task championship with a 2.25% Top-5 error rate. This
200 paper mainly verifies the interpretability of the model through image classification
201 scenarios in real applications.

202 **3 The proposed channel and spatial attention model**

203 This paper makes an deep study of image classification and recognition from the
204 most cutting-edge technical perspectives such as deep reinforcement learning,
205 self-paced learning and transfer learning. In view of the existing problems in the field
206 of image classification in complex environments, this paper considers the solutions
207 based on deep network, and further studies popular algorithms such as deep
208 reinforcement learning, self-paced learning and transfer learning, aiming at the urgent
209 model explanation problems in the framework of deep network. On this basis, object
210 detection technology and image recognition technology in complex environments are
211 specifically studied. The overall technical framework of the paper and the relationship
212 between the research contents of each part are shown in Figure 1, and the research of
213 this paper is also carried out according to these contents one by one.



214

215 Figure 1. The overview of model

216 For the convolutional neural network model, depth, width and attention
217 mechanism are the main factors affecting the accuracy of image classification. At
218 present, attention mechanism contains CAM and SAM. CAM acts on the channel

219 domain, weighting different channel features. For a $C \times H \times W$ feature graph, the
 220 C weight of channel attention is different, while the weight of $H \times W$ is the same.
 221 For CAM, the weight of each C on different channel dimensions needs to be learned.
 222 To reduce the amount of computation and improve classification accuracy, the pooling
 223 layer in the general convolutional neural network directly uses the maximum pooling
 224 method or average pooling method to compress the image information. For SAM,
 225 only the key information in the spatial features is extracted. We first introduce the
 226 general framework of proposed model with CAM and SAM in this section. Finally,
 227 we describe how to pull them together.

228 After convolution operations, an intermediate feature map $F_{in} \in \mathbb{R}^{C \times H \times W}$ is
 229 obtained. F_{in} is the input of model. One-dimensional channel map $M_c \in \mathbb{R}^{C \times 1 \times 1}$ is
 230 inferred after the channel attention module and two-dimensional spatial map $M_s \in$
 231 $\mathbb{R}^{1 \times (H \times W) \times (H \times W)}$ is inferred after the spatial attention module as shown in Fig. 1.
 232 The entire attention calculation process can be summarized as:

$$233 \quad F_1 = M_c(F_{in}) \otimes F_{in} \quad (1)$$

$$234 \quad F_{out} = M_s(F_1) \otimes F_1 \quad (2)$$

235 where \otimes denotes element-wise multiplication. After F_{in} passing the channel
 236 attention module, F_1 is achieved. F_{out} is the final output and the attention value is
 237 broadcasted during multiplication process. Attention and feature are the multiplication
 238 of element levels, which will be propagated automatically, that is, channel attention
 239 broadcasts along spatial dimension, while spatial attention broadcasts along channel
 240 dimension

241

242 3.1 Channel attention module

243 The CAM captures the relationship between channel features. It pays attention to
244 the key channel information and weakens the influence of the useless channel
245 information (Russakovsky et al., 2010 and Deng et al., 2009). It uses an attention
246 mechanism similar to the self attention mechanism (query, key, value) to get the
247 similarity between channel graphs, and then use the weight of channel graphs to
248 update. Finally, the matrix of computation attention is obtained, which can enhance
249 the key features. The CAM makes the neural network model to pay attention to the
250 channel features with the key information. On the basis of convolution, we first
251 extrude the feature graph to obtain the global feature of each channel. Then, we use
252 the global feature to get the relationship between different channels and the get the
253 weight of different channels. Finally, we multiply the weights to get the features on
254 the basis of the original feature graph.

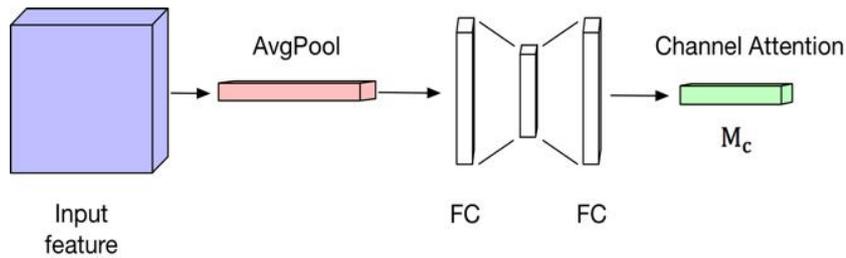
255 In a convolutional neural network, the convolution operation only performs on
256 image feature space, it is difficult for the convolution module to get the relationship
257 between different feature channels. To get an eigenmatrix, an image needs go through
258 several convolutional layers and the number of channels represents the number of
259 cores in the convolutional layers. In a normal neural network, the number of
260 convolution kernels is usually as high as 1024 or 2048. Therefore, not every channel
261 is useful for feature extraction. The CAM will help the neural network model to select
262 more informative channels. Besides, We encode spatial features using a global
263 average pool that provides feedback for each pixel on the feature map. The following
264 formula shows the global average pool calculation process.

$$265 F_{avg} = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W F_{in}(i, j) \quad (3)$$

266 where $F_{avg} \in \mathbb{R}^{C \times 1 \times 1}$ represents the result after implementing the global average
 267 pooling on the input feature map F_{in} . In order to capture the relationships between
 268 different channels when obtaining the global description features, two conditions need
 269 to be met for CAM: firstly, it must be flexible, because it needs to learn the nonlinear
 270 relationship between different channels; Secondly, the learning relationship is not
 271 mutually exclusive, because it allows for a multichannel feature instead of a hot spot
 272 form. We describe the channel attention map M_c as follows.

$$273 \quad M_c = \sigma \left(W_2 ReLU(W_1 F_{avg}) \right) \quad (4)$$

274 where $W_1 \in \mathbb{R}^{\frac{C}{r} \times C}$ and $W_2 \in \mathbb{R}^{C \times \frac{C}{r}}$, σ denotes Sigmoid function. W_1 and
 275 W_2 are fully connected layers. To improve the explanatory ability of the model, we
 276 construct two fully connected layers, namely the bottleneck structure. W_1 layer is a
 277 way of dimensionality reduction and the dimensionality reduction factor r is a super
 278 parameter. Then the Relu function is used and the W_2 layer restores the dimension to
 279 the original finally. The process is shown in Figure 2.



280

281 Figure 2. The details of channel attention module

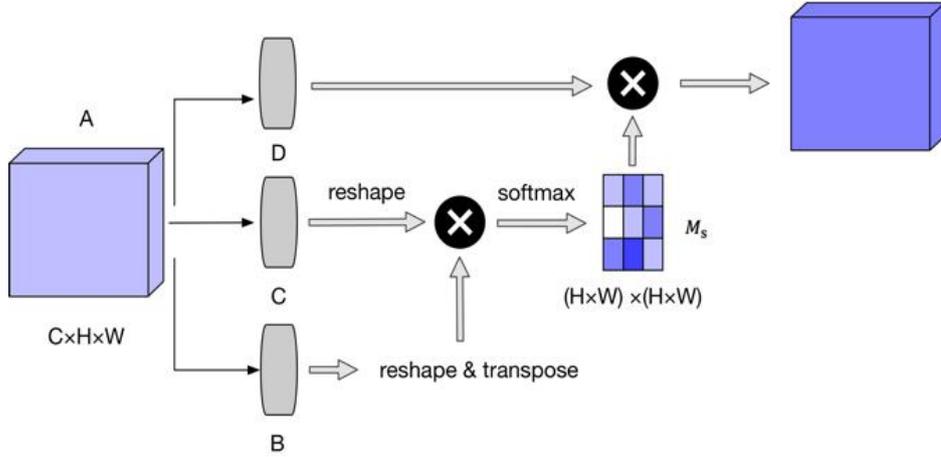
282

283 3.2 Spatial attention module

284 Unlike CAM, SAM only plays the role of distinguishing key information
 285 within a single image feature map. First of all, we use average pooling and max
 286 pooling to compress the input feature and then we use mean and max operations on

287 the input feature at channel dimensions. Finally, we get two two-dimensional
 288 features. Considering different channel size, the two two-dimensional features are
 289 spliced together to obtain the feature with channel number of 2. And they are
 290 convolved to ensure that the resulting features are consistent with the input feature in
 291 spatial dimension.

292 Spatial attention map is generated by paying attention to the internal relationship,
 293 as shown in Figure 3. $A \in \mathbb{R}^{C \times H \times W}$ represents the input of the SAM. After conducting
 294 convolutional layers, A generates feature maps B, C and D where $\{B, C, D\} \in \mathbb{R}^{C \times H \times W}$.
 295 We reshape B and C to $\mathbb{R}^{C \times (H \times W)}$ and $H \times W$ represents pixels in spatial module.



296

297 Figure 3. The details of spatial attention module

298 Then we get $\mathbb{R}^{(H \times W) \times (H \times W)}$ by operating a matrix multiplication between the
 299 transpose of B and C. We get the spatial attention map Ms when a softmax layer is
 300 applied. Ms is computed as follows.

$$301 \quad M_{s_{ij}} = \frac{\exp(B_i \times C_j)}{\sum_{i=1}^{H \times W} \exp(B_i \times C_j)} \quad (5)$$

302 where the stronger the correlation, the more similar the features of the two sites.

303 Figure 4 shows the visual result of images. It is observed that our proposed model
 304 focuses limited attention on key information, saves resources, and quickly obtains the
 305 most effective information.

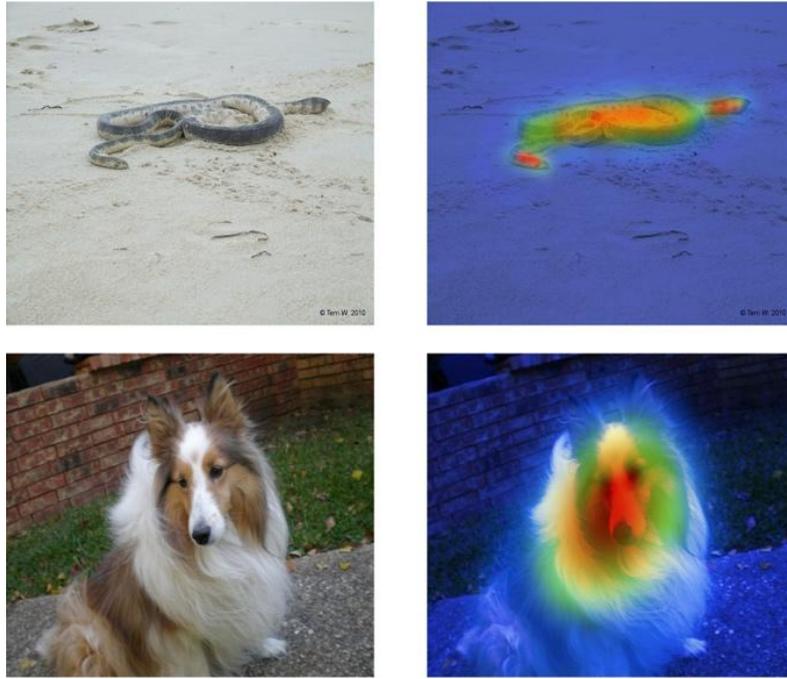
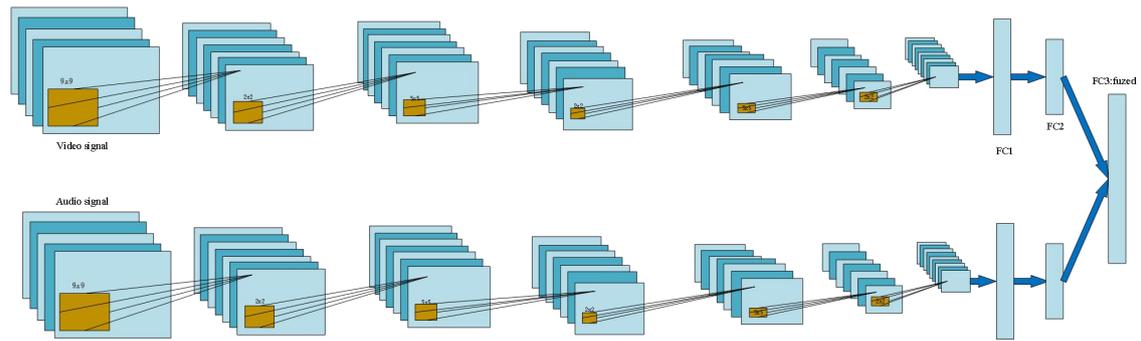


Figure 4. Diagram of the visualization results

3.3 Explainable object detection model of combing self-paced learning and deep reinforcement learning

The deep network model is designed for learning across model characteristics. This paper proposes a new application in deep network, using deep network to learn multi-mode. In particular, this paper demonstrates that cross-modal feature learning - if multimodal features are present during feature learning, better features can be learned for a modal (multimodal learning, monomodal-tested). In addition, the paper designed how to learn a shared feature between multiple modes and evaluate it on a particular task -- the classifier was trained with only audio data but tested on video-only data. Multi-mode explainable deep network model can be seen in Figure 5. This model consists of two streams, one for video information and the other for audio information. The structure of the two streams is identical, each consisting of eight layers (including the input layer).



321
322

Fig. 5. Proposed multi-mode explainable deep network model

323

324 There are two problems with the traditional multimodal model. First, there is no
 325 clear goal for the model to find correlations across modes. Some hidden layer units
 326 adjust the parameters only for voice, and others adjust the parameters only for video,
 327 so that it is possible for the model to find the desired feature. Second, there is only
 328 one mode for supervised training and testing in the cross-modal learning arrangement,
 329 which makes the model unexplainable. If there is only one modal representation, it is
 330 necessary to integrate observable variables that are not observed. Therefore, this paper
 331 proposes a deep self-coding model to solve the above problems. Inspired by the
 332 noise-reducing self-coding model, this paper proposes a training two-mode deep
 333 self-coding model (Figure 5), which uses an extended (extended single-mode input)
 334 but noisy data set. In fact, the model is still required to reconstruct the two modes
 335 when one mode uses zero as input and the other uses the original value as input when
 336 expanding. Therefore, one-third of the training data is input only by video, one-third of
 337 the training data is input only by voice, and the last third has both video and voice.
 338 This model can be viewed as an example of multitasking learning.

339 When designing the intensification strategy, the paper uses the Q network to
 340 interact with its environment during the data generation phase. The system looks at
 341 the current scene, which consists of audio and video frames, and takes actions using
 342 the -greedy strategy. This environment in turn provides scalar rewards. Interaction

343 experiences are stored in replay memory M . Replaying M preserves N recent
344 experiences, which are then used to update the network parameters during the training
345 phase. In the training stage, the network structure will use the data stored in replay
346 memory M to train the network. Assume that the superparameter n represents the
347 number of experiences replay, and for each experience replay, a mini-cache B
348 containing several interactions is randomly sampled from the finite size replay
349 memory M . The model will be trained by sampling from cache B , and the parameters
350 of the network will be updated iteratively in the direction of The Behrman target. The
351 algorithm is divided into two phases to avoid latency. Therefore, this paper divides the
352 algorithm into two stages: in the first stage, the robot collects data through limited
353 time interaction with human beings; In the second stage, it enters the stage. During
354 this rest phase, the training phase is activated to train the multimodal depth Q
355 network.

356 In this paper, for the sake of the regularization model and make it thin, it is
357 needed to make each unit has a hidden layer using the regularized the expected
358 activation function of punishment, the form of the regularized punishment is the need
359 to focus on research, it determines the cell activation function of hidden layers on the
360 sparse sex (whether the function of the activation of the hidden layer unit is activated
361 or not).

362 In order to avoid non-convex optimization problems from falling into poor local
363 solutions, the proposed network optimization method adopts multiple random
364 initializations to train the model, and then chooses the initialization network with the
365 best effect to construct the model. However, this method is too adhoc and the
366 calculation cost is too high. Self-learning is just the best solution to non-convex
367 optimization problems. The curriculum learning is to simulate the cognitive

368 mechanism of human beings by first learning simple and universal knowledge
369 structure and then gradually increasing the difficulty to learn more complex and
370 specialized knowledge. However, self-learning has been improved in course learning.
371 Instead of assigning prior knowledge to sample learning sequence in advance, the
372 learning algorithm itself determines the next learning sample in each iteration.

373 **4 Experimental results and analysis**

374 The dataset is used in the experiment including ImageNet- 1K and Cifar-10050
375 (Hu et al., 2018; Zhang et al., 2010; Sanghyun et al., 2010; Chen et al., 2017; Zhao et
376 al., 2018). The Cifar-100 contains 100 classes and every class has 600 color images.
377 But every images is only a size of 32×32 . Five hundred images in each class serve
378 as the training set and the rest as test set. For each image, it has two labels, fine-labels
379 and coarse-labels, which represent the fine-grained and coarse-grained labels of the
380 image respectively and Cifar-100 is hierarchical. In Figure 6, we extracted the images
381 from Cifar-100 as the visual example.



382

383

Figure 6. Visual examples of the Cifar-100

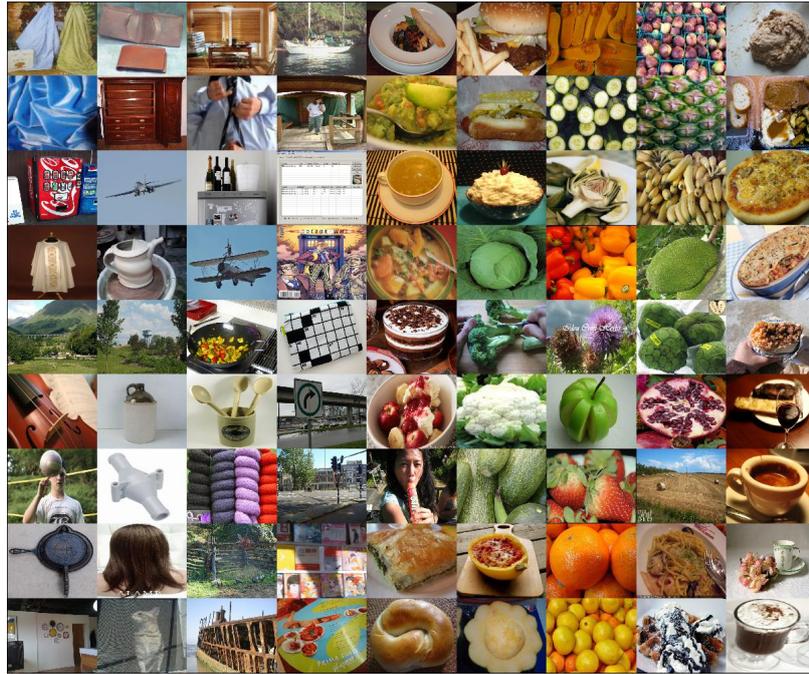


Figure 7. Visual examples of the ImageNet-1K

ImageNet-1K is an image dataset and each concept image is quality-controlled and manually tagged. At present, ImageNet-1K consists of 1,4197,122 images. The major categories include: animal, bird, fish, flower etc. In Figure 7 shows the visual examples of the ImageNet-1K.

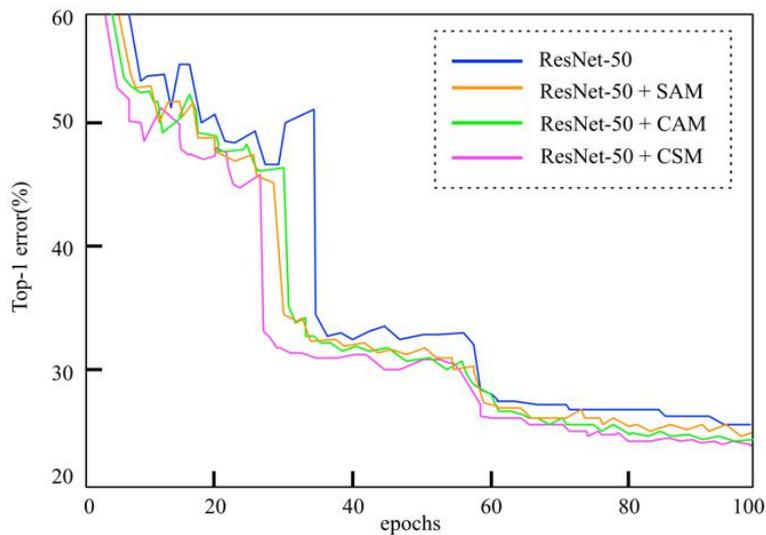


Figure 8. Comparison of different network models

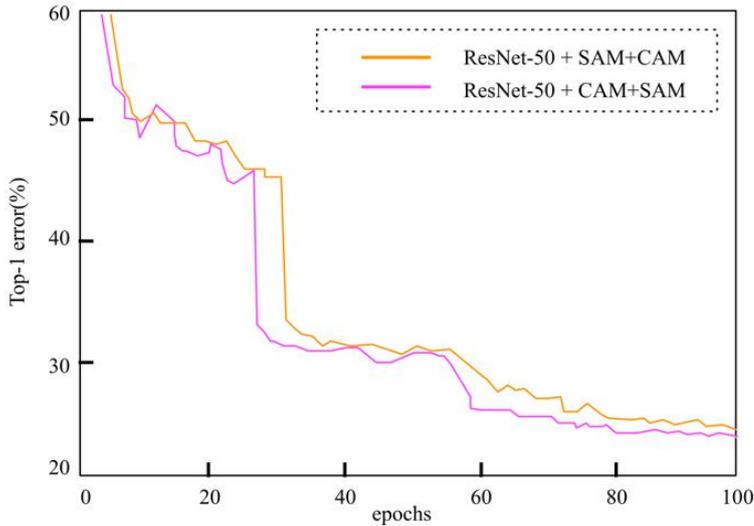


Figure 9. Comparison of different network models

In experiment, we compare the effectiveness between the CAM and the SAM. We compared 4 different network models: baseline network, baseline network with CAM, baseline network with SAM, baseline network with CSM. The result of experiment is shown in Figure 8. As is shown in Figure 8, it is easily concluded that the ResNet-50 model with CSM has achieved higher accuracy. We can observe that CAM perform better than SAM. Besides, combination of two attention modules can bring better performance. The experiment shows that it is effective to conduct CAM and SAM at the same time. The experimental results are shown in Table 1. Based on the above part of the experiment, we find that it is effective to conduct CAM and SAM at the same time for improving the expression ability of neural networks. We want to know if the CAM should be used before SAM. So, in the ablation experiment, we compare the effect of CAM and SAM in different order of use. The baseline network with CAM and SAM, baseline network with SAM and CAM are conducted respectively. Figure 9 shows the result of experiment.

Consistent with the above experiment, adding the attention module still bring improvement on image classification accuracy from Figure 8. We can observe that the

410 CAM-first combination method achieves better performance. The experimental result
 411 shows that CAM-first combination method is more effective.

412

TABLE I. COMPARISON OF DIFFERENT NETWORK MODELS

Architecture	Top-1 error (%)
ResNet-50	24.55
ResNet-50+SAM	23.47
ResNet-50+CAM	23.21
ResNet-50+CSM	22.78

413

4.1 Image Classification on ImageNet-1K

414

In this part, we use ResNet and WideResNet as the baseline model. On this basis,

415

we add attention mechanism for comparison. The extensive image classification

416

experiments are conducted based on the ImageNet-1K. The structure of ResNet with

417

adding SE module is shown in Figure 9 and the results of the experiment is shown in

418

Table 2. The experiment still prove that networks with CSM performs better than the

419

baseline module, indicating that attention mechanism can be well used on the various

420

network models. Besides, the depth and width of the neural network also greatly

421

affect image classification accuracy. SENet won the ILSVRC2017 classification task

422

championship. But CSM fuses channel features with spatial features for better

423

representation capabilities and CSM performs better than SENet.

424

TABLE II. RESULTS OF THE EXPERIMENT ON THE IMAGENET-1K DATASET

Architecture	Top-1 error (%)	Top-5 error (%)
ResNet-18	29.62	10.56
ResNet-18+SE	29.43	10.24
ResNet-18+CSM	29.29	10.12
ResNet-34	26.71	8.62
ResNet-34+SE	26.16	8.37
ResNet-34+CSM	26.03	8.28
ResNet-50	24.55	7.52
ResNet-50+SE	23.21	6.74
ResNet-50+CSM	22.78	6.57

WideResNet18(widen=1.5)	26.86	8.91
WideResNet18(widen=1.5) +SE	26.23	8.51
WideResNet18(widen=1.5) +CSM	26.14	8.49

425 4.2 Image Classification on CIFAR-100

426 Based on Cifar-100, we conduct image classification experiment to verify the
427 effectiveness of the CSM. ResNet and WideResNet are used as baseline model. Table
428 3 shows the experimental result. The experimental results prove that the combination
429 of CAM and SAM can improve classification accuracy. Besides, the depth and width
430 of the neural network also greatly affect image classification accuracy.

431 TABLE III. RESULTS OF THE EXPERIMENT ON THE CIFIR-100 DATASET

Architecture	Accuracy (%)
ResNet-18	91.7
ResNet-18+CSM	93.1
ResNet-34	92.4
ResNet-34+CSM	93.8
ResNet-50	92.9
ResNet-50+CSM	94.3
WideResNet18(widen=1.5)	92.8
WideResNet18(widen=1.5) +CSM	94.2

432 As can be seen from Table 3., the prediction density map we get after adding
433 feature self-learning is closer to the true value, so we can see from Fig.8 and Fig.9
434 below that we can more accurately describe the feature after adding feature
435 self-learning Loss reduction and corresponding changes in MAE accuracy, we also
436 give the performance of net similarity.

437 5 Conclusions

438 Different from the previous research on image classification, we propose an
439 attention module based on spatial dimension and channel dimension. This module
440 derives the attention map by CAM and SAM respectively. Then it multiplies the
441 attention map into the input feature map. The experiment shows that adding the

442 attention module to some image classification algorithms is an effective method. In
443 order to make image classification perform better, combination of CAM and SAM can
444 improve classification accuracy and the CAM should be used before SAM. The CAM
445 selectively enhances some feature channels and suppresses some feature channels by
446 learning the relational mapping. The SAM aggregates features by weighting features
447 at spatial dimension. We conducted a lot of image classification experiments for
448 comparison based on the ImageNet-1K and Cifar-100. In fact, the attention module
449 can be well embedded in different deep neural networks and it improves the
450 explainable deep neural network's ability of expression. Besides, the width and depth
451 of the neural networks are also worth considering.

452 In the following research, we will pay more attention to the improvement of
453 real-time algorithm and effectiveness of explainable deep learning, also be able to
454 obtain more accurate data collection develop an intelligent prediction machines for
455 image classification based on more effective machine learning approaches.

456

457 **Acknowledgments**

458 This work is supported by the National Natural Science Foundation of China (NO.
459 61702226); the 111 Project (B12018); open Fund of Jiangsu Key Laboratory of Image
460 and Video Understanding for Social Safety, Nanjing University of Science and
461 Technology, Nanjing (J2021-7).

462 **Conflict of interest:** The authors declare that they have no conflict of interest.

463 **Data availability:** The data that support the findings of this paper are available from
464 the corresponding author.

465 **Authorship contributions:** **Bin Wu:** Investigation, Methodology, Writing- original

466 draft, Supervision.

467 **Yuhong Fan:** Writing- Reviewing and Editing, Methodology, Visualization.

468 **Li Mao:** Validation, Visualization, Data curation.

469

470 **References**

471 Diba A, Sharma V, Pazandeh A et al. Weakly super-vised cascaded convolutional
472 networks. In: IEEE conference on computer vision and pattern recognition
473 (CVPR), July 2017, pp. 5131-5139. New York, NY: IEEE.

474 Ren, Shaoqing, Kaiming He, Ross Girshick and Jian Sun. Faster R-CNN: Towards
475 Real-Time Object Detection with Region Proposal Networks. IEEE Trans Pattern
476 Analysis And Machine Intelligence 2017; 39: 1137 - 49.

477 Girshick, Ross, Jeff Donahue, Trevor Darrell and Jitendra Malik. Rich Feature
478 Hierarchies for Accurate Object Detection and Semantic Segmentation. In: IEEE
479 conference on computer vision and pattern recognition (CVPR), 2014, pp. 580 -
480 87.

481 Xing Chen, Ming Li, Hao Zhong*, Yun Ma, Ching-Hsien Hsu. DNNOff: Offloading
482 DNN-based Intelligent IoT Applications in Mobile Edge Computing. IEEE
483 Transactions on Industrial Informatics, Publish Online, DOI:
484 10.1109/TII.2021.3075464.

485 Xing Chen, Shihong Chen, Yun Ma, Bichun Liu, Ying Zhang, Gang Huang*. An
486 Adaptive Offloading Framework for Android Applications in Mobile Edge
487 Computing. SCIENCE CHINA Information Sciences, 2019, 62(8): 82102.

488 Ying Zhang, Gang Huang, Xuanzhe Liu, Wei Zhang, Hong Mei, Shunxiang Yang.
489 Refactoring android Java code for on-demand computation offloading. ACM
490 SIGPLAN Conference on Object-Oriented Programming, Systems, Languages,
491 and Applications. 2012.

492 Xing Chen, Junqin Hu, Zheyi Chen*, Bing Lin*, Naixue Xiong, Geyong Min. A
493 Reinforcement Learning Empowered Feedback Control System for Industrial
494 Internet of Things. IEEE Transactions on Industrial Informatics, Publish Online,
495 DOI: 10.1109/TII.2021.3076393.

496 Bing Lin, Yin hao Huang, Jianshan Zhang, Junqin Hu, Xing Chen*, Jun Li*.

497 Cost-Driven Offloading for DNN-based Applications over Cloud, Edge and End
498 Devices. *IEEE Transactions on Industrial Informatics*, 2020, 16(8): 5456-5466.

499 Bing Lin, Fangning Zhu, Jianshan Zhang, Jiaqing Chen, Xing Chen*, Naixue Xiong,
500 Jaime Lloret Mauri. A Time-driven Data Placement Strategy for a Scientific
501 Workflow Combining Edge Computing and Cloud Computing. *IEEE*
502 *Transactions on Industrial Informatics*, 2019, 15(7): 4254-4265.

503 Xing Chen, Fangning Zhu, Zheyi Chen*, Geyong Min*, Xianghan Zheng, Chunming
504 Rong. Resource Allocation for Cloud-based Software Services Using
505 Prediction-Enabled Feedback Control with Reinforcement Learning. *IEEE*
506 *Transactions on Cloud Computing*, Publish Online, DOI:
507 10.1109/TCC.2020.2992537.

508 Gang Huang, Chaoran Luo, Kaidong Wu, Yun Ma and Ying Zhang, Xuanzhe Liu.
509 Software-Defined Infrastructure for Decentralized Data Lifecycle Governance:
510 Principled Design and Open Challenges. *IEEE International Conference on*
511 *Distributed Computing Systems*, 2019.

512 Schroff, Florian, Dmitry Kalenichenko, and James Philbin. FaceNet: A Unified
513 Embedding for Face Recognition and Clustering. In: *IEEE conference on*
514 *computer vision and pattern recognition (CVPR)*, 2015, pp. 815 - 23.

515 Leung Y., Ji N.N., Ma J.H., An integrated information fusion approach based on the
516 theory of evidence and group decision-making, *Information Fusion*,
517 (2013),14(4):410-422.

518 Parkhi, Omkar M, Andrea Vedaldi, and Andrew Zisserman. Deep Face Recognition.
519 In: *British machine vision conference (BMVC)*, 2015.

520 Shotton, Jamie, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard
521 Moore, Alex Kipman and Andrew Blake. Real-Time Human Pose Recognition in
522 Parts from Single Depth Images. In: *IEEE conference on computer vision and*
523 *pattern recognition (CVPR)*, 2011, pp. 1297 - 1304.

524 Newell, Alejandro, Kaiyu Yang, and Jia Deng. Stacked Hourglass Networks for
525 Human Pose Estimation. In: *Computer Vision-ECCV 2016*, pp. 483 - 99.

526 Krizhevsky, Alex, Ilya Sutskever, Geoffrey E and Hinton. ImageNet Classification
527 with Deep Convolutional Neural Networks. *Communications of The ACM* 2017;
528 60: 84 - 90.

529 Cires , an, Dan, Ueli Meier and Juergen Schmidhuber. Multi Column Deep Neural

530 Networks for Image Classification. In: IEEE conference on computer vision and
531 pattern recognition (CVPR), 2012, pp. 3642 – 49.

532 Cireş, an, Dan, Ueli Meier and Juergen Schmidhuber. Multi Column Deep Neural
533 Networks for Image Classification. In: IEEE conference on computer vision and
534 pattern recognition (CVPR), 2012. pp. 3642 – 49.

535 A.Pentina, V.Sharmanska, and C.H.Lampert. Curriculum learning of multiple tasks. In
536 Proceedings of the 28th IEEE Conference on Computer Vision and Pattern
537 Recognition, pages. 2547 – 2554, (2015).

538 L.Lin, K.Wang, D.Meng, et al. Active Self-Paced Learning for Cost-Effective and
539 Progressive Face Identification[J]. IEEE Transactions on Pattern Analysis and
540 Machine Intelligence, PP(99):7-19, (2017).

541 A.Holzinger, C.Biemann, S.P.Constantinos, and B.K.Douglas, What do we need to
542 build explainable ai systems for the medical domain?, arXiv:1411.1784,(2017).

543 A. Barredo Arrieta, N. Diaz-Rodriguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado,
544 S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, F. Herrera,
545 “Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities
546 and challenges toward responsible AI” , Information Fusion, vol. 58, pp. 82-115,
547 June (2020).

548 D. Xu, W. l. O. Yang, Xavier. Alameda-Pineda, E. Ricci, X. g. Wang, and N.cu. Sebe.
549 Learning deep structured multi-scale features using attention-gated crfs for
550 contour prediction. In NIPS, (2017), pages. 3961 – 3970.

551 Sachan, S., Yang, J. B., Xu, D. L., Benavides, D. E., Li, Y. An explainable AI
552 decision-support-system to automate loan underwriting. Expert Systems with
553 Applications, (2020), 144, 113100.

554 He K, Zhang X, Ren S and Sun J. Deep residual learning for image recognition. In:
555 IEEE conference on computer vision and pattern recognition (CVPR), 2016.

556 Simonyan K and Zisserman A. Very deep convolutional networks for large-scale
557 image recognition. arXiv preprint arXiv preprint arXiv: 1409.1556.

558 Sengupta, Abhronil, Yuting Ye, Robert Wang, Chiao Liu and Kaushik Roy. Going
559 Deeper in Spiking Neural Networks: VGG and Residual Architectures.
560 Neuroscience, 2009; 13: 95 – 95.

561 Xie S,Girshick R, Dollar P, Tu Z and He K. 2016. Aggregated residual

562 transformations for deep neural networks. arXiv preprint arXiv: 1611.05431.

563 Wang, Fei, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang,
564 Xiaogang Wang and Xiaoou Tang. Residual Attention Network for Image
565 Classification. In IEEE conference on computer vision and pattern recognition
566 (CVPR), 2017, pp. 6450 - 58.

567 Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng and Olga Russakovsky. Towards
568 Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in
569 the ImageNet Hierarchy. In: Conference on fairness, accountability and
570 transparency (FAT), 2020.

571 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma,
572 Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander
573 C. Berg and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge.
574 International journal of computer vision (IJCV), 2015.

575 J Deng, A Berg, K Li and L Fei-Fei, What does classifying more than 10,000 image
576 categories tell us? In: Computer Vision-ECCV 2010.

577 O Russakovsky and L Fei-Fei, Attribute Learning in Large scale Datasets. In:
578 Computer Vision-ECCV 2010.

579 J Deng, W Dong, R Socher, L -J Li, K Li and L Fei-Fei, ImageNet: A Large-Scale
580 Hierarchical Image Database. In: IEEE international conference on computer
581 vision and pattern recognition (CVPR) , 2009.

582 HU J, SHEN L and SUN G. Squeeze-and-excitation net works. In: IEEE conference
583 on computer vision and pattern recognition (CVPR), June 2018, pp. 7132-7141.
584 New York, USA. NY: IEEE.

585 Zhang, Yulun, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong and Yun Fu. Image
586 Super-Resolution Using Very Deep Residual Channel Attention Networks. In:
587 Computer Vision ECCV 2010. pp. 294 - 310.

588 Woo, Sanghyun, Jongchan Park, Joon-Young Lee and So Kweon. CBAM:
589 Convolutional Block Attention Module. In: Computer Vision-ECCV 2010. pp. 3
590 - 19.

591 Chen, Long, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and
592 Tat-Seng Chua. SCA-CNN: Spatial and Channel-Wise Attention in
593 Convolutional Networks for Image Captioning. In: IEEE international
594 conference on computer vision and pattern recognition (CVPR), 2017. pp. 6298

595 - 6306.

596 Zhao, Hengshuang, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin,
597 and Jiaya Jia. PSANet: Point-Wise Spatial Attention Network for Scene Parsing.
598 In: Computer Vision-ECCV 2018. pp. 267 - 83.

599



600

601 BIN WU received the B.Sc. degree from Jiangnan University, Wuxi, China, in 1996, and the
602 M.S. degree from Jiangnan University, Wuxi, China, in 2005. He is currently a lecturer with the
603 School of Internet of Things Engineering, Jiangnan University. His major research interests
604 include visual surveillance, object detection, integrated circuit design and application of embedded
605 system.



606

607 Yuhong Fan graduated from Shandong Agricultural University, Tai'an, China in 2008. In 2011,
608 she obtained a master's degree from Xihua University, Chengdu, China. She is currently a lecturer
609 in the Department of Computer Engineering, LangFang YanJing Vocational Technical College,
610 Sanhe, China. She has studied many topics and written several high-quality journal papers and
611 conference papers. His current research interests include data mining, information systems,
612 wireless networks, artificial intelligence, Internet of things and security, medical data analysis,
613 visual monitoring, scene understanding, behavior analysis, target detection and pattern analysis.



614 Li Mao received the B.Sc. degree from Southeast University (Nanjing, China) in
615 1990 and M.S. degree from Donghua University (Shanghai, China) in 2003. He is currently an
616 associate professor at the department of computer science, Jiangnan University. His major
617 research interests include visual surveillance, object detection, and data mining.

618