

# CNN Based Image Classification

Xiaoling Long  
SIST

ShanghaiTech Univ.

longxl@shanghaitech.edu.cn

Hongyu Chen  
SIST

ShanghaiTech Univ.

chenhy3@shanghaitech.edu.cn

Second Author  
SIST

ShanghaiTech Univ.

secondauthor@i2.org

## Abstract

*The convolution networks make a great success in image classification since AlexNet proposed for classification in ImageNet challenge. This End-to-End training method make researcher focus more on how to build an efficient network. The user focus more on how to use convolution networks in their application. After AlexNet, there are many ConvNet architecture proposed, such as VGGNet and ResNet. These ConvNet architectures are widely used in many actual real applications.*

## 1. Introduction

Image classification is a fundamental challenge in computer vision [6]. Consider the problem of detecting objects from a category, such as people or cars, in static images. This is a difficult problem because objects in each category can vary greatly in appearance. Variations arise from changes in illumination, viewpoint, and intra-class variability of shape and other visual properties among object instances. For example, people wear different clothes and take a variety of poses while cars come in various shapes and colors. Classification between the objects is easy task for humans but it has proved to be a complex problem for machines. Image classification refers to the labeling of images into one of a number of predefined categories. Image classification is an important and challenging task in various application domains, including biomedical imaging, biometry, video surveillance, vehicle navigation, industrial visual inspection, robot navigation, and remote sensing [12].

Image classification uses both supervised and unsupervised Traditional methods of computer vision and machine learning cannot match human performance on tasks such as the recognition of handwritten digits or traffic signs. [4]

Nonparametric classifiers such as decision tree classifier, neural network, SVM classifier and knowledge based classification are also very common in image classification.

Disadvantages about the deep learning. For the deep learning based algorithm required the ability of labeled

samples for training. The collection of labeled data is a time consuming process as well as costly.

## 2. CNN Based Image classification

The ImageNet Large Scale Visual Recognition Challenge (ILSRC) is a benchmark in object category classification and detection on 1000-classes and millions of images. After AlexNet achieved huge success in ILSRC-2012, there are various variations of AlexNet [14] and many other types of ConvNet for image classification. Since that, ConvNet is widely used for image classification. [2] illustrate briefly what is ConvNet, the components of ConvNet, the activation function in ConvNet, from LeNet to ResNet bunch of successful ConvNet and some open issues on CNN based image classification. LeNet is first proposed CNN based image classification method.

LeNet [15] is first proposed CNN based image classification method. After that, AlexNet [14] brought Convolutional neural network into ILSRC. In this implementation, it contains 8 layers 5 convolutional and 3 fully-connected. The main features of this network's architecture are ReLU [21] as activation function, overlapping pooling and skills for reducing overfitting. Based on ReLU and overlapping pooling, the network error rate has more or less reduction, and ReLU network learn several times faster than other saturating activation function such as tanh neurons network. Overfitting is common issue for machine learning, it uses data augmentation and dropout to avoid overfit. Dropout is a skill to reduce argument or increase hypothesis. There are many discuss about this. As to data augmentation, it enlarge the dataset by cropping  $224 \times 224$  patches from the original image as well as these patches' horizontal reflection. At the end averaging the predication as final score.

[5] explored the generalization ability of ConvNet features, releasing DeCAF. Traditional image classification pipeline is extracting feature, building bag of feature then put into classifier. [38] propose a CNN based feature extractor. This is an unsupervised learning ConvNet or in other words, input image is also the kind of ground truth. After feature extracting, the final result can classify by any

classifier. This strategy is used as a tool for visualizing and understanding how ConvNet works [37]. ConvNet have an impressive classification performance. However there is no clear understanding for why this work. [37] propose a architecture for visualizing and understanding how ConvNet works.

OverFeat [29] is a integrated recognition, localization and detection. This network uses CNN extract feature from image and then perform classification and localization and detection. Multi-scale classification brought up in [29] to increase accuracy.

“Networ In Network” [18] propped a new deep network structure. Different from conventional convolution layer, it brings up a new Mlpconv layer. This Mlpconv layer consist of sliding multilayer perceptron(MLP) window. In stead of fully-connected layer at the top of network, global average pooling is used to produce the resulting vector fed directly into the softmax layer. Verified by experiments, this NIN structure indeed works well on some benchmark datasets, and global average pooling can be regarded as regularizer. THis glocal average pooling has no parameter. This strately is used widely afterwards.  $1 \times 1$  convolution conception proposed in [18] is used in GoogLeNet for dimension reduction.

AlexNet make a great success in image classification. Afterwards many various Network appear. GoogLeNet [34] proposed by google is a new level of organization in the form of the “ Inception module”. This is a multi-scale arctitecture. With the limitation of computational resource, it perform a  $1 \times 1$  convolution to dimension reduction. Auxiliary classifier is also a brilliant strategy. This smart design makes a great success in ILSRC-2014. At the same time, the widely used ConvNet architecture VGGNet [31] won the first place in *Classification + Localization competition*. It adds the number of layers up to 16 – 19. Instead of  $7 \times 7$  convolution filter in [31], it uses  $3 \times 3$  as convolution filter. After multiple layers, it can get similar effect as  $7 \times 7$  one. This design significantly reduces the parameters, and then reduces the overfitting. It also means the number of layers significantly increases. Altering convolutional layers and poolint layer became a common used Network architecture.

As the depth of ConvNet increasing, training gets more and more difficult. The training of very deep network becomes a open issue in CNN. Highway Networks [32, 33] propose *information highways* which allow unimpeded information flow across several layers. The *transform gate*  $T(x, W_T)$  and the *carry gate*  $C(x, W_C)$  proposed for decided how much flow pass through to output. The new model given by

$$y_{output} = H(x, W_H) \cdot T(x, W_T) + x \cdot C(x, W_C) \quad (1)$$

For simipcty, [10] set  $C = 1 - T$ . This design make train- ing hundreds of layers be possible and the err rate just has

slightly increase. This architecture promote the success of ResNet [10]. ResNet has similar structure as deep plain network stacked by dozens of convolution layers followed by global pooling layer and 1 fully-connected layer. except shortcut connection. This design has a residual representation which called deep residual learning. This archicture keeps parameter less than VGG-19 model even the network has 152 layers. This smart design make ResNet won first place in ILSRC-15.

There are several works based on region based image classification. Such as [8], [7] [36], [1]. Region based methods are computationally expensively.

R-CNN and fast R-CNN is regioned based convolutional network method for object detection and image classification. In [7], they use the selective search to select those propose regions. Features are extracted from each proposal region. Then a SVM classifier is used for the category classification. This kind of algorithm need a lot of time to process each image. About 2000 regions are proposed from each image where there maybe several objects in the image.

Compared to R-CNN [7], Fast R-CNN [8] employs several innovations to improve training and testing speed while also increasing detection accuracy. the fast R-CNN has several advantages: (1) higher detection quality than R-CNN 2 using a multi task loss ,predict the object and its confidence. No disk storage is required for additional feature catching. The computation is shared during training. Fast achieves a near real time rates uses a very deep network.

Proposal based image classification also contains some great work. Such as [26] [23] [24] [25].

Faster R-CNN [26] is another convolutional network which based on the region proposal methods. In [26], The author show that an algorithmic change computing proposals with a deep convolutional neural network leads to an elegant and effective solution where proposal computation is nearly cost-free given the detection networks computation. They removed the select search algorithm replaced by a region proposal network. The high quality proposal is used by the Fast R-CNN network for detection and classification. For the very deep VGG-16 model the detection system has a frame rate of 5fps (including all steps) on a GPU, while achieving state-of-the-art object detection accuracy.

Fast and Faster R-CNN focus on speeding up the R-CNN framework by sharing computation and using neural networks to propose regions instead of Selective Search. While they offer speed and accuracy improvements over R-CNN, both still fall short of real-time performance. Fast R-CNN speeds up the classification stage of R-CNN but it still relies on selective search which can take around 2 seconds per image to generate bounding box proposals.

In [23], the author achieves a real time detection and classification in 45 frames per second. A smaller version of their network can be achieved by 155 frames per second.

The problem is view as a regression task. The output of the net contains five parameters.  $(x, y, w, h, confidence)$  where  $x, y, w, h$  means x location , y location, the width and the height of the target, the probabilistic of the cell contains an object respectively. YOLO shares some similarities with R-CNN. Each grid cell proposes a potential bounding boxes and scores those boxes using convolutional features.

[19],the author propose an algorithm which can runs real time object detect and classification. A key feature of the SSD algorithm is that multi-scale of the convolutional bounding boxes outputs are attached to different feature maps.SSD is faster than R-CNN and its variants.

Real-Time Detectors	Train	mAP	FPS
100Hz DPM	2007	16.0	100
30Hz DPM	2007	26.1	30
Fast YOLO	2007+2012	52.7	<b>155</b>
YOLO	2007+2012	<b>63.4</b>	45
<hr/>			
Less Than Real-Time			
Fastest DPM	2007	30.4	15
R-CNN Minus R	2007	53.5	6
Fast R-CNN	2007+2012	70.0	0.5
Faster R-CNN VGG-16	2007+2012	73.2	7
Faster R-CNN ZF	2007+2012	62.1	18

Figure 1. Real-Time Systems on P ASCAL VOC 2007.Comparing the performance and speed of fast detectors. Fast YOLO is the fastest detector on record for P ASCAL VOC detection and is still twice as accurate as any other real-time detector. YOLO is 10 mAP more accurate than the fast version while still well above real-time in speed.

Figure1 shows the most famous network on image classification. The accuracy and the speed is shown on the table.

### 3. CNN based Medical Image Classification

#### 3.1. Introduction

Convolutional neural networks (CNNs) have been used in the field of computer vision for decades. However, their true value had not been discovered until the ImageNet competition in 2012, a success that brought about a revolution through the efficient use of graphics processing units (GPUs), rectified linear units, new dropout regularization, and effective data augmentation. Acknowledged as one of the top 10 breakthroughs of 2013, CNNs have once again become a popular learning machine, now not only within the computer vision community but across various applications ranging from natural language processing to hyperspectral image processing and to medical image analysis. The main power of a CNN lies in its deep architecture, which allows for extracting a set of discriminating features at multiple levels of abstraction [35].

However, training a deep CNN with full training is very complicated. First, CNNs require a large amount of labeled

Network Index	No. of Convolutional Filters/Size			Connected Layer	Acc
	Layer 1	Layer 2	Layer 3		
CNN-01	48/7x7	72/4x4	512/5x5	512	76%
CNN-02	48/11x11	72/5x5	512/6x6	512	84%
CNN-03	24/11x11	48/5x5	1024/6x6	1024	86%
CNN-04	24/11x11	72/4x4	2048/5x5	2048	80%
CNN-05	48/11x11	72/5x5	1024/6x6	1024	87%

Figure 2. Accuracy results from different CNN configurations

training data. Second, training a deep CNN requires extensive computational and memory resources, without that the training process would be extremely time-consuming.Third, training a deep CNN is often complicated by overfitting. Therefore, deep learning from scratch can be tedious and time-consuming, demanding a great deal of diligence, patience, and expertise.

In this survey, I conducted an extensive set of experiments for 4 medical imaging applications: 1) polyp detection in colonoscopy videos [27] [39], 2) image quality assessment and classification in tissues and cells such as blood vessels videos [20] [17] [11] [13] [9] [3], 3) lung disease such as pulmonary embolism detection and so on in computed tomography (CT) images [30] [16],4) dental disease in X-ray image [22] and 5) intima-media boundary segmentation in ultrasonographic images [28].

#### 3.2. Polyp Detection

Colorectal cancer (CRC) is one of the leading causes of deathworldwide with about estimated 700 thousand deaths in 2012 [39]. Long-term follow-up studies confirmed that removal of adenomatous polyps reduces CRC mortality. Colonoscopy is the preferred technique for colon cancer screening and prevention. The goal of colonoscopy is to find and remove colonic polypsprecursors to colon cancer. But polyps can appear with substantial variations in color, shape, and size. The challenging appearance of polyps can often lead to misdetection [27]. Polyp miss-rates are estimated to be about 4% to 12%; however, a more recent clinical study is suggestive that this misdetection rate may be as high as 25%. So nowadays, there are many research groups start to use computer aided method such as CNN.

In the article [27], the author have a small dataset, which only have 100 images(75 abnormal images and 25 healthy images). After finishing the data augmentation which results in 800 images, they resized the 256\*256 image to 128\*128. In order to test the five architecture they established, he used cross validation method(56 for training and 6 for testing), the result can be seen in Figure 1, the accuracy is just 75% to 80%.

In order to improve the accuracy, in the evaluation phase, the author obtained the final decision for a 256\*256 pixel image by majority voting of the decisions of all 128\*128 pixel subimages(patches). This is a kind of fine-tuning. The redundancy of overlapping subimages can increase the

Stride	No. of Subimages	Accuracy
1	16384	90.22%
5	676	90.22%
20	49	90.21%
32	25	90.96%
48	9	89.27%
Random	16	90.31%
Random	32	90.65%
Random	64	90.49%

Figure 3. Accuracy of different strides for overlapping subimages in the evaluation.

system accuracy likewise to give the assurance of certainty for the overall decision. The result can be seen in figure 2. They also perform a random patch extraction and it can be concluded that there is not much difference between 16384 subimages or just 32 subimages (accuracy of 90.96%), saving considerable computation time and achieving good results.

In the second article [39], the author use a small datasets(PHW Database), this dataset consisted of 1104, 263 and 563 images without polyps, with hyperplasia polyps and adenomatous polyps, respectively, taken under either WL or NBI endoscopy. For fair comparison, 50 images from each class (nonpolyp, hyperplasia, and adenoma) were randomly selected as testing dataset, while the rest were treated as training dataset.

Because this dataset has an imbalanced number of images for each class, Previous study for polyp detection proposed to use an up/down sampling strategy to tackle such challenge. In this paper, the authors randomly down sampled the majority class to match the sample size of the minority class for both target tasks. The source dataset used ILSVRC and Places205 and trained for 450 000 iterations. They tested two tasks using this database, first is polyp detection and second polyp type classification. In order to do the evaluation, the authors used a feature engineering technique: bag-of-words for comparison. After finishing these tasks, we can see the results in the Figure 3.4.5. In these figures, we can see that transferring low-level CNN features gives better transfer learning performance for both target tasks and when a CNN structure is directly used for detection and classification. The performance of the proposed method is better in both tasks.

### 3.3. Tissue Detection and Classification

In this part, I will choose two typical article to discuss. In the first article [13], the author used feature vectors from several pre-trained structures, including networks with/without transfer learning to evaluate the performance of pre-trained deep features versus CNNs which have been trained by that specific dataset as well as the impact of trans-

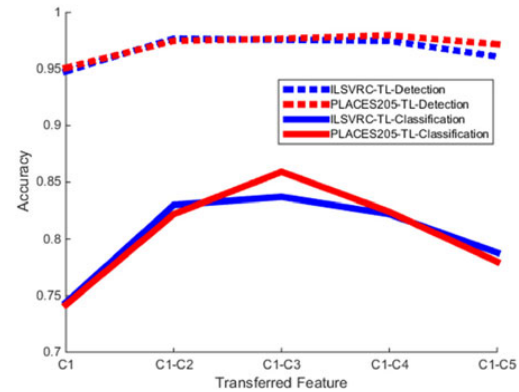


Figure 4. Average accuracy of the detection and classification tasks by transferring C1Cn features learned from ILSVRC and Places205 and using SVM as the classifier with a RBF kernel.

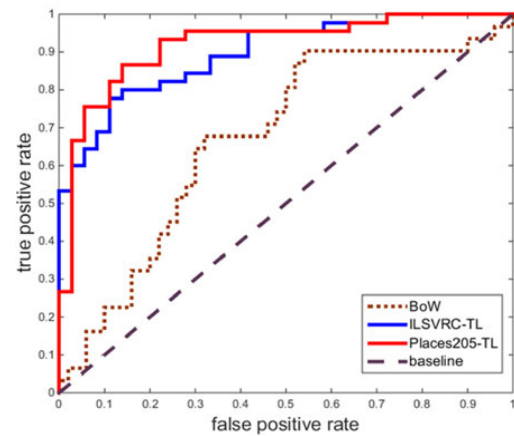


Figure 5. Typical ROC curve for polyp classification for PWH database.

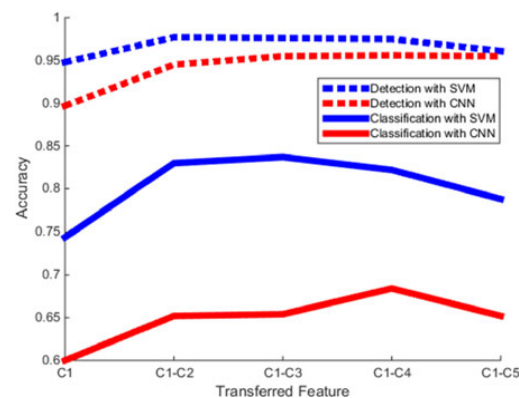


Figure 6. Average accuracy of the detection and classification tasks by transferring C1Cn features learned from ILSVRC and using either RBF kernel SVM or a fully connected CNN layer with a softmax classifier

Scheme	Approach	$\eta_p$	$\eta_w$	$\eta_{total}$
Train from scratch	CNN <sub>1</sub> [17]	64.98%	64.75%	41.80%
Pre-trained features	FE-VGG16	65.21%	64.96%	42.36%
Fine-tuning the pre-trained net	TL-VGG16	63.85%	66.23%	42.29%
Pre-trained features	FE-Inception-v3	70.94%	71.24%	50.54%
Fine-tuning the pre-trained net	TL-Inception-v3	<b>74.87%</b>	<b>76.10%</b>	<b>56.98%</b>

Figure 7. Comparing the results training from scratch, using deep features via a pre-trained network with no change (FE-VGG16), and classification after fine-tuning a pre-trained network (TL-VGG16, TL-Inception-v3). The best scores are highlighted in bold. ( $\eta_p$  means the patch-to-scan accuracy and  $\eta_n$  means whole-scan accuracy)

fer learning with a small number of samples. This experiment is done on Kimia Path24 dataset which consists of 27,055 histopathology training patches in 24 tissue texture classes along with 1,325 test patches for evaluation. In order to do this experiment, the author used fine-tuning method and a pre-trained CNN as a feature extractor and a fine-tuned CNN as a classifier.

The result shows in figure 6 that pre-trained networks are quite competitive against training from scratch. In this figure, VGG16 and CNN are quite similar, whereas the results for Inception-v3 are similar with the transfer-learned model outperforming the feature extractor. But considering Inception-v3 requires no extra effort and produces similar results with a linear SVM, one may prefer using it to train from scratch and fine-tuning a pre-trained net.

In the next article [3], the authors designed a specific CNN network which perform image-wise classification in four classes of medical relevance: normal tissue, benign lesion, in situ carcinoma and invasive carcinoma. The proposed CNN architecture is designed to integrate information from multiple histological scales, including nuclei, nuclei organization and overall structure organization. A data augmentation method is adopted to increase the number of cases in this training set. A SVM classification using the features extracted by the CNN is also used for comparison purposes.

The dataset is composed of an extended training set of 249 images, and a separate test set of 20 images. In these datasets, the four classes are balanced. The images were selected so that the pathology classification can be objectively determined from the image contents. An additional test set of 16 images is provided with images of increased ambiguity, which they denote as extended dataset.

They first normalized the images. First, the colors of the images are converted to optical density (OD) using a logarithmic transformation. Then, they used singular value decomposition (SVD) to the OD tuples to find the 2D projections with higher variance. The resulting color space transform is then applied to the original image. Finally, the image histogram is stretched so that the dynamic range covers the lower 90% of the data.

Then they do two kinds of classification: Image-wise

Dataset	Classifier	non-carcinoma		carcinoma	
		Normal	Benign	<i>in situ</i>	Invasive
Initial	CNN	61.7	69.2	83.3	91.7
	CNN+SVM	61.7	76.7	83.3	88.3
Extended	CNN	65.0	61.7	76.7	88.3
	CNN+SVM	50	72.9	58.3	66.7
Overall	CNN	54.2	66.7	43.8	56.3
	CNN+SVM	56.4	63.9	72.2	74.1
		60.2	63.9	62.0	74.1

Figure 8. Patch-wise sensitivity (%) (2 and 4 classes).

Classif.	Vote	4 Classes			2 Classes		
		Init.	Exten.	Overall	Init.	Exten.	Overall
CNN	Major.	80.0	75.0	77.8	80.0	81.3	80.6
	Max.	80.0	62.5	72.2	80.0	75.0	77.8
	Sum	80.0	68.8	75.0	80.0	75.0	77.8
	Sum	85.0	68.8	77.8	90.0	75.0	83.3
CNN+SVM	Major.	80.0	62.5	72.2	80.0	75.0	77.8
	Max.	80.0	62.5	72.2	80.0	75.0	77.8
	Sum	85.0	68.8	77.8	90.0	75.0	83.3
	Sum	85.0	68.8	77.8	90.0	75.0	83.3

Figure 9. Image-wise accuracy (%) using different voting rules (2 and 4 classes).

Dataset	Classifier	non-carcinoma		carcinoma	
		Normal	Benign	<i>in situ</i>	Invasive
Initial	CNN	70	40	100	100
	CNN+SVM	80	80	100	100
Extended	CNN	80	60	100	100
	CNN+SVM	75	50	75	75
Overall	CNN	77.8	61.1	88.9	88.9
	CNN+SVM	77.8	66.7	77.8	88.9

Figure 10. Image-wise sensitivity (%) using majority voting (2 and 4 classes).

classification and CNN patch-wise classification. Image-wise classification first divided the origin image into twelve contiguous non-overlapping patches and then use one of three different patch patch methods: majority voting, maximum probability and sum of probabilities. CNN patch-wise classification used 75% of the data to do the training and validated on the remaining images. The validation set is randomly selected for each epoch. The training process stops after the stabilization of the validation accuracy with equal weight for all the classes (50 epochs). The authors also used the features extracted by the CNN to train a SVM classifier to do the comparison. The result can be seen in figure 7, 8 and 9.

In figure 7, we can see the result similar between the CNN and CNN+SVM. But the performance of this network is lower for the extended dataset due to its increased complexity. In figure 8 and 9, we can see that CNN+SVM get the best result with the majority voting method. In comparison, CNN's performance is only better for the extended set using majority voting. In addition, we can see that maximum probability is the worst performing method in both methods, which means that this method is not suit in this case.

### 3.4. Some kinds of lung Diseases Classification

Lung cancer is notoriously aggressive with a low long-term survival rate. Quantitative analysis in lung nodules using thoracic Computed Tomography(CT) has been a central focus for early cancer diagnosis, where CT phenotype pro-



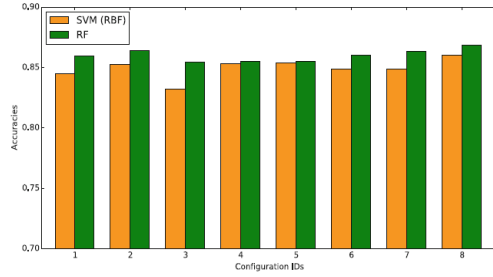


Figure 11. The classification performance of SVM with the RBF kernel and RF based on features from the MCNN using 8 different configurations. Each configuration is assigned to a unique ID for display convenience

Classifier	Scales	HOG			LBP		
		$s_w = 8$	$s_w = 16$	$s_w = 32$	$n_{pt} = 8$	$n_{pt} = 16$	$n_{pt} = 24$
SVM	32	74.18 %	63.27 %	49.82 %	64.58 %	66.40 %	67.35 %
	64	66.69 %	66.40 %	56.15 %	49.24 %	59.93 %	59.20 %
	96	64.07 %	65.16 %	56.58 %	36.00 %	52.22 %	54.84 %
RF	32	75.93 %	67.71 %	60.07 %	71.27 %	72.07 %	73.67 %
	64	73.16 %	67.78 %	62.84 %	62.54 %	62.25 %	66.55 %
	96	67.56 %	64.58 %	61.75 %	60.07 %	60.15 %	62.84 %

Figure 12. Performance using the HOG and LBP descriptors with different  $S_w$  and  $n_{pt}$

vides a powerful tool to comprehensively capture nodule characteristics. The importance of diagnostically classifying malignant and benign nodules using CT images is to facilitate radiologists for nodule staging assessment and individual therapeutic planning. [30]

In the first article [30], the authors used the LIDC-IDRI datasets, which has 1375 nodule pictures(1100 for training and 275 for testing). In order to improve the speed and accuracy, the authors introduced an Multi-scale Convolutional Neural Networks(MCNN) model to do the lung nodule diagnostic classification. This CNN model take multi-scale raw nodule patches and remove the need of any hand-crafted feature engineering work. This network can also deal with noisy data in nodule CT.

Because of the clinical fact that nodule sizes vary remarkably, this network take patches from different scales(3 layers) as inputs in parallel. The parameter is shared between these layers to reduce parameter. When doing the evaluation task, the result is decided by all the layers. The authors use the HOG and uniform LBP descriptor and SVM and RF classifier to do the classification. The result can be seen in figure 6 and figure 7. In figure 7, the  $S_w$  means the size of the cell window for SVM and  $n_{pt}$  means the number of neighbourhood points for LBP.

The second article [16] is about using CNN to classify the ILD patterns. This experiment used an ILD database which contains 113 sets of HRCT images, with 2062 2D regions indicting the ILD category. In order to augment the dataset, the CT slices were divided into half-overlapping

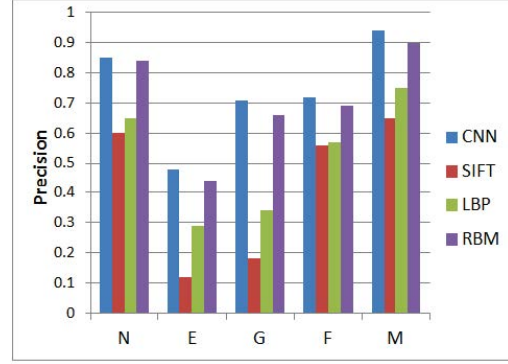


Figure 13. The classification results comparing proposed customized CNN method with SIFT, LBP and RBM

image and the only if 75% percent of its pixels falling inside the regions of interest will be adopted. The dataset thus contains 16220 image patches from 92 HRCT image sets, including 4348 norm patches, 1047 emphysema patches, 1953 ground glass patches, 2591 fibrosis patches, and 6281 micronodules patches.

The authors compared their classification results with three other feature extraction approaches: SIFT feature, LBP feature and unsupervised feature learning using RBM. The result can be seen in Figure 8. In this figure, we can see that their customized CNN method achieved the best classification performance.

### 3.5. Dental Disease Classification in X-ray Images

The author found that there is no literature for dental disease classification, so the research group start to use CNN to deal with X-ray images and make some breakthrough. Orthopantomogram (OPG) and Radiovisiography (RVG) x-ray images are the most widely used tools for the diagnoses of dental diseases. Dental caries is one of the most common dental disease worldwide and it has different stages. So the CNN network in this experiment is used to classify mainly 3 classes (dental caries, periapical infection, periodontitis) [22].

Because though the radiologists have large dataset of dental x-ray images, these x-ray images have individual privacy issues. So the dataset is very small in this experiment, just have 251 grey images of dimension 1000\*1496. So the authors use transfer learning method to do the fine tuning and improve the accuracy very much. They changed some unfrozen layers used for training in order for the pre-trained model to be more adaptive to the training data.

They first resize these picture to 500\*748, and then use 180 of 251 to do the training, 45 images for validation and 26 images for testing purpose. Because of the unavailability of the large dataset, CNN architecture could not perform well in this classification task. After they used trans-

Model	Accuracy
CNN	0.7307
Transfer learning	0.8846
Transfer learning with fine tuning	0.8846

Figure 14. The comparison of different models

Disease Name	Number of Samples	Correct results	Accuracy
Dental Caries	8	7	0.875
Periapical Infection	10	9	0.90
Periodontitis	8	7	0.875
Total	26	23	0.8846

Figure 15. Experimental results for transfer learning model

fer learning model to do the fine tuning, the accuracy is increased by 15.39% compared to pure CNN model, and achieved 88.46% accuracy, which is very encouraging.

### 3.6. Intima-media Boundary Segmentation

Automated classification of human anatomy is an important prerequisite for many computer-aided diagnosis systems. The spatial complexity and variability of anatomy throughout the human body makes classification difficult. So the authors want to use CNN to do this classification. In this paper, the authors choose to use 4298 separate axial 2D key images to learn 5 anatomical classes(neck, lungs, liver, pelvis and legs) [28].

When applying the CNN to build the anatomy-specific classifier for CT images, because the authors want to classify these picture to 5 classes, so they choose 5 cascaded layers. All the convolutional filter kernel elements are trained from the data in a supervised fashion. In order to avoid overfitting, the fully-connected layers are constrained, using the *DropOut* method. The datasets are from the Picture Archiving and Communication System (PACS) of the Clinical Center of the National Institutes of Health. In order to enrich their data, they use spatial deformations to each image, using random translation, rotations and non-rigid deformations, which lead their datasets from hundred's picture to near 100 thousand pictures. Before import into the CNN, the author resize all the picture to 256\*256 pixels. The authors use 80% of their total dataset to train the CNN and reserve 20% to do the test. After doing the experiments, the accuracy of this net can reach 94.1%, which can be seen in figure 4. This classification result is achieved in less than 1 minute on a modern desktop computer and GPU card (Dell Precision T7500, 24GB RAM, NVIDIA Titan Z).

		prediction				
		legs	pelvis	liver	lung	neck
actual	legs	90	0	0	0	0
	pelvis	0	24	2	0	1
	liver	0	6	484	42	0
	lungs	0	0	28	93	5
	neck	0	0	0	0	102
error		9.6%				

		prediction				
		legs	pelvis	liver	lung	neck
actual	legs	90	0	0	0	0
	pelvis	0	27	0	0	0
	liver	0	0	518	14	0
	lungs	0	0	38	88	0
	neck	0	0	0	0	102
error		5.9%				

Figure 16. Confusion matrices on the original test images before and after data augmentation.

### 3.7. Conclusion

In this part, I aimed to address to know how the CNN can be used on the medical image classification and the result these experiments made. My experiment, based on 4 distinct medical imaging applications from different imaging modality systems, have demonstrated that deep CNN are useful for medical image analysis. If the training data is limited, the fine-tuned CNN can perform better than fully trained CNN. I think the potential of CNNs for medical imaging applications is confirmed because both deeply fine-tuned CNNs and fully trained CNNs can outperform the corresponding handcrafted alternatives. We can also see that the speed is depend on the devices, the more powerful the graphics is, the quicker the CNN network use to train.

### 4. Conclusion

In this paper we have discussed about the different types of image classification techniques and many CNN based medical application. So this paper will help us in selecting an appropriate classification technique among all the available techniques.

## References

- [1] B. Alexe, T. Deselaers, and V. Ferrari. Measuring the objectness of image windows. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(11):2189–2202, 2012.
- [2] N. Aloysius and M. Geetha. A review on deep convolutional neural networks. In *Communication and Signal Processing (ICCSP), 2017 International Conference on*, pages 0588–0592. IEEE, 2017.
- [3] T. Araújo, G. Aresta, E. Castro, J. Rouco, P. Aguiar, C. Eloy, A. Polónia, and A. Campilho. Classification of breast cancer histology images using convolutional neural networks. *PLoS one*, 12(6):e0177544, 2017.
- [4] D. Ciregan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. *computer vision and pattern recognition*, pages 3642–3649, 2012.
- [5] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655, 2014.
- [6] P. F. Felzenszwalb, R. B. Girshick, D. A. Mcallester, and D. Ramanan. Visual object detection with deformable part models. *Communications of The ACM*, 56(9):97–105, 2013.
- [7] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, June 2014.
- [8] R. B. Girshick. Fast r-cnn. *international conference on computer vision*, pages 1440–1448, 2015.
- [9] O. Hadad, R. Bakalo, R. Ben-Ari, S. Hashoul, and G. Amit. Classification of breast lesions using cross-modal deep learning. In *Biomedical Imaging (ISBI 2017), 2017 IEEE 14th International Symposium on*, pages 109–112. IEEE, 2017.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [11] Y. Huang, H. Zheng, C. Liu, X. Ding, and G. K. Rohde. Epithelium-stroma classification via convolutional neural networks and unsupervised domain adaptation in histopathological images. *IEEE journal of biomedical and health informatics*, 21(6):1625–1632, 2017.
- [12] P. Kamavisdar, S. Saluja, and S. Agrawal. A survey on image classification approaches and techniques. *International Journal of Advanced Research in Computer and Communication Engineering*, 2(1):1005–1009, 2013.
- [13] B. Kieffer, M. Babaie, S. Kalra, and H. Tizhoosh. Convolutional neural networks for histopathology image classification: Training vs. using pre-trained networks. *arXiv preprint arXiv:1710.05726*, 2017.
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [15] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel. Handwritten digit recognition with a back-propagation network. In *Advances in neural information processing systems*, pages 396–404, 1990.
- [16] Q. Li, W. Cai, X. Wang, Y. Zhou, D. D. Feng, and M. Chen. Medical image classification with convolutional neural network. In *Control Automation Robotics & Vision (ICARCV), 2014 13th International Conference on*, pages 844–848. IEEE, 2014.
- [17] X. Li, W. Li, X. Xu, and W. Hu. Cell classification using convolutional neural networks in medical hyperspectral imagery. In *Image, Vision and Computing (ICIVC), 2017 2nd International Conference on*, pages 501–504. IEEE, 2017.
- [18] M. Lin, Q. Chen, and S. Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.
- [19] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg. Ssd: Single shot multibox detector. *European conference on computer vision*, pages 21–37, 2016.
- [20] S. McIlroy, Y. Kubo, T. Trappenberg, J. Toguri, and C. Lehmann. In vivo classification of inflammation in blood vessels with convolutional neural networks. In *Neural Networks (IJCNN), 2017 International Joint Conference on*, pages 3022–3027. IEEE, 2017.
- [21] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [22] S. A. Prajapati, R. Nagaraj, and S. Mitra. Classification of dental diseases using cnn and transfer learning.
- [23] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. *computer vision and pattern recognition*, pages 779–788, 2016.
- [24] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. pages 6517–6525, 2016.
- [25] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv*, 2018.
- [26] S. Ren, K. He, R. B. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017.
- [27] E. Ribeiro, A. Uhl, and M. Häfner. Colonic polyp classification with convolutional neural networks. In *Computer-Based Medical Systems (CBMS), 2016 IEEE 29th International Symposium on*, pages 253–258. IEEE, 2016.
- [28] H. R. Roth, C. T. Lee, H.-C. Shin, A. Seff, L. Kim, J. Yao, L. Lu, and R. M. Summers. Anatomy-specific classification of medical images using deep convolutional nets. In *Biomedical Imaging (ISBI), 2015 IEEE 12th International Symposium on*, pages 101–104. IEEE, 2015.
- [29] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- [30] W. Shen, M. Zhou, F. Yang, C. Yang, and J. Tian. Multi-scale convolutional neural networks for lung nodule classification. In *International Conference on Information Processing in Medical Imaging*, pages 588–599. Springer, 2015.



- [31] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [32] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387*, 2015.
- [33] R. K. Srivastava, K. Greff, and J. Schmidhuber. Training very deep networks. In *Advances in neural information processing systems*, pages 2377–2385, 2015.
- [34] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, et al. Going deeper with convolutions. *Cvpr*, 2015.
- [35] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, and J. Liang. Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE transactions on medical imaging*, 35(5):1299–1312, 2016.
- [36] J. R. R. Uijlings, K. E. A. V. De Sande, T. Gevers, and A. W. M. Smeulders. Selective search for object recognition. *International Journal of Computer Vision*, 104(2):154–171, 2013.
- [37] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [38] M. D. Zeiler, G. W. Taylor, and R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2018–2025. IEEE, 2011.
- [39] R. Zhang, Y. Zheng, T. W. C. Mak, R. Yu, S. H. Wong, J. Y. Lau, and C. C. Poon. Automatic detection and classification of colorectal polyps by transferring low-level cnn features from nonmedical domain. *IEEE journal of biomedical and health informatics*, 21(1):41–47, 2017.