# Transferring CNNs to Multi-instance Multi-label Classification on Small Datasets

Mingzhi Dong[1], Kunkun Pang[2], **Yang Wu**[3], Jinghao Xue[1], Timothy Hospedales[2], Tsukasa Ogasawara[3]

[1]University College London  [2] University of Edinburgh
[3]**Nara Institute of Science and Technology**

ICIP, 2017

# Outline

# Multi-instance Multi-Label (MIML) Problems

- **Multi-label problems**
  Each instance is associated with *a set of labels*.

- **Multi-instance problems**
  The learner receives *a set of bags, each containing many instances*.
  Instead of receiving the labels of the instances, in multi-instance
  problems, the bags are labelled.

- **Multi-instance Multi-label (MIML) problems**
  The learner receives *a set of bags* which are associated with *a set of
  labels*.

- **Image Tagging as a MIML Problem**
  **Multi-label**: each image is likely worth many words;
  **Multi-instance**: each image can be viewed as a bag of local regions,
  the labels are assigned to a whole image but not its specific regions.
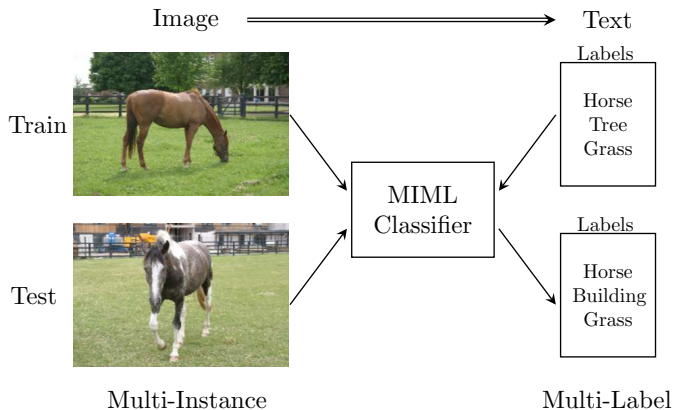
# Image Tagging as a MIML Problems



Figure: Multi-instance multi-label learning for image tagging.

# Outline

# CNNs and MIML Problems

**CNNs are suitable for solving MIML problems**

- The way to select candidate instances inside a bag (**Mulit-instance**)
  **Traditional algorithms**: extracting a bag of instances from each image and then applying multi-instance classifiers.
  **CNNs**: convolutional layers slide through an image (the bag) and create the candidate instances, and the max-pooling layers select the most representative instances inside the bag.

- The way to depict the relationship between the labels (**Multi-label**)
  **Traditional algorithms**: solve multi-label problems by explicitly or implicitly building the relationship between the labels;
  **CNNs**: by utilizing a deep hierarchical structure, different levels of representation of all the labels have already been embedded in the network, and the relationship of the labels can be explored via checking the sharing of the neurons.

# CNNs and MIML Problems

**How to apply CNNs to MIML problems?**

- **A simple way** to adapt CNNs for MIML problems
  Change the original multi-class classifier in the Softmax Layer into a
  multi-label classifier, such as *the multiple binary-class logistic regression (LR)*.

- **Shortage**
  Most existing MIML datasets only have
  a relatively small amount of training data, as the annotation costs for
  them are generally much higher. Hence, it will be hard to directly
  train an effective CNN model for such cases.

# Outline

# Transfer CNNs to small MIML tasks

1. **Extract Features from pre-trained CNN models**
   Based on the VGG 16 layers network (VGG16) pre-trained on Imagenet, we extract features from each group of its layers, which enable a depiction of the *multi-level relationship between the labels*.
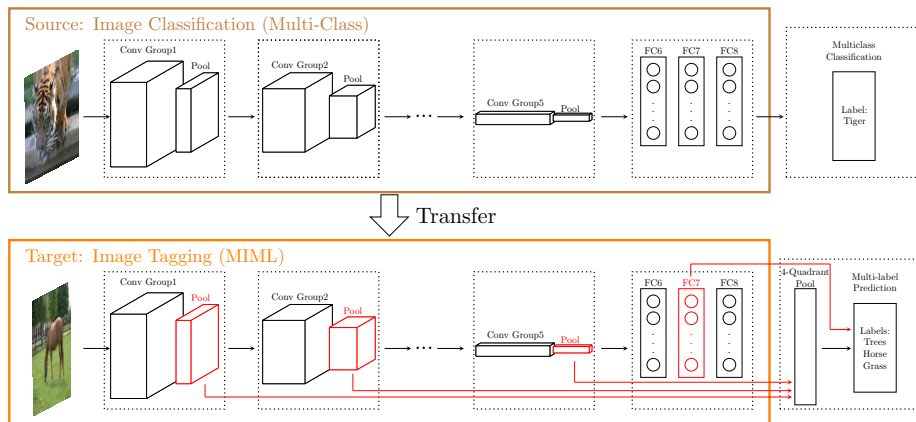
2. $L_1$-**norm regularized Logistic Regression ($L_1$LR)**
   Then an $L_1$-norm regularized Logistic Regression ($L_1$LR) is adopted to learn one classifier for each label. The aim of utilizing the sparsity regularization is to encourage the classifiers to select only much smaller subsets of "effective" features for specific labels.
   Learning the $j$-th classifier can be expressed as

$$\min_{\mathbf{w}_j} \frac{1}{m} \sum_{i=1}^{m} \log(1 + \exp(-y_{ij}(\mathbf{x}_i^T \mathbf{w}_j + b))) + \lambda |\mathbf{w}_j|_1. \qquad (1)$$

# Framework Illustration



Figure: The framework of Transfer CNN to small MIML tasks. Source Task: Multi-class Classification (Image Classification), outputting one class label; Target Task: MIML Recognition (Image Tagging), outputting multiple labels.

# Outline

## Settings

- **MIML Datasets**
  **MSRC v2**, a subset of the Microsoft Research Cambridge (MSRC) image dataset with 591 images and 23 labels (on average 2.5 labels per image),
  **Scene**, a natural scene dataset with 2000 images and 5 different labels (on average 1.2 labels per image);
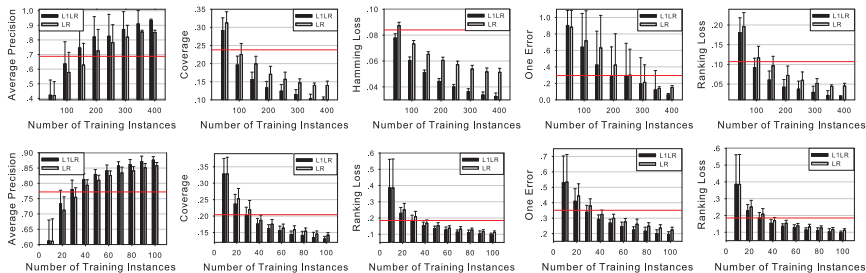
- The VGG features are extracted via **Caffe**;

- LR and $L_1$LR are implemented via **Liblinear**.

- The shrinkage parameter $\lambda$ of $L_1$LR is fixed as default of 1.

# Results

Table: Performance comparison (mean $\pm$ std). The symbol $\uparrow$ ($\downarrow$) indicates that the larger (smaller) the value, the better the performance.

| | **VGG+$L_1$LR** | **VGG+LR** | MIMLfast | DBA | KISAR |
|---|---|---|---|---|---|
| MSRC v2 | | | | | |
| a.p. $\uparrow$ | .933 $\pm$ .010 | .851 $\pm$ .016 | .688 $\pm$ .017 | .326 $\pm$ .016 | .666 $\pm$ .018 |
| co. $\downarrow$ | .102 $\pm$ .006 | .140 $\pm$ .012 | .238 $\pm$ .014 | .837 $\pm$ .018 | .254 $\pm$ .015 |
| h.l. $\downarrow$ | .033 $\pm$ .003 | .051 $\pm$ .003 | .100 $\pm$ .007 | .140 $\pm$ .006 | .086 $\pm$ .004 |
| o.e. $\downarrow$ | .060 $\pm$ .017 | .149 $\pm$ .024 | .295 $\pm$ .025 | .415 $\pm$ .026 | .341 $\pm$ .031 |
| r.l. $\downarrow$ | .018 $\pm$ .003 | .045 $\pm$ .007 | .108 $\pm$ .009 | .675 $\pm$ .017 | .131 $\pm$ .010 |
| Scene | | | | | |
| a.p. $\uparrow$ | .948 $\pm$ .006 | .926 $\pm$ .004 | .770 $\pm$ .015 | .600 $\pm$ .013 | .772 $\pm$ .012 |
| co. $\downarrow$ | .082 $\pm$ .006 | .096 $\pm$ .004 | .207 $\pm$ .012 | .334 $\pm$ .011 | .204 $\pm$ .008 |
| h.l. $\downarrow$ | .070 $\pm$ .004 | .090 $\pm$ .004 | .188 $\pm$ .009 | .269 $\pm$ .009 | .194 $\pm$ .005 |
| o.e. $\downarrow$ | .082 $\pm$ .010 | .114 $\pm$ .007 | .351 $\pm$ .023 | .386 $\pm$ .025 | .351 $\pm$ .020 |
| r.l. $\downarrow$ | .038 $\pm$ .005 | .055 $\pm$ .004 | .189 $\pm$ .014 | .348 $\pm$ .012 | .185 $\pm$ .010 |

Figure: The performance with a relatively small number of training instances. The red line indicates the best performance of the 6 state-of-the-art algorithms with 2/3 training instances (394 training instances on MSRC dataset on the upper 5 panels, and 1333 training instances on Scene dataset on the lower 5 panels). The performance is based on 30 random repetition and the bar indicates the standard deviation.

# Conclusion

- We propose a CNN-based transfer learning framework with sparsity regularization for multi-instance multi-label problems with small datasets;

- The proposal has achieved substantial improvement in classification performance when compared with the state-of-the-art MIML algorithms.