

# Luo's part about the computer vision survey

Ruiqi Luo

16064455

Shanghaitech University

Huanke Rd No.199, Zhangjiang, Pudong district, Shanghai, China

luorq@shanghaitech.edu.cn

## 1. Part III: CNN in medical image classification

### 1.1. 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[12].

However, training a deep CNN with full training is very complicated. First, CNNs require a large amount of labeled 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 [9][13], 2) image quality assessment and classification in tissues and cells such as blood vessels videos[7][6][3][4][2][1], 3) lung disease such as pulmonary embolism detection and so on in computed tomography (CT) images[11][5], 4) dental disease in X-ray image[8] and 5) intima-media boundary segmentation in ultrasonographic images[10].

| 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 1. Accuracy results from different CNN configurations

### 1.2. Polyp detection

Colorectal cancer (CRC) is one of the leading causes of death worldwide with about estimated 700 thousand deaths in 2012[13]. 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 polyps precursors 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[9]. 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[9], 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 system accuracy likewise to give the assurance of certainty for the overall decision. The result can be seen in figure 2.

| 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 2. Accuracy of different strides for overlapping subimages in the evaluation.

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[13], the author use a small dataset-(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.

### 1.3. tissue detection and classification

In this part, I will choose two typical article to discuss. In the first article[4], 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 transfer learning with a small number of samples. This experiment is done on Kimia Path24 dataset which consists of 27,055 histopathol-

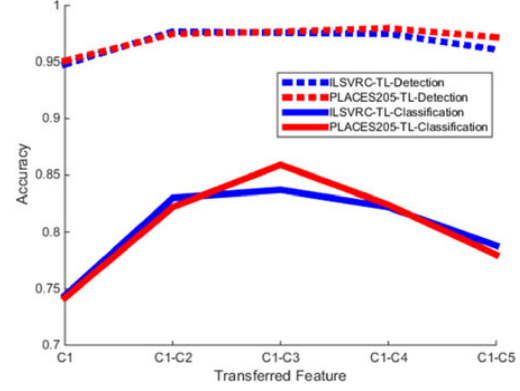


Figure 3. Average accuracy of the detection and classification tasks by transferring C1Cn features learned from ILSVRC and Places205 and using SVM as the classifier with a RBF kernel.

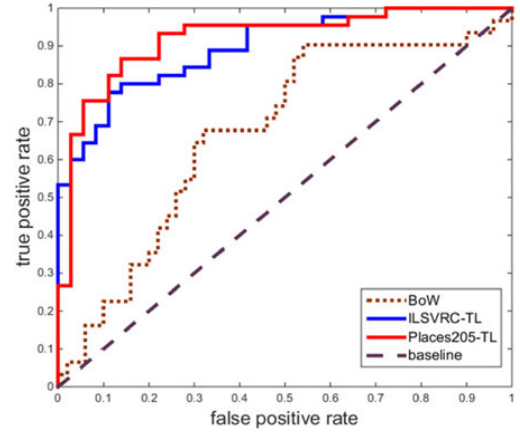


Figure 4. Typical ROC curve for polyp classification for PWH database.

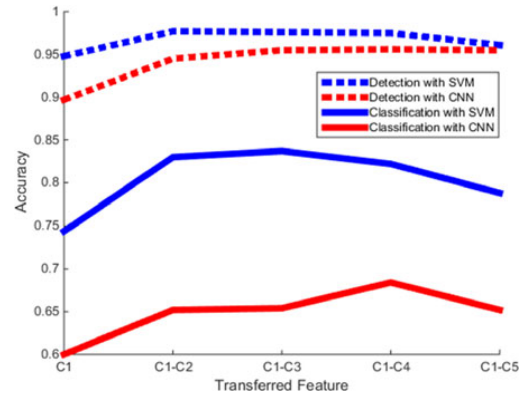


Figure 5. Average accuracy of the detection and classification tasks by transferring C1Cn features learned from ILSVRC and using either RBF kernel SVM or a fully connected CNN layer with a softmax classifier

| Scheme                          | Approach              | $\eta_p$      | $\eta_n$      | $\eta_{total}$ |
|---------------------------------|-----------------------|---------------|---------------|----------------|
| 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 6. 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)

ogy 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 training from scratch and fine-tuning a pre-trained net.

In the next article[1], 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 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

| 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          | 81.3   | 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 7. Patch-wise sensitivity (%) (2 and 4 classes).

| Classif. | Vote | 4 Classes |        |         | 2 Classes |        |         |
|----------|------|-----------|--------|---------|-----------|--------|---------|
|          |      | Init.     | Exten. | Overall | Init.     | Exten. | Overall |
| CNN      | Max. | 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  | Max. | 80.0      | 62.5   | 72.2    | 80.0      | 75.0   | 77.8    |
|          | Max. | 80.0      | 68.8   | 77.8    | 90.0      | 75.0   | 83.3    |
|          | 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 8. 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 9. Image-wise sensitivity (%) using majority voting (2 and 4 classes).

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.

#### 1.4. 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 provides a powerful tool to comprehensively capture nodule characteristics. The importance of diagnostically classifying malignant and benign nodules using CT images is to

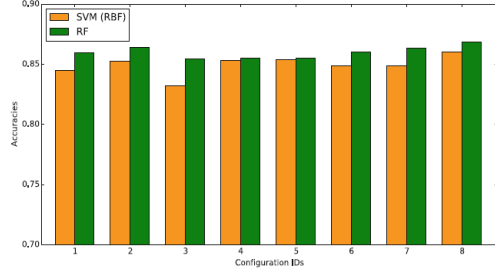


Figure 10. The classification performance of SVM with the RBF kernel and RF based on features from the MCNN using 8 different configurations. Each configuration is assigned to a unique ID for display convenience

| 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 11. Performance using the HOG and LBP descriptors with different  $S_w$  and  $n_{pt}$

facilitate radiologists for nodule staging assessment and individual therapeutic planning.[11]

In the first article[11], 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[5] 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 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

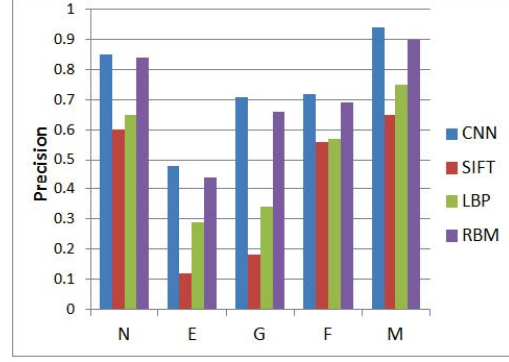


Figure 12. The classification results comparing proposed customized CNN method with SIFT, LBP and RBM

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.

## 1.5. 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)[8].

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 transfer 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.



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

Figure 13. 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 14. Experimental results for transfer learning model

## 1.6. 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)[10].

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).

## 1.7. Conclusion

In this part, I aimed to address to know how the CNN can be used on the medical image classification and the re-

|        |        | 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 15. Confusion matrices on the original test images before and after data augmentation.

sult 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.

## References

- [1] 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.
- [2] 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.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] S. A. Prajapati, R. Nagaraj, and S. Mitra. Classification of dental diseases using cnn and transfer learning.

- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.