



Fine-Grained Image Classification via Combining Vision and Language Xiangteng He and Yuxin Peng*, Peking University

Institute of Computer Science & Technology, Peking Universit

1. Introduction



> Motivation

- Not all the parts which obtained through the part detection models are beneficial and indispensable for classification.
- Fine-grained image classification requires more detailed visual descriptions which could not be provided by the part locations or attribute annotations.

> Contribution

- The *vision stream* learns deep representations from the original visual information via deep convolutional neural network.
- The *language stream* utilizes the natural language descriptions which could point out the discriminative parts or characteristics for each image, and provides a flexible and compact way of encoding the *salient visual aspects* for distinguishing subcategories.
- Since the two streams are *complementary*, combining the two streams can further achieves better classification accuracy.

2. Our CVL Approach

Language

(1) This bird has a grey body, a white head with an orange bill and black wings, tarsus

(2)A white bird with black wings and orange (3) This white bird has a bright orange bill

) This white bird has grey wings and tai

2) This bird has a mellow vellow b coloration and deep red feet (3) The bird has a yellow beak with a white head and orange web feet

This colorful bird has an orange crow primaries being rimmed in yellow 2) This bird is grey with red and has a long

3) The bird has a tan spiked crown and short

> Vision Stream

Given an image I, its object region b is generated via saliency extraction and co-segmentation, then the object region is clipped from the original image and saved as image I'. We take the original image I and its object image I' as the inputs of the CNN model to obtain the prediction, which is the result of the vision stream.

> Language Stream

We apply the deep structured joint embedding method, jointly embedding images and fine-grained visual descriptions. We use the inner product of features generated by deep neural encoders ($\theta(v)$ for images and $\varphi(t)$ for texts), and maximize the compatibility between a description and its matching image as well as minimize compatibility with images from other classes.

 $F(v,t) = \theta(v)^{\mathrm{T}} \varphi(t)$



3. Experimental Results

Category	Image	
Sooty Albatross		(1)This bir (2)This bir (3)This bir
California Gull		(1)This bir (2)This bir (3)This bir
Cerulean Warbler		(1)A little (2)The bird (3)This bir

> Comparisons with state-of-the-art methods

We compare with 12 state-of-the-art fine-grained image classification methods on CUB-200-2011 to verify the effectiveness of our CVL approach, which jointly integrates the vision and language streams to exploit the correlation between visual feature and nature language descriptions, and enhance their complementarity.

N loth o d		Train Anno.		nno.	
Method	Bbox	Parts	Bbox	Parts	ACC. (%)
Our CVL Approach					85.55
PD[Zhang et al, CVPR 2016]					84.54
Spatial Transformer[Jaderberg et al, NIPS 2015]					84.10
Bilinear-CNN [Lin et al, ICCV 2015]					84.10
NAC[Simon et al, ICCV 2015]					81.01
TL Atten[Xiao et al, CVPR 2015]					77.90
VGG-BGLm[Zhou et al, CVPR 2016]					75.90
PG Alignment[Krause et al, CVPR 2015]	V		V		82.80
Triplet-A (64)[Cui et al, CVPR 2016]	V		V		80.70
VGG-BGLm[Zhou et al, CVPR 2016]	V		V		80.40
SPDA-CNN[Zhang et al, CVPR 2016]	V	v	V		85.14
Part-based R-CNN[Zhang et al, ICCV 2014]	V	V	V	V	76.37
POOF[Berg et al, CVPR 2013]	√	V	V	V	73.30



IEEE 2017 Conference on **Computer Vision and Pattern** Recognition



Text Rank List(Top3)

rd has wings that are grey and has a black bill. rd is gray in color, with a large curved beak d is white and brown in color, and has a black beak.

rd has large feet, a short yellow bill, and a black and white body. rd has wings that are grey and has a white belly and yellow bill. rd has a yellow beak as well as a white belly.

bird with a short, grey bill, blue crown, nape, white breast. d has a white abdomen, black breast and white throat, blue specks. rd is **blue** and **white** in color with a black beak, and black eye rings.