

# Detection of DCIS and IDC in Whole-Slide H&E Stained Breast Histopathology Images Using Stacked Convolutional Neural Networks

Guido C. A. Zuidhof\*, Babak Ehteshami Bejnordi, Geert Litjens, and Jason Farquhar

**Abstract**—This paper presents and evaluates a method for detection and localization of breast malignant lesions in histopathological breast tissue images. The goal of this method is to best classify digitized whole-slide hematoxylin and eosin (H&E) stained whole-slide images (WSIs) into three classes: benign, ductal carcinoma in situ (DCIS) and invasive ductal carcinoma (IDC). A convolutional neural network (CNN) was trained on small patches of these images, labeled with the class of the center pixel. To distinguish between the DCIS and IDC class the information available in these small patches is often not sufficient. A second CNN was trained on the output of the aforementioned network with a larger patch size to capture a larger context. Sliding this classifier across a whole slide image yields a probability map. From this probability map structural and statistical features were extracted, which were used to train a final classifier to predict the label for the whole-slide image. The system is evaluated on a dataset containing X WSIs of breast tissue using FROC analysis. The results show a...

**Index Terms**—Computer-aided diagnosis, DCIS and IDC detection, deep learning, whole-slide imaging.

## I. INTRODUCTION

**B**REAST cancer is the most common cancer in women [1]. Cancers of the breast kill more women than any other form of cancer in all parts of the developing world [2]. An important tool for the detection and management of breast cancer is analysis of tissue samples under the microscope by a pathologist.

Looking at cancer cells under the microscope, the pathologist searches for certain features that can help predict how likely the cancer is to grow and spread. These features include the spatial arrangement of the cells, morphometric characteristics of the nuclei, whether they form tubules, and how many of the cancer cells are in the process of dividing (mitotic count). These features taken together determine the extent or spread of cancer at the time of diagnosis.

This paper was written as part of the master's thesis in Artificial Intelligence of G. C. A. Zuidhof. This project was supervised by B. Ehteshami Bejnordi, G. Litjens and J. Farquhar. *Asterisk indicates corresponding author.*

\*G. C. A. Zuidhof is a student in the Artificial Intelligence Master's Programme at the Radboud University, 6525HP Nijmegen, The Netherlands (e-mail: guido.zuidhof@student.ru.nl).

B. Ehteshami Bejnordi is with the Diagnostic Image Analysis Group, Radboud University Medical Center, 6500HB Nijmegen, The Netherlands.

G. Litjens is with the Department of Pathology, Radboud University Medical Center, 6500HB Nijmegen, The Netherlands.

J. Farquhar is with the Department of Artificial Intelligence, Radboud University, 6525HP Nijmegen, The Netherlands and the Donders Institute for Brain, Cognition and Behaviour, 6525EN Nijmegen, The Netherlands.

First draft, work in progress; October 1, 2016.

Increasingly, slides of human tissue are digitized instead of being analysed only under a microscope. This has spawned the relatively new field of *digital pathology*. Digital images as opposed to glass slides allow for easier consultations between pathology experts. An additional advantage of digital images is the opportunity to analyze these images automatically using computer algorithms.

Visual microscopic interpretation of tissue sections is laborious and prone to subjectivity. *Computer-aided diagnosis (CAD)* has a huge potential in alleviating shortcomings of human interpretation and will reduce the workload of the pathologists. As a result, more accurate diagnostic information may be extracted, helping clinicians in selecting the most optimal treatment for individual patients. CAD can facilitate diagnosis by sieving out obviously benign slides and providing quantitative characterization of suspicious areas.

### A. Patch-based classification

Convolutional neural networks, as well as other statistical machine learning methods, are bound by computational and memory constraints. The tissue slides are scanned at a high magnification (up to 40X), resulting in very large images. A typical image can be 100,000 by 200,000 pixels, and have a file size of around 20GB uncompressed.

To create a classification for the whole-slide image a patch based method will be employed. A patch is a small sub-image of the whole image. By applying the classifier in a sliding window approach over the original image we can create a prediction map for the whole image.

TODO more about domain, dataset

## II. DENSE PREDICTION

TODO a lot

### A. Architectures

We will apply two different architectures to the problem of patch classification, with patches of size 224x224 pixels.

1) *Wide Residual Networks*: Residual network was the winner of the 2015 ImageNet challenge [3]. Other than previous networks that entered this competition, it features no fully connected layer at the end of the network. It was the deepest successful architecture to date featuring up to 152 layers. It can get away with this depth by using residual learning blocks, where the layers learn residual functions with relation to

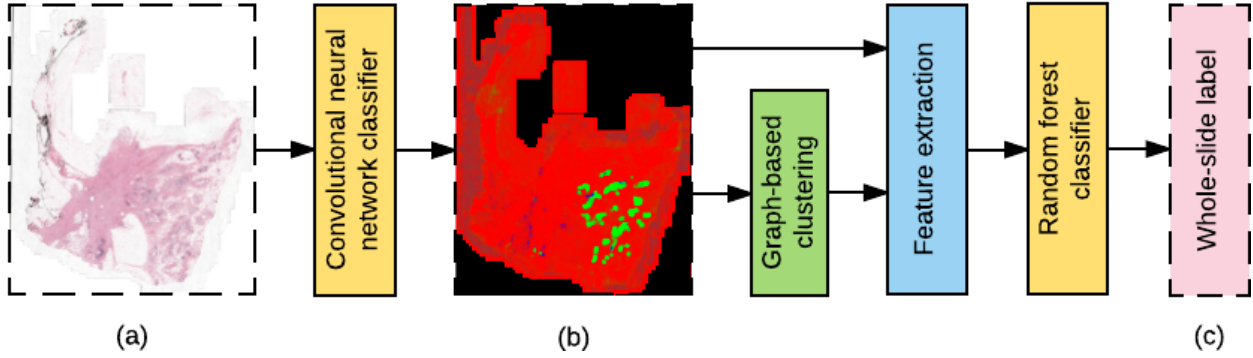


Fig. 1. Overview of the proposed system for labeling whole-slide images. (a) Original WSI of breast tissue. (b) Resulting probability map of applying the CNN in a patch-based fashion. (c) Output of the system; a single label for the whole-slide image (either benign, DCIS or IDC).

the input instead of learning unreferenced functions. To the extreme, it makes it easier to learn the identity function if that is optimal for that layer.

In this study we will apply an adaptation of this architecture, the wide residual network as proposed by Zagoruyko et al. [4]. They showed that wider and less deep residual networks outperform their deeper and thinner counterparts both in accuracy and efficiency.

The architecture is a recipe that has three main design choices, the  $N$  value that determines the depth, the  $k$  value that determines the width (the amount of filters) and the type of ResNet block. Here, we use  $N = 4$ ,  $k = 2$ , which was empirically evaluated to fit a decent batch size and train fast. The ResNet block type is detailed in figure 2b, TODO more about this block. The resulting architecture is shown in figure 2a.

2) *VGG-16*: VGG-16 is a popular architecture because of its simplicity. It features small 3x3 convolution filters and 2x2 pooling throughout the network and has a depth of 16 layers. The architecture is detailed in figure 2c. Unlike the ResNet architecture it does not feature skip connections, also it features a stack of two dense fully connected layers at the end.

### B. Preprocessing

As a preprocessing step, the images are normalized by dividing their pixel value by 255. Then, the mean value of each color channel is subtracted to zero center the data.

### C. Learning rate decay

The learning rate reduction policy is as follows. The learning rate is multiplied by 0.2 after no better validation accuracy has been observed for 8 epochs, this value is called the *patience*. This patience is increased by 20% after every reduction in learning rate (rounded up).

### D. Data augmentation

Artificially increasing the amount of data by adding variations of the original data can further help train a model that generalizes well. Augmentation consists of applying (random) perturbations to the samples in ways that do not change the

TABLE I  
BEST EPOCH PATCH-LEVEL ACCURACY OF 224x224 NETWORKS

| LABELS                   | ARCHITECTURE | ACCURACY |
|--------------------------|--------------|----------|
| <i>Benign, Cancer</i>    | VGG-16       | 0.9237   |
|                          | WRN-4-2      | 0.9241   |
| <i>Benign, DCIS, IDC</i> | VGG-16       | 0.8245   |
|                          | WRN-4-2      | 0.7995   |

label of the samples. It has a regularizing effect which helps prevent overfitting. Especially effective are augmentations that are also realistic examples of real world data.

The augmentation methods used are as follows:

1) *Flips*: The images are randomly mirrored in the X and/or Y direction, both with a 0.5 probability.

2) *Rotations*: The images are randomly rotated 0, 90, 180, or 270 degrees with equal probability.

3) *HSV jittering*: HSV is a colorspace that represents a color image in three channels; a *hue* (color), *saturation* (vibrance) and a *value* (brightness) channel. There exists a large variability in the staining of the slides, which shows itself mostly in the hue and saturation channel. We randomly jitter the hue and saturation of an image with a random value between -0.075 and 0.075. The values of all pixels are then clipped between 0 and 1.

Another commonly used data augmentation method is also zooming. This was not used here, as the size of the nuclei plays a role in determining the class label. Elastic distortions were successfully applied as part of the data augmentation strategy in other applications of convolutional neural networks [5], also in the medical imaging domain (cytology) [6]. The irregularity of the sizes of the nuclei, as well as the architectural features are important biomarkers that can be used to distinguish between benign and cancerous lesions. Applying these elastic distortions would likely contaminate this information, and as such, it was not used here.

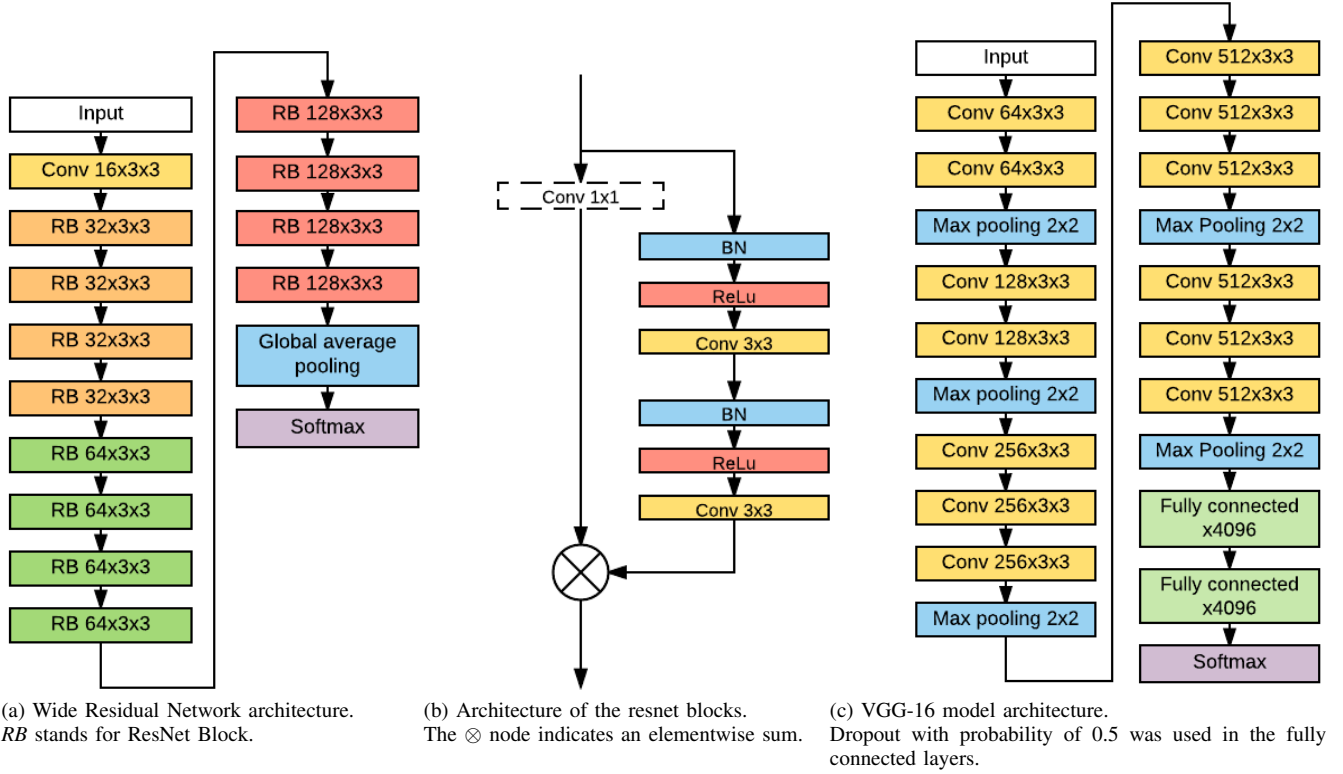


Fig. 2. Architectures used for 224x224 patch classification. The 1x1 convolution layer in the resnet blocks is only present in the blocks where the input is downsampled, which is every first block with a higher amount of filters in the Wide ResNet architecture (figure 2a).

TABLE II  
BEST EPOCH PATCH-LEVEL ACCURACY OF 768x768 STACKED NETWORKS

| LABELS                   | STACKED ON      | ACCURACY |
|--------------------------|-----------------|----------|
| <i>Benign, Cancer</i>    | 2 Class WRN-4-2 | 123      |
|                          | 3 Class WRN-4-2 | 123      |
| <i>Benign, DCIS, IDC</i> | 2 Class WRN-4-2 | 123      |
|                          | 3 Class WRN-4-2 | 123      |

- [4] S. Zagoruyko and N. Komodakis, "Wide residual networks," *CoRR*, vol. abs/1605.07146, 2016. [Online]. Available: <http://arxiv.org/abs/1605.07146>
- [5] J. P. Patrice Y. Simard, Dave Steinkraus, "Best practices for convolutional neural networks applied to visual document analysis." Institute of Electrical and Electronics Engineers, Inc., August 2003.
- [6] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," *CoRR*, vol. abs/1505.04597, 2015. [Online]. Available: <http://arxiv.org/abs/1505.04597>

### III. WHOLE-SLIDE IMAGE LABELING

### IV. RESULTS

### V. CONCLUSION

The conclusion goes here.

### ACKNOWLEDGMENT

The authors would like to thank...

### REFERENCES

- [1] R. L. Siegel, K. D. Miller, and A. Jemal, "Cancer statistics, 2016," *CA: A Cancer Journal for Clinicians*, vol. 66, no. 1, pp. 7–30, 2016. [Online]. Available: <http://dx.doi.org/10.3322/caac.21332>
- [2] P. Porter, "Westernizing women's risks? breast cancer in lower-income countries," *New England Journal of Medicine*, vol. 358, no. 3, pp. 213–216, 2008, pMID: 18199859. [Online]. Available: <http://dx.doi.org/10.1056/NEJMp0708307>
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>

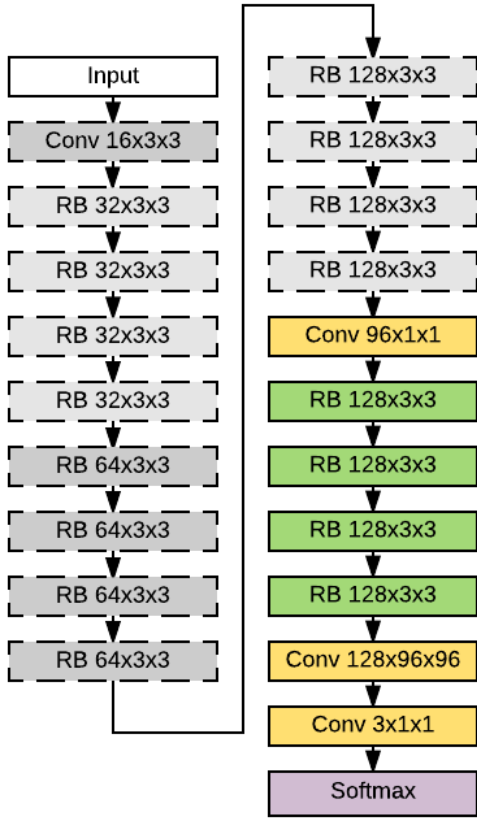


Fig. 3. Architecture of the stacked network. The weights of the components with the dotted outlines are taken from the previously trained 224x224 patch model, and are no longer updated (these are frozen).

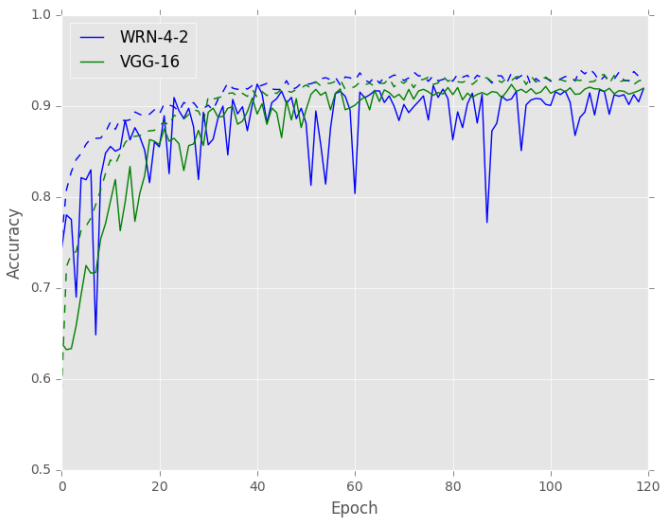


Fig. 4. Performance of networks trained on the two class problem of benign versus cancerous patches. The dashed line shows the performance on the images from the train set it was shown that epoch.

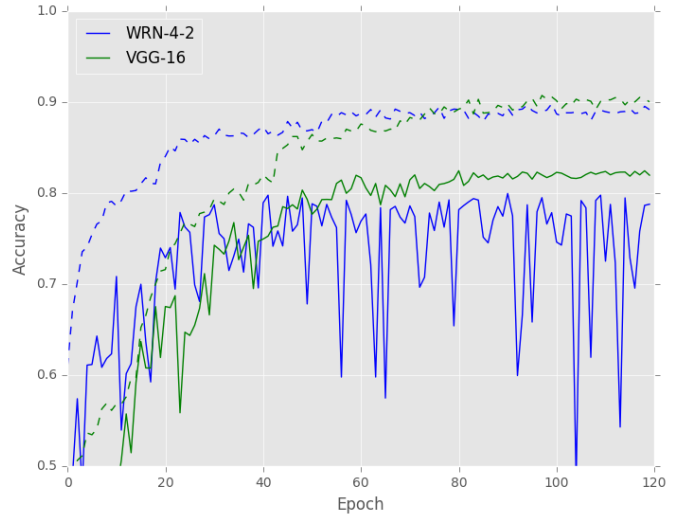


Fig. 5. Performance of networks trained on all three classes. The dashed line shows the performance on the images from the train set it was shown that epoch. The validation accuracy of the wide resnet is especially noisy.