# HEp-2 Cell Image Classification With Deep Convolutional Neural Networks

**4 authors:**

Zhimin Gao
University of Wollongong
30 PUBLICATIONS   1,094 CITATIONS

SEE PROFILE

Jianjia Zhang
University of Wollongong
17 PUBLICATIONS   382 CITATIONS

SEE PROFILE

Luping Zhou
University of Wollongong
105 PUBLICATIONS   2,651 CITATIONS

SEE PROFILE

Lei Wang
Sichuan University
231 PUBLICATIONS   5,678 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:

Video Coding View project

Action Recognition View project

# HEp-2 Cell Image Classification with Deep Convolutional Neural Networks

Zhimin Gao[a], Lei Wang[a,*], Luping Zhou[a], Jianjia Zhang[a]

[a]*School of Computer Science and Software Engineering, University of Wollongong, NSW 2522, Australia*

## Abstract

Efficient Human Epithelial-2 (HEp-2) cell image classification can facilitate the diagnosis of many autoimmune diseases. This paper presents an automatic framework for this classification task, by utilizing the deep convolutional neural networks (CNNs) which have recently attracted intensive attention in visual recognition. This paper elaborates the important components of this framework, discusses multiple key factors that impact the efficiency of training a deep CNN, and systematically compares this framework with the well-established image classification models in the literature. Experiments on benchmark datasets show that i) the proposed framework can effectively outperform existing models by properly applying data augmentation; ii) our CNN-based framework demonstrates excellent adaptability across different datasets, which is highly desirable for classification under varying laboratory settings. Our system is ranked high in the cell image classification competition hosted by ICPR 2014.

*Keywords:* Indirect immunofluorescence, staining patterns classification, deep convolutional neural networks, data augmentation

## 1. Introduction

Indirect immunofluorescence (IIF) on Human Epithelial-2 (HEp-2) cells is a recommended methodology to diagnose autoimmune diseases (Rigon et al.,

---

*Corresponding author

*Email addresses:* `zg126@uowmail.edu.au` (Zhimin Gao), `leiw@uow.edu.au` (Lei Wang), `lupingz@uow.edu.au` (Luping Zhou), `jz163@uowmail.edu.au` (Jianjia Zhang)

2007). However, manual analysis of IIF images leads to crucial limitations, such as the subjectivity of result, the inconsistence across laboratories, and the low efficiency in processing a large number of cell images (Meroni and Schur, 2010; Foggia and Vento, 2013). To improve this situation, automatic and reliable cell images classification has become an active research topic.

Many methods have been recently proposed for this topic, especially during the HEp-2 cell classification competitions (Foggia and Vento, 2013; Foggia et al., 2014; Lovell et al., 2014). Most of them treat feature extraction and classification as two separate stages. For the former, a variety of hand-crafted features are adopted, including local binary pattern (LBP) (He and Wang, 1990; Nosaka and Fukui, 2014; Theodorakopoulos et al., 2014b), scale-invariant feature transform (SIFT) (Lowe, 2004), histogram of oriented gradients (Dalal and Triggs, 2005), discrete cosine transform, and the statistical features like gray-level co-occurrence matrix (Haralick et al., 1973) and gray-level size zone matrix (Thibault et al., 2014). For the latter, nearest-neighbor classifier, boosting, support vector machines (SVM) and multiple kernel SVM have been employed (Wiliem et al., 2014). As a result, the performance of these classifiers relies highly on the appropriateness of the empirically chosen hand-crafted features. Moreover, because features and classifier are treated separately, they cannot work together to maximally identify and retain discriminative information.

Very recently, deep convolutional neural networks (CNNs) have consistently achieved outstanding performance on generic visual recognition tasks (Krizhevsky et al., 2012) and this has revived extensive research interest in CNN-based classification model (Razavian et al., 2014). The CNNs consist of multi-stage processing of an input image to extract hierarchical and high-level feature representations. Many hand-crafted features and the corresponding classification pipelines can be regarded as an approximation to or a special case of the CNNs, by sharing some basic building blocks. Nevertheless, these features and pipelines have to be carefully designed and integrated in order to preserve discriminative information. The excellent performance achieved by deep CNNs on generic visual recognition and the high demand for full automation of HEp-2 cell image classification motivate us to research the CNNs for this classification task.

To this end, we propose an automatic feature extraction and classification framework for HEp-2 staining patterns based on deep CNNs (LeCun et al., 1998). This framework extracts features from the raw pixels of cell images

2

and avoids using hand-crafted features. Feature representations for each kind of staining patterns are learned and optimized via training the multi-layer network. Also, the classification layer is jointly learned with this network to predict the probability of a cell image for each class. The highly non-linear and high-capacity properties (LeCun et al., 2012) make the multi-layer CNNs difficult to train, especially when the number of training samples is not sufficiently large. We explore multiple important aspects in this CNN-based classification system, including network architecture, image preprocessing, hyper-parameters selection, and data augmentation, which are important for CNNs to achieve effective and reliable cell classification. Furthermore, we conduct rigorous experimental comparison with two state-of-the-art hand-designed shallower image representation models, i.e., bag-of-features (BoF) and Fisher Vector (FV), to investigate the advantages and disadvantages of our CNN-based framework on cell image classification. Our system has participated in the *Contest on Performance Evaluation on Indirect Immunofluorescence Image Analysis Systems* hosted by ICPR 2014[1] and won the fourth place among 11 international teams.

The rest of the paper is organized as follows. Section 2 reviews the classification models of BoF, FV and deep CNNs. In Section 3, our CNN-based framework for cell images classification is presented and a set of key factors are discussed. Section 4 reports the experimental investigation and comparison, and the conclusions are drawn in Section 5.

We were invited by the ICPR 2014 contest organizers to report our system in a workshop short paper (Gao et al., 2014). This paper significantly extends that workshop paper in the following aspects: i) a more detailed description of our deep CNN-based classification framework for HEp-2 cell images is presented and multiple key factors for effectively training a reliable deep CNN are discussed and experimentally demonstrated; ii) the role of image rotation as a data augmentation method in helping the deep CNN to achieve robust representations in this classification task is investigated and analyzed; iii) systematic experimental comparisons of our CNN-based framework and the state-of-the-art hand-designed classification models are conducted; iv) the excellent adaptability of our cell classification system with respect to different laboratory settings is demonstrated by transferring the learned network across two datasets with easy implementation, which makes

---

[1]Contest website is at `http://i3a2014.unisa.it/?page_id=91`.

our system attractive for practical clinical applications.

## 2. Related Work

### 2.1. Bag-of-features and Fisher Vector Models

The BoF model (Csurka et al., 2004) generally consists of four stages: local feature extraction, dictionary learning, feature encoding, and feature pooling. The dictionary is composed of a set of visual words describing the common visual patterns shared by local descriptors. The relationship between local descriptors and visual words is characterized by feature encoding. A variety of coding methods have been proposed in the literature (Liu et al., 2011; Wang et al., 2010; Jegou et al., 2010; Boiman et al., 2008). On top of these, spatial pyramid matching (SPM) (Lazebnik et al., 2006) is usually utilized to incorporate the spatial information of an image. The BoF model has been applied to staining patterns classification (Wiliem et al., 2014; Kong et al., 2014; Shen et al., 2014; Stoklasa et al., 2014), in which one or more of the above four stages are tailored to obtain better cell image representations for classification. Readers are referred to the review Foggia et al. (2014) for more details.

In the past several years, FV model has shown superior performance to the BoF model (Perronnin and Dance, 2007; Perronnin et al., 2010; Sánchez et al., 2013). Their main differences lie at dictionary learning and feature encoding. The dictionary in FV is generated by a probabilistic model, e.g., the Gaussian mixture model (GMM), that characterizes the distribution of local descriptors. Each local descriptor is then encoded by the first- and second-order gradients with respect to the model parameters. FV model has also been applied to cell image classification (Faraki et al., 2014; Han et al., 2014).

### 2.2. Deep Convolutional Neural Networks

CNNs belong to a class of learning models inspired by the multi-stage processes of visual cortex (Hubel and Wiesel, 1962). A pioneering work of CNNs was Fukushima's "neocognitron" (Fukushima, 1980). It has a structure similar to the hierarchical model of the visual nervous system discovered by Hubel and Wiesel (Hubel and Wiesel, 1959). Each stage of the network imitates the functions of simple and complex cells in the primary visual cortex. Later on, LeCun et al. (1998) extended the neocognitron by utilizing backpropagation algorithm to train the model parameters of CNNs and achieved excellent performance in hand-written digit recognition.

With the advent of fast parallel computing, better regularization strategies, and large-scale datasets, deep CNNs models have recently significantly outperformed the models with hand-crafted features on generic object classification, detection and retrieval (Razavian et al., 2014), as well as other visual recognition tasks, such as face verification (Taigman et al., 2014) and mitosis detection in breast cancer histopathology images (Veta et al., 2015). As for cell images classification, Malon et al. (Foggia and Vento, 2013) adopted a CNN to classify HEp-2 cell images. Buyssens et al. (2013) designed a multiscale CNN for cytological pleural cancer cells classification. Our CNN framework presented in this paper is different from their works in terms of both image preprocessing method and network architecture. Moreover, our CNN performs better than the CNN reported in Foggia and Vento (2013) on ICPR 2012 HEp-2 cell classification.

Although CNNs have been initially applied to cell image classification, the following issues have not been systematically investigated and thus remain unclear: i) what are the key issues when adopting deep CNNs for cells classification? ii) how is the performance of the CNN-based classification model when compared with the well-established classification models in the literature, especially the BoF and FV models? These issues will be carefully investigated and addressed in this work.

## 3. Proposed Framework

The proposed deep CNN-based HEp-2 cell image classification framework consists of three components: image preprocessing, network training, and feature extraction and classification, which are elaborated in this section. Also, data augmentation which plays an important role in this classification framework will be described and analyzed.

### 3.1. Network Architecture

A proper selection of network architecture is crucial to CNNs. Usually, deep CNNs are composed of multiple convolutional layers interlaced with subsampling (pooling) layers, as shown in Fig. 1. Each layer outputs a set of two-dimensional feature maps, each of which represents a specific feature detected from all positions of the input. These feature maps are in turn used as the input of the next layer. Fully-connected layers are usually stacked on the top of the network to conduct classification.
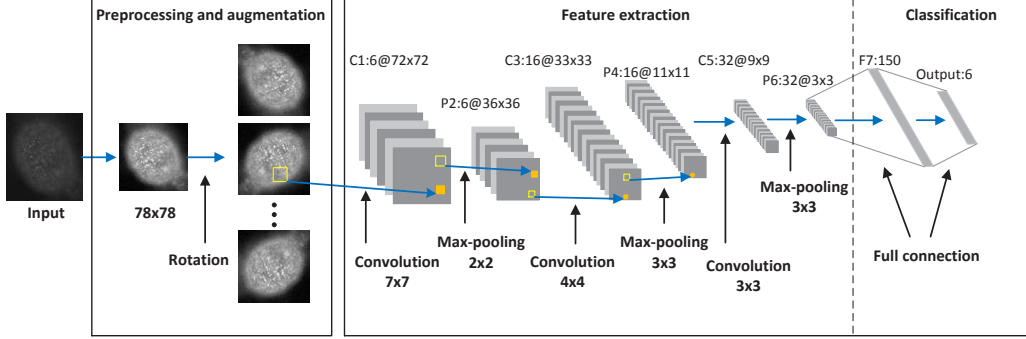
Figure 1: The architecture of our deep convolutional neural network classification system for HEp-2 cell images. Each plane within the feature extraction stage denotes a feature map. The convolutional layer and max-pooling layer is abbreviated as C and P respectively. C1:6@72 × 72 means that this is a convolutional layer, and is the first layer of the network. This layer is comprised of six feature maps, each of which has size of 72 × 72. The symbols and number above the feature maps of other layers have the similar meaning, whereas F7:150 means that this is a fully-connected layer. It is the seventh layer of the network and has 150 neurons. The words and number between two layers stand for: the operation, i.e., convolution or max-pooling, applied to the feature maps of the previous layer in order to obtain the feature maps of this layer; and the size of each filter or the size of pooling region.

Our deep CNN shares the basic architecture as the classical LeNet-5 (LeCun et al., 1998). Specifically, it contains eight layers. Among them, the first six layers are convolutional layers alternated with pooling layers, and the remaining two are fully-connected layers for classification.

### 3.1.1. Convolutional Layer

Let's assume that it is the $l$th layer. Let $N^l$ denote the number of feature maps at this layer, where $l$ is used as a superscript. Accordingly, each feature map is denoted as $\mathbf{h}_j^l$ ($j = 1, 2, ..., N^l$). This convolutional layer is parametrized by an array of two-dimensional filters $\mathbf{W}_{ij}^l$ associating the $i$th feature map $\mathbf{h}_i^{l-1}$ in the $(l-1)$th layer with the $j$th feature map $\mathbf{h}_j^l$ in the $l$th layer and the bias $b_j$. Each filter acts as a feature detector to detect one particular kind of feature by convolving with every location of the input feature map. To obtain $\mathbf{h}_j^l$, each input feature map $\mathbf{h}_i^{l-1}$ ($i = 1, 2, ..., N^{l-1}$) is firstly convolved with the corresponding filter $\mathbf{W}_{ij}^l$. The results are summed and appended with the bias $b_j^l$. After that, a non-linear activation function $\phi(\cdot)$, which can be sigmoid, tanh or rectified linear function (Krizhevsky et al.,

6

2012), is applied in an element-wise manner. Mathematically, the feature maps of the $l$th layer can be expressed as follows:

$$\mathbf{h}_j^l = \phi(\sum_{i=1}^{N^{l-1}} \mathbf{h}_i^{l-1} * \mathbf{W}_{ij}^l + b_j^l), \ j = 1, 2, ..., N^l. \tag{1}$$

where $*$ denotes the convolution operation.

### 3.1.2. Pooling Layer

A pooling layer down-samples a feature map. This will greatly reduce the computation of training a CNN and also introduces invariance to small translations of input images. Max-pooling or average-pooling is usually applied. The former selects the maximum activation over a small pooling region, while the latter uses the average activation over this region. Max-pooling generally performs better than average-pooling (Boureau et al., 2010).

### 3.1.3. Classification Layer

Classification layers usually involve one or more fully-connected layers at the top of a CNN. Our network contains two fully-connected layers. The first fully-connected layer (F7 in Fig. 1) takes the cascade of all the feature maps of the sixth layer (denoted as $\mathbf{h}^6$) as input. This layer is parametrized by weights $\mathbf{W}^7$ and biases $\mathbf{b}^7$. The output of this layer $\mathbf{h}^7$ is obtained as $\mathbf{h}^7 = \phi(\mathbf{W}^7\mathbf{h}^6 + \mathbf{b}^7)$. The last fully-connected layer is the output layer and parametrized by weights $\mathbf{W}^8$ and biases $\mathbf{b}^8$. It contains $n$ neurons corresponding to $n$ classes of staining patterns, and outputs the probabilities $\hat{\boldsymbol{y}} = [\hat{y}_1, \hat{y}_2, ..., \hat{y}_n]^\top \in \mathbb{R}^n$ via softmax regression as follows:

$$\mathbf{h}^8 = \mathbf{W}^8\mathbf{h}^7 + \mathbf{b}^8, \ \mathbf{h}^8 \in \mathbb{R}^n \tag{2}$$

$$\hat{y}_j = \frac{\exp(h_j^8)}{\sum_{i=1}^n \exp(h_i^8)}, \ j = 1, 2, ..., n. \tag{3}$$

where $\hat{y}_j$ is the output probability of the $j$th neuron.

The network architecture of our deep CNN is illustrated in Fig. 1. Specifically, the first layer convolves an input image with each of the six filters of size $7 \times 7$ with a stride of one pixel, and then adds a bias to each of them after convolution. We adopt the hyperbolic tangent function $\phi(x) = 1.7159 \tanh(\frac{2}{3}x)$ (LeCun et al., 1998) as the activation function. The second layer takes the output of the first layer as input, and applies max-pooling over non-overlapping regions of size $2 \times 2$ for each feature map. The third layer adopts

7

filters of size $4 \times 4$, and has 16 feature maps. The fourth layer then applies max-pooling over non-overlapping pooling regions of size $3 \times 3$. The fifth layer employs filters of size $3 \times 3$ and includes 32 feature maps. The sixth layer employs $3 \times 3$ non-overlapping max-pooling to the output maps of the fifth layer. After that, the resulting 32 feature maps of size $3 \times 3$ are cascaded and passed to the first fully-connected layer containing 150 neurons.

When a cell image is fed into the network, the spatial resolution of each feature map decreases as the features are extracted hierarchically from one layer to next. The spatial information of each cell is extracted by the feature maps because of the spatial convolution and pooling operations, which are important to distinct different staining pattern types. The features obtained are invariant to small translation or shift of cell images, because the filter weights of the convolutional layers are uniform for different regions of the input maps and max-pooling is robust to small variations.

## 3.2. Image Preprocessing

An appropriate image preprocessing method that takes the characteristic of images into consideration is necessary for deep CNNs to obtain good internal feature representation and classification performance.

The brightness and contrast of the HEp-2 cell images provided by the ICPR 2014 contest (ICPR2014 dataset in short) vary greatly. To reduce this variance and enhance the contrast, we normalize each image by first subtracting the minimum intensity value of the image. The resulting intensity is then divided by the difference between the maximum and minimum intensity values. Furthermore, each image is resized to $78 \times 78$ to guarantee a uniform scale of all the images used for training. This size is approximately the average size of all the cell images. Examples of six staining patterns in ICPR2014 dataset and the preprocessed images are shown in Fig. 2. In addition, we just use the preprocessed whole cell images to train our network instead of adopting a mask to only keep the foreground within each cell as Malon et al. in (Foggia and Vento, 2013), because the mask information of each cell is usually unavailable in practice, and we find that the classification performance of our system is adversely affected by using cell masks.

## 3.3. Data Augmentation

Deep CNNs are high-capacity architecture having a large number of parameters to be learned. It will be difficult to effectively train a CNN when training images are insufficient. Data augmentation (Krizhevsky et al., 2012)
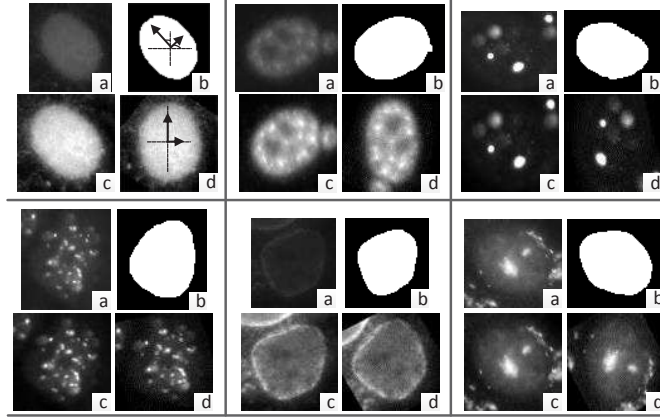
8

Figure 2: Example cells of six classes in ICPR2014 dataset and their corresponding preprocessed and aligned images. There are four images for each cell: (a) the original image; (b) the mask of this cell image (we do not take advantage of it for training the CNN); (c) the preprocessed image when the original image is applied contrast normalization and resized; (d) the aligned image when the contrast normalized image is aligned by PCA and then resized.

has been regarded as a simple and effective way to generate more samples to train a CNN and gain robustness against a variety of variances.

For data augmentation in the cell image classification, we identify the following two points: i) generating new training images by rotating existing ones can effectively boost the classification performance of the CNNs; ii) instead of merely increasing the robustness of the CNNs against the global orientation of a cell, the extra samples generated via such rotation-based augmentation help to show the intrinsic distribution of the staining patterns belonging to each cell category, which is a more important factor contributing to the improvement of the classification performance.

To demonstrate the first point, we keep rotating each training image with respect to its center by a step of $\theta$ degree. The newly generated images inherit the class label from the original training image, because rotating a cell image does not change its class label. By doing so, the original training set is enlarged by a factor of $m = \frac{360}{\theta}$, and this augmented training set is used to train the CNN.

To demonstrate the second point, we pre-align each cell image to approximately have the same global orientation. In this way, if the global orientation variance is really the main factor affecting the training performance of the

9

CNN, we shall observe some improvement by using the pre-aligned training set. Also, augmenting this pre-aligned training set with rotated images shall not lead to significantly better classification performance.

To investigate our hypothesis, we apply principal component analysis (PCA) to each cell's mask to obtain the principal direction of its shape. Each contrast normalized cell is rotated to make this principal direction to be vertical and then is resized. Applying this process to all training cell images makes them pre-aligned. These operations are illustrated in the upper left portion (as indicated) in Fig.2, followed by more examples of cell images before and after alignment. After that, we use the pre-aligned training images to train the CNN and then classify test images which are also pre-aligned.

We find that the CNN trained in this manner does not show better performance than the CNN trained with the preprocessed training images without alignment. However, when data augmentation is applied to the pre-aligned training set images, the performance of the trained CNN increases greatly. This indicates that, in terms of cell classification, adequately demonstrating the staining patterns within a cell image is more important than removing the global orientation variance[2]. Detailed experimental results will be presented in Section 4.

*3.4. Network Training*

Due to the non-convex property of the cost surface of CNNs, it is essential to select appropriate network training parameters, e.g., learning rate, and regularization methods, e.g., weight decay and dropout (Hinton et al., 2012) to make the network converge to good solutions fast.

Our deep CNN is parameterized by the weights and biases of different convolutional layers and fully-connected layers $\{\mathbf{W}^l, \mathbf{b}^l\}$, where $l = 1, 3, 5, 7, 8$. The total number of parameters is over $50,000$. The network is trained by minimizing the cross-entropy between the output probability vector $\hat{\boldsymbol{y}} = [\hat{y}_1, \hat{y}_2, ..., \hat{y}_n]^\top$ and the binary class label vector $\boldsymbol{y} = [y_1, y_2, ..., y_n]^\top$ with one non-zero entry "1" corresponding to the true class, which is expressed as

---

[2]A good example in contrast is human facial image, for which pre-alignment is generally helpful for recognition. This is because the patterns within a facial image, e.g., eyes, nose and mouth, have a rigid geometric association with the global orientation of the face. Pre-aligning the faces with respect to their global orientations effectively makes the patterns inside align with each other. Nevertheless, it is not such a case for cell images.

follows.

$$E(\boldsymbol{y}, \hat{\boldsymbol{y}}) = -\sum_{j=1}^{n} y_j \log(\hat{y}_j) \tag{4}$$

The weights are initialized from a uniform distribution and the biases are initialized to zero. All these trainable parameters are updated periodically via stochastic gradient descent (SGD) (LeCun et al., 1998) after evaluating the cost function. Let $w^l$ denote a weight of the $l$th layer, i.e., an element of $\mathbf{W}^l$. Let $b^l$ be a bias of the $l$th layer (an element of $\mathbf{b}^l$). Each weight $w^l$ and bias $b^l$ are updated by the following rules:

$$w^l := w^l - \eta \cdot \frac{\partial E}{\partial w^l}; \quad b^l := b^l - \eta \cdot \frac{\partial E}{\partial b^l} \tag{5}$$

where $\eta$ is the learning rate, and $\frac{\partial E}{\partial w^l}$ and $\frac{\partial E}{\partial b^l}$ are the partial derivatives of the cost function with respect to $w^l$ and $b^l$ respectively. They are calculated and updated via back-propagating the output error to the $l$th layer (LeCun et al., 1989) after a number of training images (a mini-batch (Bengio, 2012)) feed into the network.

To smooth the directions of gradient descent and make the network converge fast, we employ momentum (Bengio, 2012) to speed up the learning by guiding the descent direction with past gradients. The update rules of $w^l$ and $b^l$ become as the follows:

$$
\begin{aligned}
v_w^l &:= \alpha \cdot v_w^l - \beta \cdot \eta \cdot w^l - \eta \cdot \frac{\partial E}{\partial w^l}; \quad w^l := w^l + v_w^l \\
v_b^l &:= \alpha \cdot v_b^l - \eta \cdot \frac{\partial E}{\partial b^l}; \quad b^l := b^l + v_b^l
\end{aligned}
\tag{6}
$$

where $v_w^l$ and $v_b^l$ are the momentum variables for $w^l$ and $b^l$ respectively; $\alpha$ and $\beta$ are the coefficients of momentum term and weight decay term, and their optimal values are experimentally tuned, as shown in Section 4. When training error rate becomes stabilized, the learning rate $\eta$ will be reduced to achieve finer learning. The whole training process terminates after the classification error rates of both training set and validation set (which is held out from the given training images) plateau at some epochs.

In addition, another newly developed regularization strategy, dropout (Hinton et al., 2012), is also investigated in the network training. It randomly sets a fraction of the activations in the hidden layers to zero to force

the hidden units to learn more independent and robust features that could generalize well and to prevent overfitting.

### 3.5. Feature Extraction and Classification

When classifying a test image, the same preprocessing and rotation in Section 3.2 and 3.3 are applied. This results in $m$ rotated variants in total. Each of them is forward-propagated through the network, and the probability of this image for each of the $n$ classes is obtained. To further improve the robustness of classification, we select four similar CNNs after the training process becomes stable and use them collectively for classification following Krizhevsky et al. (2012). The predicted class is the one having the maximum output probability averaged over the $4m$ probabilities, that is,

$$\hat{l} = \arg\max_j \hat{y}_j = \arg\max_j \frac{1}{4m} \sum_{k=1}^{m} \sum_{i=1}^{4} \hat{y}_{ik}, \; j = 1, 2, ..., n. \tag{7}$$

## 4. Experimental Results

We evaluate our CNN classification system on two datasets of HEp-2 cell classification competition held by ICPR 2014 and 2012. The evaluation criterion is the mean class accuracy (MCA) newly adopted by ICPR 2014 competition. It is the average of the per-class accuracies (Lovell et al., 2014) defined as follows:

$$\text{MCA} = \frac{1}{n} \sum_{k=1}^{n} \text{CCR}_k \tag{8}$$

where $\text{CCR}_k$ is the classification accuracy of class $k$ and $n$ is the number of cell classes.

The average classification accuracy (ACA), which is the overall correct classification rate of all the cell images, used by the previous competition is also calculated for the ease of comparison.

### 4.1. Introduction of the HEp-2 Cell Datasets

**ICPR2014 cell dataset.** This dataset contains $13,596$ training cell images, and the test set is reserved by the competition organizers and not published yet. The cell images are extracted from 83 specimen images captured by monochrome high dynamic range cooled microscopy camera fitted on a microscope with a plane-Apochromat $20\times/0.8$ objective lens and an LED

illumination source (Lovell et al., 2014). These specimen images have been automatically segmented by using the DAPI channel and manually annotated by specialists. Each image belongs to one of the six staining patterns: *Homogeneous, Speckled, Nucleolar, Centromere, Nuclear Membrane* and *Golgi*, as shown in the top row of Fig. 3.

**ICPR2012 cell dataset.** It consists of $1,455$ cell images extracted from 28 specimens, which are acquired with a fluorescence microscope (40-fold magnification) coupled with 50W mercury vapor lamp and with a digital camera (Foggia and Vento, 2013). The dataset is pre-partitioned into training set (721 images) and test set (734 images). Each image belongs to one of the six classes: *Homogeneous, Coarse Speckled, Nucleolar, Centromere, Fine Speckled* and *Cytoplasmic*, as shown in the bottom row of Fig. 3.

Comparing the two datasets shows that two of the six classes are different. Specifically, two sub-categories of ICPR2012 dataset (*Fine Speckled* and *Coarse Speckled*) are merged into one category (*Speckled*) in ICPR2014 dataset, and two less frequent staining patterns appearing in daily clinical cases, *Golgi* and *Nuclear Membrane* are introduced in ICPR2014 dataset for developing more realistic HEp-2 cell classification systems. **Moreover, because the images in the two datasets are captured with different laboratory settings, a classification system that can be easily transferred from one dataset to the other one will be highly desired.**
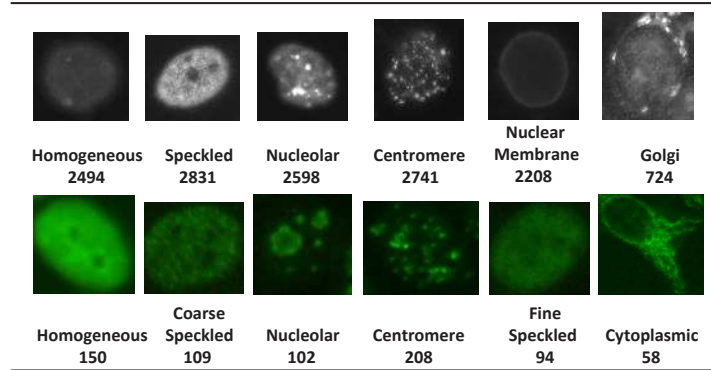


Figure 3: Comparison of HEp-2 cell images of ICPR2014 dataset (top row) and ICPR2012 dataset (bottom row). The number below the name of each cell is the total number of this kind of cells in the training set of each dataset.

## 4.2. Experiments of Hyper-parameters Optimization

This experiment demonstrates the importance of properly tuning the hyper-parameters in the CNN-based system. We categorize the hyper-parameters into two groups: model-relevant and training-relevant, as listed in Tables 1 and 2.

| Layer Number | Layer Type | Hyper-parameter |
|---|---|---|
| Layer 1 | Convolution | Filter size: $7 \times 7$ |
| | | Feature map number: 6 |
| | | Activation function: |
| | | hyperbolic tangent $\phi(x) = 1.7159 \tanh(\frac{2}{3}x)$ |
| Layer 2 | Pooling | Pooling region size: $2 \times 2$ |
| | | Pooling method: max-pooling |
| Layer 3 | Convolution | Filter size: $4 \times 4$ |
| | | Feature map number: 16 |
| | | Activation function: |
| | | hyperbolic tangent $\phi(x) = 1.7159 \tanh(\frac{2}{3}x)$ |
| Layer 4 | Pooling | Pooling region size: $3 \times 3$ |
| | | Pooling method: max-pooling |
| Layer 5 | Convolution | Filter size: $3 \times 3$ |
| | | Feature map number: 32 |
| | | Activation function: |
| | | hyperbolic tangent $\phi(x) = 1.7159 \tanh(\frac{2}{3}x)$ |
| Layer 6 | Pooling | Pooling region size: $3 \times 3$ |
| | | Pooling method: max-pooling |
| Layer 7 | Full connection | Neurons number: 150 |
| | | Activation function: |
| | | hyperbolic tangent $\phi(x) = 1.7159 \tanh(\frac{2}{3}x)$ |

Table 1: Model-relevant hyper-parameters obtained

| Hyper-parameter | Initial learning rate | Mini-batch size | Momentum coefficient | Weight decay coefficient | Dropout ratio |
|---|---|---|---|---|---|
| Value | 0.01 | 113 | 0.9 | 0.0005 | 0 |

Table 2: Training-relevant hyper-parameters obtained

To tune these hyper-parameters, we randomly partition the $13,596$ cell images of ICPR2014 dataset into three subsets, that is, 64% for training (8701

images), 16% for validation (2175 images), and 20% for test (2720 images). This partition is utilized by all experiments on ICPR2014 dataset (multiple partitions could be certainly implemented when the computational resource is not an issue.). Data augmentation is not used when tuning hyper-parameters. Following Bengio (2012), the parameters are tuned until the error rate of not only the training set but also the validation set become sufficiently small and stabilized. The hyper-parameters obtained by this tuning process are summarized in Tables 1 and 2.

We highlight that training-relevant hyper-parameters can significantly affect the convergence of cost function, the learning speed and the generalization capability of the network. Their impacts are demonstrated via the learning curves of MCA on training, validation and test sets shown from Fig. 4 to Fig. 8. In each figure, we focus on one hyper-parameter while the others are set to their optimal values in Table 2.

Fig. 4 (a) indicates that when learning rate is small, e.g., 0.001, the learning process is so slow that the MCA of the three sets have not become stable in 100 epochs. Properly increasing the learning rate effectively improves learning efficiency and the MCA becomes stable in 35 epochs, as shown in Fig. 4 (b). At the same time, an over-large learning rate, e.g., 0.1, will destabilize the learning process and degrade the classification performance. Also, Fig. 5, 6 and 7 demonstrate the impacts of mini-batch size, momentum and weight decay, respectively.

The comparison in Fig. 8 shows that the dropout strategy (Hinton et al., 2012) shall be used cautiously. When dropout with ratio of 0.5 (randomly setting the activations to zero with probability of 0.5) is applied to the first fully-connected layer of our CNN system, the learning process becomes slow and fluctuated on ICPR2014 cell dataset. A stabler and faster learning process without overfitting on the test set is gained when removing dropout, as well as better classification performance. This indicates that the neurons at the first fully-connected layer may have to work together to distinguish different staining patterns. In light of this, we decide not to employ dropout when training our network on ICPR2014 dataset.

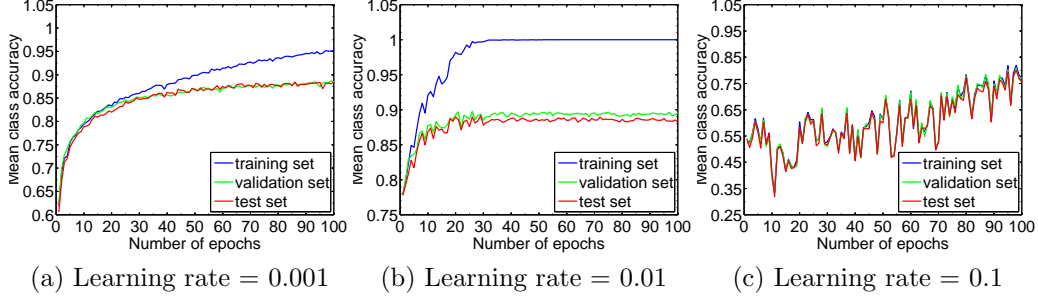(a) Learning rate = 0.001   (b) Learning rate = 0.01   (c) Learning rate = 0.1

Figure 4: Demonstration of the impact of learning rate. It shows that an over-small learning rate, e.g., 0.001, slows down the learning process, whereas an over-large learning rate, e.g., 0.1, destabilizes the learning process and degrades the classification performance. A better classification result can be obtained by properly tuning the learning rate, as shown in (b).



(a) Mini-batch size = 11   (b) Mini-batch size = 77

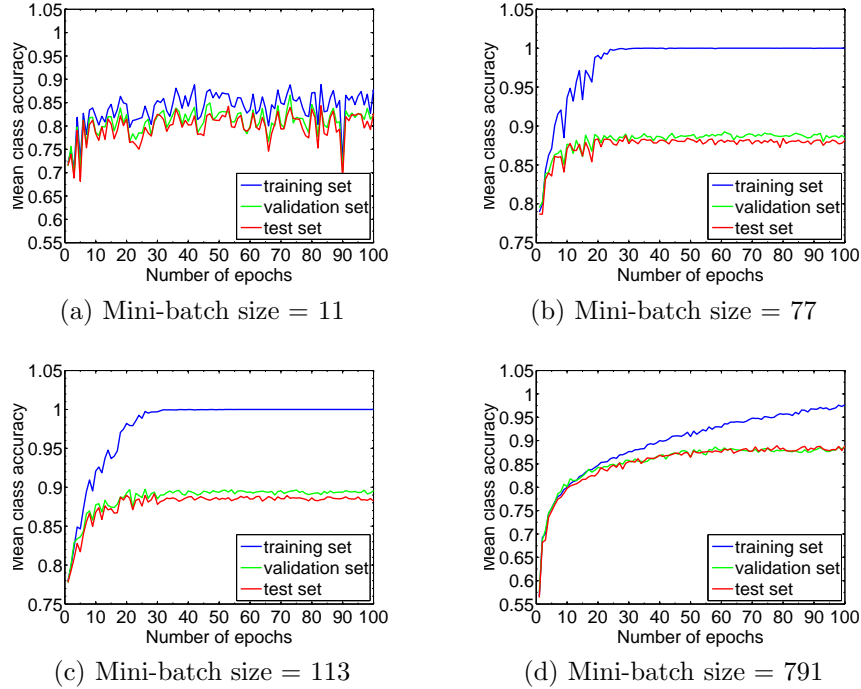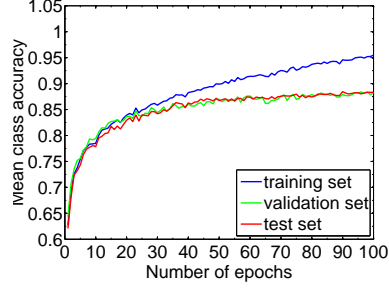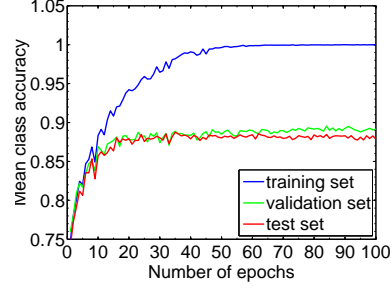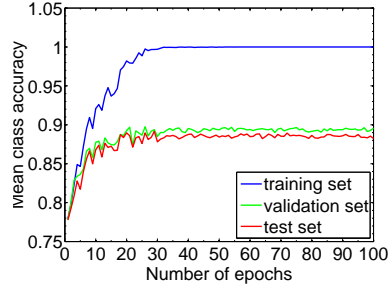(c) Mini-batch size = 113   (d) Mini-batch size = 791

Figure 5: Demonstration of the impact of mini-batch size. It shows that when mini-batch size is unnecessarily small, the learning process becomes bumpy and does not lead to the best result. On the other hand, when the mini-batch size is too large, the learning process becomes less responsive and the learning efficiency is decreased.
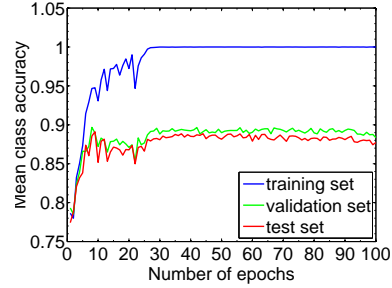
16

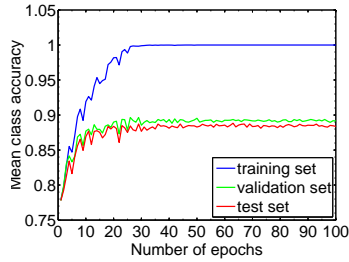(a) Momentum coefficient = 0      (b) Momentum coefficient = 0.8
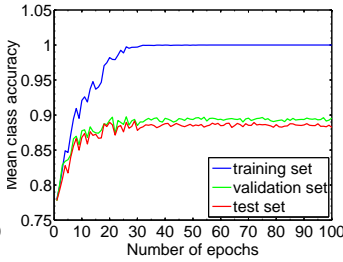
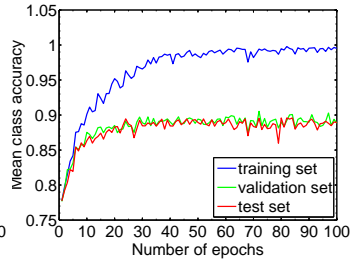(c) Momentum coefficient = 0.9      (d) Momentum coefficient = 0.97

Figure 6: Demonstration of the impact of momentum. It shows that using momentum can well accelerate the learning process. Meanwhile, a large momentum coefficient, e.g., 0.97, makes the descent direction dominated by the previous ones and causes oscillation at the initial stage. Also, it decreases the classification performance at the later stage.



(a) Weight decay coefficient = 0.00005    (b) Weight decay coefficient = 0.0005    (c) Weight decay coefficient = 0.005

Figure 7: Demonstration of the impact of weight decay. It shows that a smaller weight decay coefficient seems to be a safer choice, while a larger coefficient, e.g., 0.005, could destabilize the learning process.

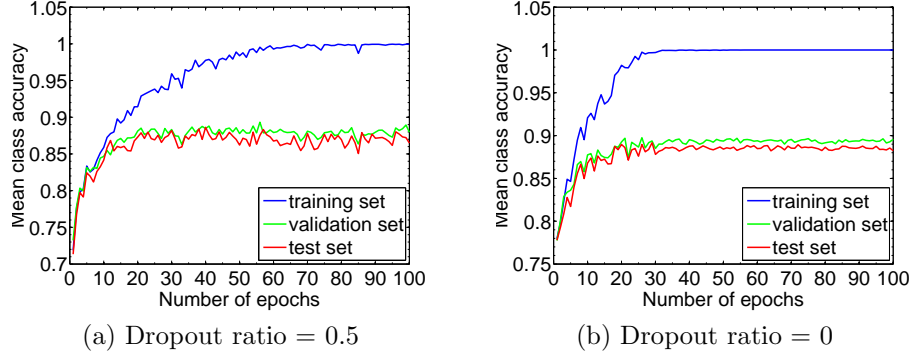(a) Dropout ratio = 0.5          (b) Dropout ratio = 0

Figure 8: Demonstration of the impact of dropout. It shows that the dropout strategy shall be used cautiously. As seen in (a), the learning process becomes slow and fluctuated on ICPR2014 cell dataset, when dropout is applied. A better learning process is obtained in (b) after removing dropout.

In sum, among the hyper-parameters of a CNN, the learning rate, mini-batch size, momentum coefficient, and weight decay coefficient can significantly impact the network training process. They have to be carefully tuned before satisfactory classification performance is obtained. For our deep CNN system, with the hyper-parameters set in Table 2, we can achieve the MCA of 89.17% on the test set of ICPR2014 dataset without using data augmentation.

### 4.3. Experiments on Data Augmentation

This experiment demonstrates the two points presented in Section 3.3, which are recapped as follows: i) the performance of the CNN can be greatly boosted by generating new training images via rotation; ii) the extra samples generated via such rotation-based augmentation help to enrich our observations of the staining patterns of each cell category for training the CNN, which is a more important factor contributing to the improvement of the classification performance than increasing robustness of the CNN against the global orientation of cells.

**Effectiveness of data augmentation.** We augment the training set by rotating each cell image for 360°, with the step of 36°, 18° and 9°, respectively. In this way, the training set is expanded by 10, 20 and 40 times, and they are used to train the CNNs, respectively. To improve the robustness of our system, we select four CNNs corresponding to the 75th, 85th,

95th and 100th epochs after the network learning becomes stable[3] as in Krizhevsky et al. (2012). A test image will go through the same rotation process as the training images and be jointly classified by the four CNNs as in Eq.(7). This system is named as "CNN". As shown in the first row of Table 3, the MCA is significantly improved (by more than 7 percentage points) from "No data augmentation" to "Augmentation by a rotation angle step of 36°". Furthermore, applying a smaller angle step to generate more training data pushes the MCA even higher, reaching 96.76%. Similar results can be observed on the ACA values. These consistent and continuous improvements well demonstrate the effectiveness and efficiency of data augmentation on cell image classification.

| Method | Accuracy (on test set) | No data augmentation | Augmentation by a rotation angle step of 36° | Augmentation by a rotation angle step of 18° | Augmentation by a rotation angle step of 9° |
|---|---|---|---|---|---|
| CNN | MCA(%) | 88.58 | 95.99 | 96.71 | **96.76** |
| | ACA(%) | 89.04 | 96.51 | 97.10 | **97.24** |
| CNN-Align | MCA(%) | 88.86 | 95.13 | 96.50 | **96.52** |
| | ACA(%) | 88.71 | 95.33 | 96.84 | **96.84** |

Table 3: Classification accuracy of our deep CNN on ICPR2014 dataset

**Data augmentation vs pre-alignment.** To gain more insight on the rotation-based data augmentation, we pre-align all the cell images with PCA as described in Section 3.3 to train the CNNs. We call this method "CNN-Align". Two experiments are conducted: i) only using these aligned images to train the CNNs without performing data augmentation; and ii) as a comparison, we further rotate each aligned training image by 360°, also with an angle step of 36°, 18° and 9°, respectively. The augmented training set is used for training. As previous, augmentation (or no augmentation) is equally applied to test images.

As shown in Table 3, when no augmentation is performed, CNN-Align does not achieve any improvement over CNN. This indicates that pre-alignment does not help here. In contrast, when training data are augmented by rotation (even with the largest angle step of 36°), CNN-Align improves significantly. This sharp change clearly demonstrates that through the rotation-based augmentation, the network can access more examples showing the diverse staining patterns within cell images. This is a more important factor

---

[3]This strategy is adopted as a model average. Different number of CNNs may be chosen, e.g. 3 or 5, to compromise between the computational expense and performance, which leads to similar classification accuracy in our experiments.

contributing to the performance improvement compared with pre-alignment that only tackles the global orientation variance of cells.

The features (filters) learned by the first and second convolutional layers of CNN corresponding to the 100th epoch trained with 9° rotated cell images are depicted as Fig. 9. It can be seen that the filters of the first convolutional layer are stain-like texture detectors. Some of the second convolutional layer filters are edge-like detectors, and most of them are also stain-like texture extractors.



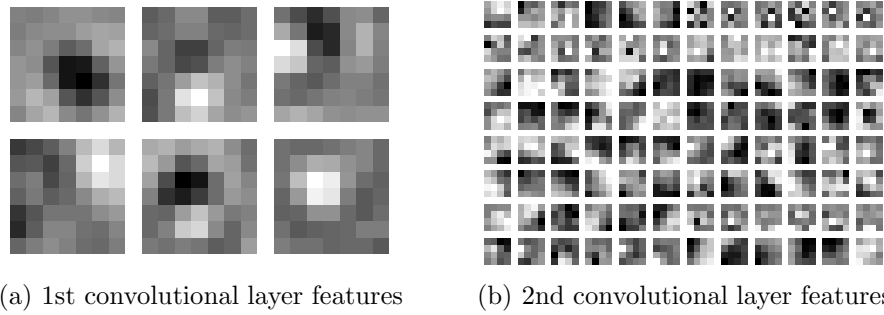(a) 1st convolutional layer features      (b) 2nd convolutional layer features

Figure 9: The features learned by the first and second convolutional layers. In general, most of the filters are stain-like texture detectors, and some are edge-like extractors.

| | Homogeneous | Speckled | Nucleolar | Centromere | Nuclear Membrane | Golgi |
|---|---|---|---|---|---|---|
| Homogeneous | 98.47 | 1.15 | 0.00 | 0.00 | 0.00 | 0.38 |
| Speckled | 2.43 | 94.26 | 1.39 | 0.52 | 1.22 | 0.17 |
| Nucleolar | 0.38 | 0.57 | 98.87 | 0.00 | 0.00 | 0.19 |
| Centromere | 0.00 | 1.03 | 1.03 | 97.95 | 0.00 | 0.00 |
| Nuclear Membrane | 0.23 | 0.23 | 0.23 | 0.00 | 98.87 | 0.45 |
| Golgi | 0.00 | 0.61 | 4.85 | 0.00 | 2.42 | 92.12 |

Figure 10: Confusion matrix of our best CNN (9° rotation) (%).

In addition, the confusion matrix of the best CNN (trained with the rotation angle step of 9°) is shown in Fig. 10. The overall classification

20

performance is very promising. The staining patterns *Nucleolar* and *Nuclear Membrane* obtain the highest classification accuracy (both 98.87%), which means that they are well separated from the others. The maximum misclassification rate (4.85%) happens to *Golgi* cells. They are easy to be misclassified as *Nucleolar* cells, because both patterns consist of a few large dots within the cells (see misclassification examples in Fig. 11). Also, *Golgi* can be confused with *Nuclear Membrane*. This may be because when the large dots within *Golgi* cells are at the edge, they will look like the *Nuclear Membrane* cells having ring-like edges. In addition, the *Speckled* cells are easy to be misclassified as *Homogeneous* cells, probably because the densely distributed speckles are the main signatures for both patterns. Misclassification examples of these staining patterns are shown in Fig. 11.



Figure 11: Misclassification examples of the three highest misclassification rates in the confusion matrix of Fig. 10. Every two rows form a group, and the first row shows cells that are misclassified to the cell type of the second row.

### 4.4. Comparison with the BoF and Fisher Vector Models

**Experimental setting.** To ensure a fair comparison, the same image preprocessing in our CNN model is equally used in both models. For each cell image, SIFT descriptors are extracted from densely sampled patches with a

stride of two pixels. The visual dictionary is generated by applying the $k$-means clustering to the descriptors extracted from training images. Local soft-assignment coding (LSC) (Van Gemert et al., 2008; Liu et al., 2011) is employed to encode the SIFT descriptors. SPM is used to partition each image into $1 \times 1$, $2 \times 2$ and $1 \times 3$ regions, and max-pooling is applied to extract representations from each region.

A similar setting is applied to the FV model. In addition, the 128-dimensional SIFT descriptors are decorrelated and reduced to dimensions of 64 by PCA as in Sánchez et al. (2013). A GMM is then estimated to represent the visual dictionary. Afterwards, each PCA-reduced SIFT descriptor is encoded with the improved Fisher encoding (Perronnin et al., 2010), where the signed square-root and $l^2$-normalization are applied to the coding vector. SPM with four regions ($1 \times 1$ and $1 \times 3$) are adopted (Sánchez et al., 2013). Following the literature, a multi-class linear SVM classifier is used in the BoF and FV models. In our implementation of BoF and FV, the publicly available VLFeat toolbox (Vedaldi and Fulkerson, 2010) is used.

**Parameter setting.** There are two primary parameters in the BoF and FV models: patch size and dictionary size (or equally, the number of components of the GMM in the FV model). We tune these parameters by five-fold cross-validation on the union of training and validation sets, with the criterion of MCA. The candidate patch sizes are $9 \times 9$, $11 \times 11$, $13 \times 13$, $15 \times 15$ and $20 \times 20$, while the candidate dictionary sizes are $1,000$, $2,000$, $3,000$, $4,000$, $5,000$ and $10,000$. Also, the number of Gaussian components will be chosen from 64, 128, 256, 512 and 1024 for FV. Through the cross-validation, the patch size and the dictionary size in the BoF model are selected as $15 \times 15$ and $10,000$. With the use of SPM, this results in a $80,000$-dimensional representation for each cell image. For the FV model, the patch size is chosen as $20 \times 20$ and the number of GMM components is 512. With the use of SPM, this leads to a $262,144$-dimensional representation for each image.

**Comparison results.** The BoF, FV and CNN models are compared on the same training and test sets. Also, both of the cases, i.e., with and without data augmentation, are investigated. To be fair, when data augmentation is used, the visual dictionary in the BoF and FV models will be built with the augmented training set. Also, to keep consistent with the setting of our deep CNN system, each test image in this case will be equally augmented and its label is predicted in the way similar to Eq.(7), except that the probabilities are replaced by the decision values of the linear SVM classifier.

As shown in Table 4, FV is consistently better than BoF, regardless of

whether data augmentation is applied or not. This agrees well with the literature. Furthermore, both BoF and FV can well benefit from data augmentation, with an average performance increase of about 4 percentage points. Compared with BoF and FV, CNN system shows slightly lower performance (88.85% vs 89.83% for BoF and 91.60% for FV), when there is no augmentation. However, CNN outperforms both BoF and FV once data augmentation is applied. In specific, the highest MCA, 96.76%, is obtained by our CNN, while BoF and FV achieve only 94.23% and 95.73% respectively. Similar situation can be observed from the ACA values. These results suggest that i) when training samples are not sufficient, the high-capacity CNN is more difficult to train than the shallower, hand-designed models such as BoF and FV; and ii) by properly using data augmentation to generate more training data, the CNN can be better trained and are able to achieve better performance than the BoF and FV models.

| Accuracy (on test set) | Methods | No data augmentation | Augmentation by a rotation angle step of 36° | Augmentation by a rotation angle step of 18° | Augmentation by a rotation angle step of 9° |
|---|---|---|---|---|---|
| MCA (%) | BoF | 89.83 | 94.23 | 93.98 | 94.14 |
| | FV | **91.60** | 95.41 | 95.73 | 95.53 |
| | CNN | 88.58 | **95.99** | **96.71** | **96.76** |
| ACA (%) | BoF | 90.70 | 94.30 | 94.19 | 94.38 |
| | FV | **92.65** | 95.78 | 96.07 | 95.81 |
| | CNN | 89.04 | **96.51** | **97.10** | **97.24** |

Table 4: Comparison of classification accuracy among the methods of BoF, FV and our deep CNN on ICPR2014 datatset

### 4.5. Experiments on the Adaptability across Datasets

As previously mentioned, HEp-2 cell image classification varies with laboratory settings, the types of staining patterns involved, and the size of dataset. Such differences can be well seen from the ICPR2014 and ICPR2012 datasets. As a result, it is highly desired that a cell classification system trained with one dataset can be conveniently adapted to another one. Owning this feature not only improves the efficiency of system building, but also can take full advantages of the image data in different datasets. To demonstrate this feature for our CNN-based system, we compare the CNN purely trained on ICPR2012 dataset (called CNN-Standard in short) with the other CNN which is an adapted version of the CNN pre-trained on ICPR2014 dataset to ICPR2012 dataset (called CNN-Finetuning).

Following previous experimental settings, CNN-Standard is trained with the 721 training images predefined in ICPR2012 dataset. Only the green

channel of each image is kept and the same preprocessing in Section 3.2 is performed. The dropout strategy (with ratio of 0.5) is used, because it can benefit network training and classification performance on this small dataset. CNN-Standard is trained by 100 epochs and then used to classify the predefined test images by following Eq.(7).

To train CNN-Finetuning, we first select a basic CNN system learned with the ICPR2014 dataset. It is the one obtained at the 100th epoch when the system is trained with an augmented (rotation with an angle step of 9°) training set of ICPR2014. Afterwards, this basic system is fine-tuned with the training set of ICPR2012 dataset, with or without data augmentation. All the trainable network parameters of different layers are updated during this fine-tuning process. To demonstrate the efficiency, we only fine-tune this basic system by 10 epochs, which takes significantly less time than the 100 epochs spent in training CNN-Standard.
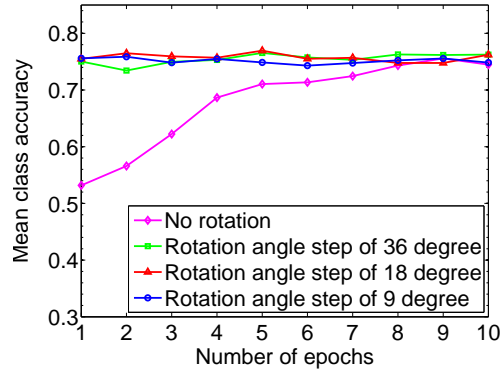


Figure 12: The MCA of test set obtained by CNN-Finetuning at each of the 10 epochs. Data augmentation with various angle steps is investigated.

The evolution of the MCA on test set with the 10 epochs is plotted in Fig. 12. As shown by the line of "No rotation", CNN-Finetuning does not work well at the beginning. Nevertheless, it catches up quickly in a couple of epochs and reaches a satisfying performance in 10 epochs. Furthermore, the adaption stage is significantly shortened, by applying data augmentation to the small training set of ICPR2012 to increase training samples. These results demonstrate the high efficiency of the adaptability of our CNN-based system, especially considering that there are two different classes of staining patterns across these datasets. Comparison of CNN-Standard and CNN-

24

Finetuning is shown in Table 5. It is interesting to note that CNN-Finetuning consistently outperforms CNN-Standard, even though it is only fine-tuned for a few epochs. We attribute its superiority to the good initialization of the network obtained from the training process on ICPR2014 dataset. Based on the above results, we believe that our CNN-based system will be a better option for practical applications.

| Accuracy (on test set) | Methods | No data augmentation | Augmentation by a rotation angle step of 36° | Augmentation by a rotation angle step of 18° | Augmentation by a rotation angle step of 9° |
|---|---|---|---|---|---|
| MCA (%) | CNN-Standard | 63.1 | 72.4 | 72.4 | 73.2 |
|  | CNN-Finetuning | 74.5 | 76.3 | 76.2 | 74.9 |
| ACA (%) | CNN-Standard | 64.3 | 70.2 | 70.0 | 70.1 |
|  | CNN-Finetuning | 72.9 | 74.8 | 74.7 | 73.3 |

Table 5: Classification accuracy of our CNN-based system on ICPR2012 dataset

At last, we compare our CNN-Finetuning (rotation with an angle step of 36°) with other methods reported in the literature in Table 6. As seen, it outperforms the best-performing method of that contest and the CNN at the ICPR2012 contest. For that CNN, a $100 \times 100$ pixels area of the green channel centered at the largest connected component of each cell is taken via the mask and then is normalized by mapping the first and 99th percentile values to 0 and 1. The architecture of that CNN is composed of two sequences of convolution, absolute value rectification and subtractive normalization, one average pooling layer, one max pooling layer and one fully connected layer[4], which is also quite different from our architecture. The better performance of our CNN may benefit from these differences as well as our effective data augmentation. Also, our CNN-Finetuning is just slightly inferior to the method in Theodorakopoulos et al. (2014b). That method combines two kinds of hand-crafted features: the distribution of SIFT and gradient-oriented co-occurrence LBP, and a dissimilarity representation of an image is created with them.

---

[4]Please refer to the contest report available at `http://mivia.unisa.it/hep2contest/HEp-Contest_Report.pdf` for the detailed presentation of the contest CNN.

| Method | Average classification accuracy (ACA) |
|---|---|
| 2012 contest best-performing method (Foggia and Vento, 2013) | 68.7% |
| 2012 contest CNN (Foggia and Vento, 2013) | 59.8% |
| Nosaka and Fukui (2014) | 68.5% |
| Shen et al. (2014) | 74.4% |
| Faraki et al. (2014) | 70.2% |
| Larsen et al. (2014) | 71.5% |
| Theodorakopoulos et al. (2014b) | **75.1**% |
| Our CNN-Finetuning | **74.8**% |

Table 6: Comparison with other methods on the ICPR2012 dataset

In addition, it is worth mentioning that in the ICPR2014 contest (Lovell et al., 2014), the three methods that perform better than or comparable to our deep CNNs system (87.10%, 83.64% and 83.33% vs 83.23% with the MCA criterion) are all built on two-stage frameworks: hand-designed feature representation and classification. The top-ranked method utilizes multi-scale and multiple types of local descriptors (Manivannan et al., 2014); the second-ranked method adopts the hand-crafted rotation invariant dense scale local descriptor (Gragnaniello et al., 2014); and the third method combines morphological features and different local texture features (Theodorakopoulos et al., 2014a). In contrast, our CNN system generates discriminative features from raw pixels directly by utilizing class label information and jointly learns the classifier in a single architecture without learning extra dictionaries as these methods.

*4.6. Discussion on Computational Issues*

For the CNN-based classification system, training the network is the most time-consuming step in the whole pipeline. However, this process can be well accelerated by utilizing GPU programming. Also, as previously shown, an existing CNN-based system can be efficiently transferred to a new but related task via a short training process. Once the networks are trained, a test cell image only needs to go through the four networks and then is classified within 1.2 seconds in total with Matlab implementation on a computer with 3.30GHz Intel CPU and 16GB RAM.

For the BoF and FV models, building visual dictionary or the GMM is computationally intensive, especially when there are a large number of training images, e.g., due to the use of data augmentation. For example, building a dictionary of $10,000$ visual words and the GMM of $512$ components

takes more than 4 days and 2 days in our implementation, when the training set of ICPR2014 dataset is augmented by rotation with an angle step of 9°. Also, a large dictionary in the BoF model could slow down the encoding process, e.g., around 78 seconds per image in our experiment. Although the time for this process can be reduced in the FV model, it still takes about three seconds per image. In addition, SPM is usually needed to attain better classification performance. In this case, the dimensions of the resulting image representation are much higher than that in the CNN-based system ($80,000$ or $262,144$ vs $150$ only).

## 5. Conclusion

This paper proposes an automatic HEp-2 cell staining patterns classification framework with deep convolutional neural networks. We give a detailed description on various aspects of this framework and carefully discuss a number of key issues that could affect its classification performance. Extensive experimental study on two benchmark datasets demonstrates i) the advantages of our framework over the well-established image classification models on cell image classification; ii) the importance and effectiveness of data augmentation, especially when training images are not sufficient; iii) the desirable adaptability of our CNN-based system across different datasets, which makes our system attractive for practical tasks. Much future work can be done to further improve the performance of the proposed system. In particular, a super-CNN trained with a large-scale generic image benchmark, ImageNet (Deng et al., 2010), has recently prevailed on many generic visual recognition tasks. We would like to explore the effectiveness of the features generated by this CNN for HEp-2 cell image and the adaption of this CNN to cell image classification. These issues will be of significance considering the substantial differences between generic images and HEp-2 cell images.

## References

## References

Bengio, Y., 2012. Practical recommendations for gradient-based training of deep architectures, in: Neural Networks: Tricks of the Trade. Springer, pp. 437–478.

Boiman, O., Shechtman, E., Irani, M., 2008. In defense of nearest-neighbor based image classification, in: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, IEEE. pp. 1–8.

Boureau, Y.L., Ponce, J., LeCun, Y., 2010. A theoretical analysis of feature pooling in visual recognition, in: Proceedings of the 27th International Conference on Machine Learning (ICML-10), pp. 111–118.

Buyssens, P., Elmoataz, A., Lézoray, O., 2013. Multiscale convolutional neural networks for vision–based classification of cells, in: Computer Vision–ACCV 2012. Springer, pp. 342–352.

Csurka, G., Dance, C., Fan, L., Willamowski, J., Bray, C., 2004. Visual categorization with bags of keypoints, in: Workshop on statistical learning in computer vision, ECCV, pp. 1–2.

Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human detection, in: Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, pp. 886–893 vol. 1.

Deng, J., Berg, A.C., Li, K., Fei-Fei, L., 2010. What does classifying more than 10,000 image categories tell us?, in: Computer Vision–ECCV 2010. Springer, pp. 71–84.

Faraki, M., Harandi, M.T., Wiliem, A., Lovell, B.C., 2014. Fisher tensors for classifying human epithelial cells. Pattern Recognition 47, 2348–2359.

Foggia, P., Percannella, G., Saggese, A., Vento, M., 2014. Pattern recognition in stained hep-2 cells: Where are we now? Pattern Recognition 47, 2305–2314.

Foggia, P., P.G.S.P., Vento, M., 2013. Benchmarking hep-2 cells classification methods. Medical Imaging, IEEE Transactions on 32, 1878–1889.

Fukushima, K., 1980. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. Biological Cybernetics 36, 193–202.

Gao, Z., Zhang, J., Zhou, L., Wang, L., 2014. Hep-2 cell image classification with convolutional neural networks, in: Pattern Recognition Techniques for Indirect Immunofluorescence Images (I3A), 2014 1st Workshop on, pp. 24–28.

Gragnaniello, D., Sansone, C., Verdoliva, L., 2014. Biologically-inspired dense local descriptor for indirect immunofluorescence image classification, in: Pattern Recognition Techniques for Indirect Immunofluorescence Images (I3A), 2014 1st Workshop on, pp. 1–5.

Han, X.H., Wang, J., Xu, G., Chen, Y.W., 2014. High-order statistics of microtexton for hep-2 staining pattern classification. Biomedical Engineering, IEEE Transactions on 61, 2223–2234.

Haralick, R., Shanmugam, K., Dinstein, I., 1973. Textural features for image classification. Systems, Man and Cybernetics, IEEE Transactions on SMC-3, 610–621.

He, D.C., Wang, L., 1990. Texture unit, texture spectrum, and texture analysis. Geoscience and Remote Sensing, IEEE Transactions on 28, 509–512.

Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.R., 2012. Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580 .

Hubel, D.H., Wiesel, T.N., 1959. Receptive fields of single neurones in the cat's striate cortex. The Journal of physiology 148, 574–591.

Hubel, D.H., Wiesel, T.N., 1962. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of physiology 160, 106.

Jegou, H., Douze, M., Schmid, C., Perez, P., 2010. Aggregating local descriptors into a compact image representation, in: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, pp. 3304–3311.

Kong, X., Li, K., Cao, J., Yang, Q., Wenyin, L., 2014. Hep-2 cell pattern classification with discriminative dictionary learning. Pattern Recognition 47, 2379–2388.

Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks., in: NIPS, p. 4.

Larsen, A., Vestergaard, J., Larsen, R., 2014. Hep-2 cell classification using shape index histograms with donut-shaped spatial pooling. Medical Imaging, IEEE Transactions on 33, 1573–1580.

Lazebnik, S., Schmid, C., Ponce, J., 2006. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories, in: Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on, IEEE. pp. 2169–2178.

LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W., Jackel, L.D., 1989. Backpropagation applied to handwritten zip code recognition. Neural computation 1, 541–551.

LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. Proceedings of the IEEE 86, 2278–2324.

LeCun, Y.A., Bottou, L., Orr, G.B., Müller, K.R., 2012. Efficient backprop, in: Neural networks: Tricks of the trade. Springer, pp. 9–48.

Liu, L., Wang, L., Liu, X., 2011. In defense of soft-assignment coding, in: Computer Vision (ICCV), 2011 IEEE International Conference on, IEEE. pp. 2486–2493.

Lovell, B.C., Percannella, G., Vento, M., Wiliem, A., 2014. Performance evaluation of indirect immunofluorescence image analysis systems. ICPR 2014 URL: `http://i3a2014.unisa.it/`.

Lowe, D.G., 2004. Distinctive image features from scale-invariant keypoints. International journal of computer vision 60, 91–110.

Manivannan, S., Li, W., Akbar, S., Wang, R., Zhang, J., McKenna, S., 2014. Hep-2 cell classification using multi-resolution local patterns and ensemble svms, in: Pattern Recognition Techniques for Indirect Immunofluorescence Images (I3A), 2014 1st Workshop on, pp. 37–40.

Meroni, P.L., Schur, P.H., 2010. Ana screening: an old test with new recommendations. Annals of the rheumatic diseases 69, 1420–1422.

Nosaka, R., Fukui, K., 2014. Hep-2 cell classification using rotation invariant co-occurrence among local binary patterns. Pattern Recognition 47, 2428–2436.

Perronnin, F., Dance, C., 2007. Fisher kernels on visual vocabularies for image categorization, in: Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on, IEEE. pp. 1–8.

Perronnin, F., Sánchez, J., Mensink, T., 2010. Improving the fisher kernel for large-scale image classification, in: Computer Vision–ECCV 2010. Springer, pp. 143–156.

Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S., 2014. Cnn features off-the-shelf: An astounding baseline for recognition, in: Computer Vision and Pattern Recognition Workshops (CVPRW), 2014 IEEE Conference on, pp. 512–519.

Rigon, A., Soda, P., Zennaro, D., Iannello, G., Afeltra, A., 2007. Indirect immunofluorescence in autoimmune diseases: assessment of digital images for diagnostic purpose. Cytometry Part B: Clinical Cytometry 72, 472–477.

Sánchez, J., Perronnin, F., Mensink, T., Verbeek, J., 2013. Image classification with the fisher vector: Theory and practice. International journal of computer vision 105, 222–245.

Shen, L., Lin, J., Wu, S., Yu, S., 2014. Hep-2 image classification using intensity order pooling based features and bag of words. Pattern Recognition 47, 2419–2427.

Stoklasa, R., Majtner, T., Svoboda, D., 2014. Efficient k-nn based hep-2 cells classifier. Pattern Recognition 47, 2409–2418.

Taigman, Y., Yang, M., Ranzato, M., Wolf, L., 2014. Deepface: Closing the gap to human-level performance in face verification, in: Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on, pp. 1701–1708.

Theodorakopoulos, I., Kastaniotis, D., Economou, G., Fotopoulos, S., 2014a. Hep-2 cells classification using morphological features and a bundle of local gradient descriptors, in: Pattern Recognition Techniques for Indirect Immunofluorescence Images (I3A), 2014 1st Workshop on, pp. 33–36.

Theodorakopoulos, I., Kastaniotis, D., Economou, G., Fotopoulos, S., 2014b. Hep-2 cells classification via sparse representation of textural features fused into dissimilarity space. Pattern Recognition 47, 2367–2378.

Thibault, G., Angulo, J., Meyer, F., 2014. Advanced statistical matrices for texture characterization: Application to cell classification. Biomedical Engineering, IEEE Transactions on 61, 630–637.

Van Gemert, J.C., Geusebroek, J.M., Veenman, C.J., Smeulders, A.W., 2008. Kernel codebooks for scene categorization, in: Computer Vision–ECCV 2008. Springer, pp. 696–709.

Vedaldi, A., Fulkerson, B., 2010. Vlfeat: An open and portable library of computer vision algorithms, in: Proceedings of the international conference on Multimedia, ACM. pp. 1469–1472.

Veta, M., van Diest, P.J., Willems, S.M., Wang, H., Madabhushi, A., Cruz-Roa, A., Gonzalez, F., Larsen, A.B., Vestergaard, J.S., Dahl, A.B., Cirean, D.C., Schmidhuber, J., Giusti, A., Gambardella, L.M., Tek, F.B., Walter, T., Wang, C.W., Kondo, S., Matuszewski, B.J., Precioso, F., Snell, V., Kittler, J., de Campos, T.E., Khan, A.M., Rajpoot, N.M., Arkoumani, E., Lacle, M.M., Viergever, M.A., Pluim, J.P., 2015. Assessment of algorithms for mitosis detection in breast cancer histopathology images. Medical Image Analysis 20, 237 – 248.

Wang, J., Yang, J., Yu, K., Lv, F., Huang, T., Gong, Y., 2010. Locality-constrained linear coding for image classification, in: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, IEEE. pp. 3360–3367.

Wiliem, A., Sanderson, C., Wong, Y., Hobson, P., Minchin, R.F., Lovell, B.C., 2014. Automatic classification of human epithelial type 2 cell indirect immunofluorescence images using cell pyramid matching. Pattern Recognition 47, 2315–2324.