

Sparse Autoencoder for Unsupervised Nucleus Detection and Representation in Histopathology Images

Le Hou¹, Vu Nguyen¹, Dimitris Samaras¹, Tahsin M. Kurc^{1,2}, Yi Gao¹, Tianhao Zhao¹, Joel H. Saltz^{1,3}
¹Stony Brook University, ²Oak Ridge National Laboratory, ³Stony Brook University Hospital

Abstract

Histopathology images are crucial to the study of complex diseases such as cancer. The histologic characteristics of nuclei play a key role in disease diagnosis, prognosis and analysis. In this work, we propose a sparse Convolutional Autoencoder (CAE) for fully unsupervised, simultaneous nucleus detection and feature extraction in histopathology tissue images. Our CAE detects and encodes nuclei in image patches in tissue images into sparse feature maps that encode both the location and appearance of nuclei. Our CAE is the first unsupervised detection network for computer vision applications. The pretrained nucleus detection and feature extraction modules in our CAE can be fine-tuned for supervised learning in an end-to-end fashion. We evaluate our method on four datasets and reduce the errors of state-of-the-art methods up to 42%. We are able to achieve comparable performance with only 5% of the fully-supervised annotation cost.

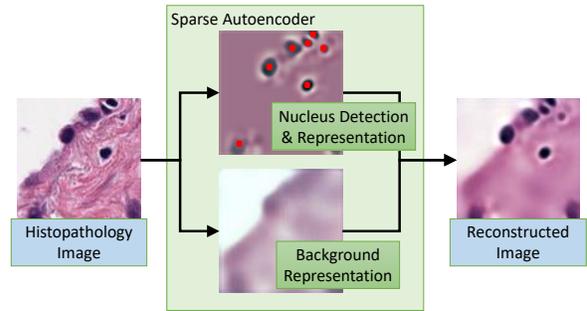


Figure 1: Our fully unsupervised autoencoder. It first decomposes an input histopathology image patch into foreground (e.g. nuclei) and background (e.g. cytoplasm). It then detects nuclei in the foreground by representing the locations of nuclei as a sparse feature map. Finally, it encodes each nucleus to a feature vector. Our autoencoder is trained end-to-end, minimizing the reconstruction error.

1. Introduction

Pathologists routinely examine glass tissue slides for disease diagnosis in healthcare settings. Nuclear characteristics, such as size, shape and chromatin pattern, are important factors in distinguishing different types of cells and diagnosing disease stage. Manual examination of glass slides, however, is not feasible in large scale research studies which may involve thousands of slides. Automatic analysis of nuclei can provide quantitative measures and new insights to disease mechanisms that cannot be gleaned from manual, qualitative evaluations of tissue specimens.

Collecting a large-scale supervised dataset is a labor intensive and challenging process since it requires the involvement of expert pathologists who’s time is a very limited and expensive resource [1]. Thus existing state-of-the-art nucleus analysis methods are semi-supervised [23, 5, 18, 4, 35]: 1). Pretrained an autoencoder for unsupervised representation learning; 2). Construct a CNN from the pretrained autoencoder; 3). Fine-tune the constructed CNN for supervised nucleus classification. To better cap-

ture the visual variance of nuclei, one usually trains the unsupervised autoencoder on image patches with nuclei in the center [12, 20]. This requires a separate nucleus detection step [30] which in most cases needs tuning to optimize the final classification performance.

Instead of tuning the detection and classification modules separately, recent works [11, 25, 24, 15] successfully trained end-to-end CNNs to perform these tasks in a unified pipeline. Prior work has developed and employed supervised networks. To utilize unlabeled data for unsupervised pretraining, a network that can be trained end-to-end must perform unsupervised nucleus detection. Such unsupervised detection networks do not exist in any visual application domains, despite the success of unsupervised learning in other tasks [22, 8].

We design a novel Convolutional Autoencoder (CAE) that unifies nuclei detection and feature/representation learning in a single network and can be trained end-to-end without supervision. We also show that with existing labeled data, our network can be easily fine-tuned with supervision to improve the state-of-the-art performance of nuclei

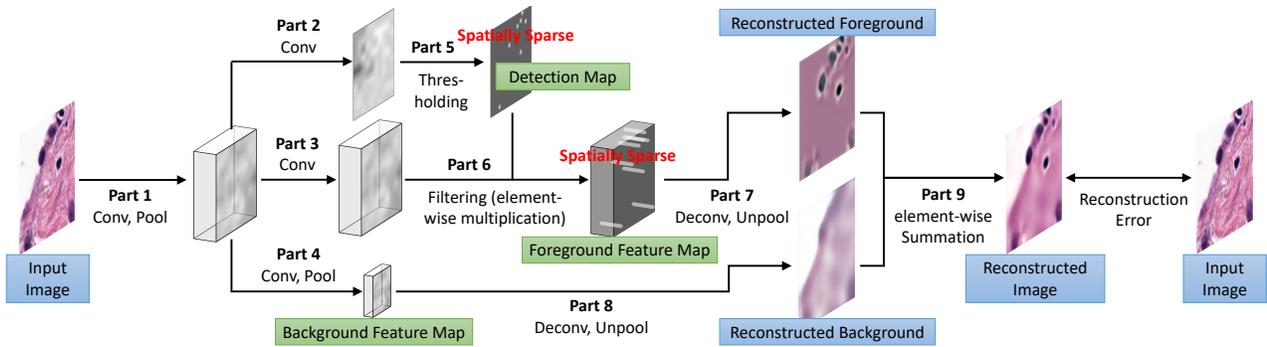


Figure 2: The architecture of our sparse Convolutional Autoencoder (CAE). The CAE minimizes image reconstruction error. The reconstructed image patch is a pixel-wise summation of two intermediate reconstructed image patches: the background and the foreground. The background is reconstructed from a set of small feature maps (background feature map) that can only encode large scale color and texture. The foreground is reconstructed from a set of crosswise sparse feature maps (foreground feature map). These foreground feature maps capture local high frequency signal: nuclei. We define *crosswise sparsity* as a type of sparsity with the following constraint: when there is no detected nucleus at a location, neurons in all foreground feature maps at the same location should not be activated. Details of network parts 1-8 are in Tab. 1.

classification and segmentation.

To perform unsupervised representation learning and detection, we modify the conventional CAE to encode not only appearance, but also spatial information in feature maps. To this end, our CAE first learns to separate background (*e.g.* cytoplasm) and foreground (*e.g.* nuclei) in an image patch, as shown in Fig. 2. We should note that an image patch is a rectangular region in a whole slide tissue image. We use image patches, because a tissue image can be very large and may not fit in memory. It is common in tissue image analysis to partition tissue images into patches and process the patches. We will refer to the partitioned image patches simply as the images. The CAE encodes the input image in a set of low resolution feature maps (background feature maps) with a small number of encoding neurons. The feature maps can only encode large scale color and texture variations because of their limited capacity. Thus these feature maps encode the image background. The high frequency residual between the input image and the reconstructed background is the foreground that contains nuclei.

The set of nuclei is often sparse in the image. We utilize this sparse property to design the foreground learning sub-network. Specifically, we design our network to learn the foreground feature maps in a “crosswise sparse” manner: neurons across all feature maps are not activated (output zero) in most feature map locations. Only neurons in a few feature map locations can be activated. Since the non-activated neurons have no influence in the later decoding layers, the image foreground is reconstructed using only the non-zero responses in the foreground encoding feature maps. This means that the image reconstruction error will be minimized only if the activated encoding neurons at different locations capture salient objects- the detected nuclei.

Learning a set of crosswise sparse foreground encoding feature maps is not straightforward. Neurons at the same location across all foreground feature maps should be synchronized: they should be activated or not activated at the same time depending on the presence of nuclei. In order to achieve this synchronization, the CAE needs to learn the locations of nuclei by optimizing the reconstruction error. Hence, the nucleus detection and feature extraction models are learned simultaneously during optimization. To represent the inferred nuclear locations, we introduce a special binary feature map: the nucleus detection map. We make this map sparse by thresholding neural activations. After optimization, a neuron in the nucleus detection map should output 1, if and only if there is a detected nucleus at the neuron’s location. The foreground feature maps are computed by element-wise multiplications between the nucleus detection map and a set of dense feature maps (Fig. 2).

In summary, our contributions are as follows.

1. We propose a CAE architecture with crosswise sparsity that can *detect* and represent nuclei in histopathology images with the following advantages:
 - As far as we know, this is the first unsupervised detection network for computer vision applications.
 - Our method can be fine-tuned for end-to-end supervised learning.
2. Our experimental evaluation using multiple datasets shows the proposed approach performs significantly better than other methods. The crosswise constraint in our method boosts the performance substantially.
 - Our method reduces the error of the best performing

baseline by up to 42% on classification tasks, and reduces the error of the U-net method [26] by 20% in nucleus segmentation.

- Our method achieves comparable results with 5% of training data needed by other methods, resulting in considerable savings in the cost of label generation.

2. Crosswise Sparse CAE

In this section we first introduce the The Convolutional Autoencoder (CAE) then describe our crosswise sparse CAE in detail. Our CAE reconstructs an image as the pixel-wise summation of a reconstructed foreground image and a reconstructed background image. The sparse foreground encoding feature maps represent detected nucleus locations and extracted nuclear features. The background feature maps are not necessarily sparse.

2.1. CAE for Semi-supervised CNN

An autoencoder is an unsupervised neural network that learns to reconstruct its input. The main purpose of this model is to learn a compact representation of the input as a set of neural responses [7]. A typical feedforward autoencoder is composed of an encoder and a decoder, which are separate layers. The encoding layer models the appearance information of the input. The decoder reconstructs the input from neural responses in the encoding layer. The CAE [18] and sparse CAE [23, 5, 35, 21] are autoencoder variants. One can construct a CNN with a trained CAE. Such semi-supervised CNNs outperform fully supervised CNNs significantly in many applications [14, 12].

2.2. Overall Architecture of Our CAE

The architecture of our CAE is shown in Fig. 2. We train the CAE to minimize the input image reconstruction error. The early stages of the CAE network consists of six convolutional and two average-pooling layers. The network then splits into three branches: the nucleus detection branch, the foreground feature branch, and the background branch. The detection branch merges into the foreground feature branch to generate the foreground feature maps that represent nuclei. The foreground and background feature maps are decoded to generate the foreground and background reconstructed images. Finally the two intermediate images are summed to form the final reconstructed image.

2.3. Background Encoding Feature Maps

We first model the background (tissue, cytoplasm etc.) then extract the foreground that contains nuclei. The biggest part of tissue images is background; its texture and color vary usually in a larger scale compared to the foreground. Also, usually a major portion of a tissue image is background. Thus, we encode an input image to a few small

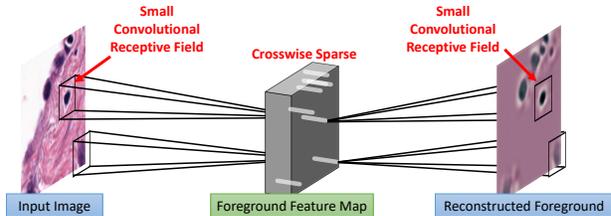


Figure 3: An illustration of how each nucleus is encoded and reconstructed. First, the foreground feature map must be crosswise sparse (Eq. 2). Second, the size of the receptive field of each encoding neuron should be small enough that it contains only one nucleus in most cases.

feature maps and assume those feature maps mostly contain the background information. In practice we represent the background of a 100×100 image by five 5×5 maps.

2.4. Foreground Encoding Feature Maps

Once the background is encoded and then reconstructed, the residual between the reconstructed background and the input image will be the foreground. The foreground consists of nuclei which are roughly of the same scale and often dispersed throughout the image. The foreground encoding feature maps encode everything about the nuclei, including their locations and appearance. A foreground feature map can be viewed as a matrix, in which each entry is a vector (a set of neuron responses) that encodes an image patch (the neurons’ receptive field). The vectors will encode nuclei, if there are nuclei at the center of the image patches. Otherwise the vectors contain zeros only. Since a small number of non-zero vectors encode nuclei, the foreground feature map will be sparse.

2.4.1 Crosswise Sparsity

We formally define crosswise sparsity as follows: We denote a set of f feature maps as X_1, X_2, \dots, X_f . Each feature map is a matrix. We denote the i, j -th entry of the l -th feature map as $X_l^{i,j}$, and the size of a feature map is $s \times s$. A conventional sparsity constraint requires:

$$\frac{\sum_{i,j,l} \mathbb{1}(X_l^{i,j} \neq 0)}{f s^2} \ll 1, \quad (1)$$

where $\mathbb{1}(\cdot)$ is the function that returns 1 if its input is true and 0 otherwise. Crosswise sparsity requires:

$$\frac{\sum_{i,j} \mathbb{1}\left(\sum_l \mathbb{1}(X_l^{i,j} \neq 0) > 0\right)}{s^2} \ll 1. \quad (2)$$

In other words, in most locations in the foreground feature maps, neurons across *all* the feature maps should *not* be

activated. This sparsity definition, illustrated in Fig. 3, can be viewed as a special form of group sparsity [19, 10].

If a foreground image is reconstructed by feature maps that are crosswise sparse, as defined by Eq. 2, the foreground image is essentially reconstructed by a few vectors in the feature maps. As a result, those vectors must represent salient objects in the foreground image- nuclei, since the CAE aims to minimize the reconstruction error.

2.4.2 Ensuring Crosswise Sparsity

Crosswise sparsity defined by Eq. 2 is not achievable using conventional sparsification methods [21] that can only satisfy Eq. 1. We introduce a binary matrix D with its i, j -th entry $D^{i,j}$ indicating if $X_l^{i,j}$ are activated for any l or not:

$$D^{i,j} = \mathbb{1}\left(\sum_l \mathbb{1}(X_l^{i,j} \neq 0) > 0\right). \quad (3)$$

Therefore Eq. 2 becomes:

$$\frac{\sum_{i,j} D^{i,j}}{s^2} \ll 1. \quad (4)$$

The foreground feature maps X_1, X_2, \dots, X_f are crosswise sparse, *iff* there exists a matrix D that satisfies Eq. 3 and Eq. 4. To satisfy Eq. 3, we design the CAE to generate a binary sparse feature map that represents D . The CAE computes X_l based on a dense feature map X'_l and D by element-wise multiplication:

$$X_l = X'_l \odot D. \quad (5)$$

We call the feature map D the detection map, shown in Fig. 2. The dense feature map X'_l is automatically learned by the CAE by minimizing the reconstruction error.

The proposed CAE also computes the D that satisfies Eq. 4. Notice that Eq. 4 is equivalent to the conventional sparsity defined by Eq. 1, when the total number of feature maps $f = 1$ and X_f is a binary feature map. Therefore, Eq. 4 can be satisfied by existing sparsification methods. A standard sparsification method is to add a sparsity penalty term in the loss function [21]. This method requires parameter tuning to achieve the desired expected sparsity. The k -sparse method [17] guarantees that exactly k neurons will be activated in D , where k is a predefined constant. However, in tissue images, the number of nuclei per image varies; the sparsity rate also should vary.

In this paper, we propose to use a threshold based method that guarantees an overall expected predefined sparsity rate, even though the sparsity rate for each CAE's input can vary. We compute the binary sparse feature map D as output from an automatically learned input dense feature map D' :

$$D^{i,j} = \text{sig}(r(D'^{i,j} - t)), \quad (6)$$

where $\text{sig}(\cdot)$ is the sigmoid function, r is a predefined slope, and t is an automatically computed threshold. We choose $r = 20$ in all experiments, making D a binary matrix in practice. Different r values do not affect the performance significantly based on our experience. Our CAE computes a large t in the training phase, which results in a sparse D . We define the expected sparsity rate as $p\%$, which can be set according to the average number of nuclei per image (we use $p = 1.6$ in all experiments), we compute t as:

$$t = E[\text{percentile}_p(D'^{i,j})], \quad (7)$$

where $\text{percentile}_p(D'^{i,j})$ is the p -th percentile of $D'^{i,j}$ for all i, j , given a particular CAE's input image. We compute t using the running average method:

$$t \leftarrow (1 - \alpha)t + \alpha \text{percentile}_p(D'^{i,j}). \quad (8)$$

We set the constant $\alpha = 0.1$ in all experiments. This running average approach is also used by batch normalization [13]. To make sure the running average of t converges, we also use batch normalization on $D'^{i,j}$ to normalize the distribution of $D'^{i,j}$ in each stochastic gradient descent batch. In total, three parameters are introduced in our CAE: r , p , and α . The sparsity rate p can be decided based on the dataset easily. The other two parameters do not affect the performance significantly in our experiments.

With crosswise sparsity, each vector in the foreground feature maps can possibly encode multiple nuclei. To achieve one-on-one correspondence between nuclei and encoded vectors, we simply reduce the size of the encoding neurons' receptive fields, such that a vector encodes a small region the size of which is close to the size of a nucleus.

3. Experiments

We initialize CNNs with our crosswise sparse CAEs. We empirically evaluate this approach on four datasets: a self-collected lymphocyte-rich region classification dataset, a self-collected individual lymphocyte classification dataset, the nuclear shape and attribute classification dataset, and the MICCAI 2015 nucleus segmentation challenge dataset [3]. The results show that the proposed method achieves better results than other methods.

3.1. Datasets

Dataset for Unsupervised Learning. We collected 0.5 million unlabeled small images randomly cropped from 400 lung adenocarcinoma histopathology images obtained from the public TCGA repository [2]. The cropped images are 100×100 pixels in 20X (0.5 microns per pixel). We will refer to cropped images simply as images in the rest of this section.

Dataset for Lymphocyte Classification Experiments (Sec 3.6). Lymphocytes and plasma cells are types of white

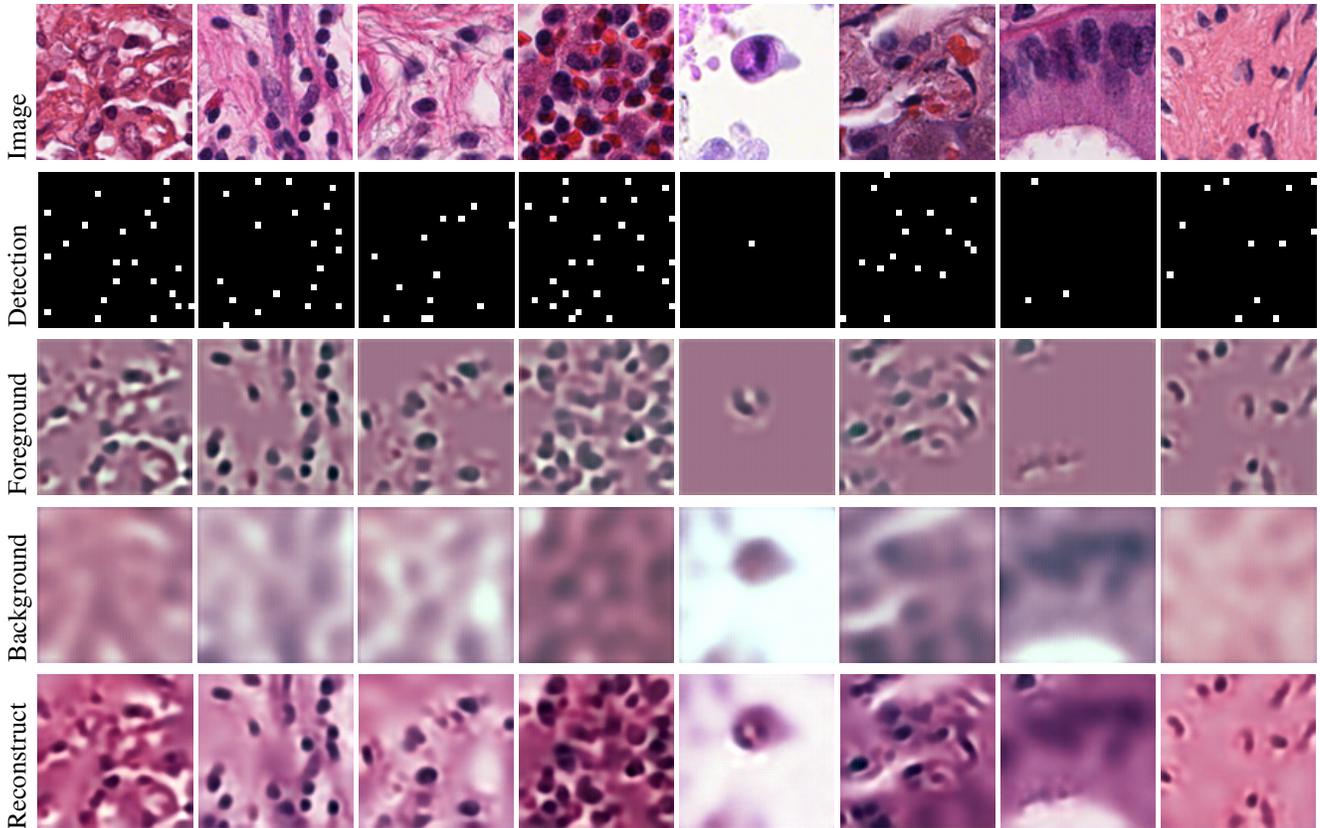


Figure 4: Randomly selected examples of unsupervised nucleus detection and foreground, background image representation (reconstruction) results. We show the detection map in Fig. 2 in the second row. The reconstructed image (last row) is the pixel-wise summation of the reconstructed foreground and background images. We can see that our CAE can decompose input images to foreground and background, and detect and represent (reconstruct) nuclei in the foreground.

blood cells in the immune system. Automatic recognition of lymphocytes and plasma cells is very important in many situations including the study of cancer immunotherapy [9, 28, 33]. We collected a dataset of 23,356 images labeled as lymphocyte (including plasma cells) rich or not. These images are cropped from lung adenocarcinoma tissue images in the publicly available TCGA dataset [2]. Each image is 300×300 pixels in $20X$ (0.5 microns per pixel). A pathologist labeled these images, according to the percentage of lymphocytes in the center 100×100 pixels of the image. The peripheral pixels provide context information to the pathologist and to automatic classification models. Overall, around 6% of the images are labeled as lymphocyte rich. We show examples of the training set in Fig. 5.

Dataset for Classifying Individual Lymphocytes Experiments (Sec. 3.7). We collected a dataset of 1785 images of individual objects that were labeled lymphocyte or non-lymphocyte by a pathologist. These 1785 images were cropped from 12 representative lung adenocarcinoma whole slide tissue images from the TCGA repository. We use la-

beled images cropped from 10 whole slide tissue images as the training data and the rest as the test dataset.

Dataset for Nuclear Shape and Attribute Classification Experiments (Sec. 3.8). We apply our method on an existing dataset [20] for nuclear shape and attribute classification. The dataset consists of 2000 images of nuclei labeled with fifteen morphology attributes and shapes.

Dataset for Nucleus Segmentation Experiments (Sec. 3.9). We test our method for nucleus segmentation using the MICCAI 2015 nucleus segmentation challenge dataset [3] which contains 15 training images and 18 testing images with a typical resolution of 500×500 . The ground truth masks of nuclei are provided in the training dataset.

3.2. CAE Architecture

The CAEs in all three classification experiments (Sec. 3.6, Sec. 3.7 and Sec. 3.8) are trained on the unlabeled dataset with the same architecture illustrated in Fig. 2 and Tab. 1. Note that we apply batch normalization [13] before

Part	Layer	Kernel size	Stride	Output size
1	Input	-	-	$100^2 \times 3$
	Conv	5×5	1	$100^2 \times 100$
	Conv	5×5	1	$100^2 \times 120$
	Pool	2×2	2	$50^2 \times 120$
	Conv	3×3	1	$50^2 \times 240$
	Conv	3×3	1	$50^2 \times 320$
	Pool	2×2	2	$25^2 \times 320$
	Conv	3×3	1	$25^2 \times 640$
	Conv	3×3	1	$25^2 \times 1024$
2	Conv	1×1	1	$25^2 \times 100$
	Conv	1×1	1	$25^2 \times 1$
3	Conv	1×1	1	$25^2 \times 640$
	Conv	1×1	1	$25^2 \times 100$
4	Conv	1×1	1	$25^2 \times 128$
	Pool	5×5	5	$5^2 \times 128$
	Conv	3×3	1	$5^2 \times 64$
	Conv	1×1	1	$5^2 \times 5$
5	Thres.	Defined by Eq. 6		$25^2 \times 1$
6	Filter.	Defined by Eq. 5		$25^2 \times 100$
7	Deconv	3×3	1	$25^2 \times 1024$
	Deconv	3×3	1	$25^2 \times 640$
	Deconv	4×4	0.5	$50^2 \times 640$
	Deconv	3×3	1	$50^2 \times 320$
	Deconv	3×3	1	$50^2 \times 320$
	Deconv	4×4	0.5	$100^2 \times 320$
	Deconv	5×5	1	$100^2 \times 120$
	Deconv	5×5	1	$100^2 \times 100$
	Deconv	1×1	1	$100^2 \times 3$
8	Deconv	3×3	1	$5^2 \times 256$
	Deconv	3×3	1	$5^2 \times 128$
	Deconv	9×9	0.2	$25^2 \times 128$
	Deconv	3×3	1	$25^2 \times 128$
	Deconv	3×3	1	$25^2 \times 128$
	Deconv	4×4	0.5	$50^2 \times 128$
	Deconv	3×3	1	$50^2 \times 64$
	Deconv	3×3	1	$50^2 \times 64$
	Deconv	4×4	0.5	$100^2 \times 64$
	Deconv	5×5	1	$100^2 \times 32$
	Deconv	5×5	1	$100^2 \times 32$
	Deconv	1×1	1	$100^2 \times 3$

Table 1: Layer setup of different parts in our CAE. Please refer to Fig. 2 for the overall network architecture.

the leaky ReLU activation functions [16] in all layers.

The average nucleus size in the dataset for nucleus segmentation experiments (Sec. 3.9) is around 20×20 pixels. Therefore pooling layers can discard important spatial information which is important for pixel-wise segmentation. The U-net [26] addresses this issue by adding skip connections. However, we find in practice that eliminating pooling

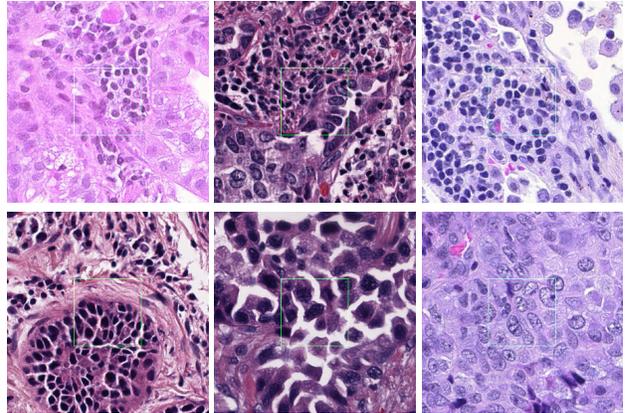


Figure 5: Examples of the lymphocyte rich region dataset. Top: lymphocyte rich images. Bottom: non-lymphocyte-rich images. A pathologist labeled 23,356 images depending on the percentage of lymphocytes in the center 100×100 pixels of the image (framed in green). The peripheral pixels provide context information to the pathologist and automatic classification models.

layers completely yields better performance. The computation complexity is very high for a network without any pooling layers. Thus, compared to Tab. 1, we use smaller input dimensions (40×40) and fewer (80 to 200) feature maps in the CAE. Other settings of the CAE for segmentation remain unchanged.

3.3. CNN Architecture

We construct supervised CNNs based on trained CAEs. For classification tasks, the supervised CNN contain Parts 1-6 of the CAE. We initialize the parameters in these layers to be the same as the parameters in the CAE. We attach four 1×1 convolutional layers after the foreground encoding layer and two 3×3 convolutional layers after the background encoding layer. Each added layer has 320 convolutional filters. We then apply global average pooling on the two branches. The pooled features are then concatenated together, followed by a final classification layer with sigmoid activation function. For the segmentation task, the supervised CNN only contains Parts 1, 3 of the CAE. We attach six 3×3 convolutional layers followed by a segmentation layer. The segmentation layer is the same to the patch-CNN’s [34] segmentation layer which is a fully-connected layer with sigmoid activation function followed by reshaping. For all tasks, we randomly initialize the parameters of these additional layers. We train the parameters of the added layers until convergence before fine-tuning the whole network.

3.4. Learning Details

We train our CAE on the unlabeled dataset, minimizing the pixel-wise root mean squared error between the input images and the reconstructed images. We use stochastic gradient descent with batch size 32, learning rate 0.03 and momentum 0.9. The loss converges after 6 epochs. We show randomly selected examples of the nucleus detection feature map as well as the reconstructed foreground and background images in Fig. 4.

For the CNN (constructed from the CAE) training, we use stochastic gradient descent with batch size, learning rate, and momentum selected for each task independently. For all tasks, we divide the learning rate by 10 when the error has plateaued. We use sigmoid as the nonlinearity function in the last layer and log-likelihood as the loss function. We apply three types of data augmentation. First, the input images are randomly cropped from a larger image. Second, the colors of the input images are randomly perturbed. Third, we randomly rotate and mirror the input images. During testing, we average the predictions of 25 image crops. We implemented our CAE and CNN using Theano [31]. We trained the CAE and CNN on a single Tesla K40 GPU.

3.5. Methods Tested

We describe our method (abbreviated as CSP-CNN) and other tested methods below:

CSP-CNN CNN initialized by our proposed crosswise sparse CAE shown in Fig. 2. The exact CNN construction is described in Sec. 3.3. We set the sparsity rate to 1.6%, such that the number of activated foreground feature map locations roughly equals to the average number of nuclei per image in the unsupervised training set.

SUP-CNN A fully supervised CNN. Its architecture is similar to our CSP-CNN except that: 1). There is no background representation branch (no Part 4, 8 in Fig. 2). 2). There is no nucleus detection branch (no Part 2, 5 in Fig. 2). The SUP-CNN has a very standard architecture, at the same time similar to our CSP-CNN.

U-NET We use the authors’ U-net architecture and implementation [26] for nucleus segmentation. We test five U-nets with the same architecture but different number of feature maps per layer and select the best performing network. All five U-nets perform similarly.

DEN-CNN CNN initialized by a conventional Convolutional Autoencoder (CAE) without the sparsity constraint. Its architecture is similar to our CSP-CNN except that there is no nucleus detection branch. In particular, there is no Part 2 and Part 5 in Fig. 2 and Part 6 is an identity mapping layer.

SP-CNN CNN initialized by a sparse CAE without the crosswise constraint. Its architecture is similar to our CSP-CNN except that it has no nucleus detection branch and uses the conventional sparsity constraint defined by Eq. 1. In particular, there is no Part 2 and Part 5 in Fig. 2 and Part 6 is a thresholding layer: define its input as D' , its output $D = \text{ReLU}(D' - t)$, where t is obtained in the same way defined by Eq. 7. We set the sparsity rate to 1.6% which equals to the rate we use in CSP-CNN.

VGG16 We finetune the VGG 16-layer network [29] which is pretrained on ImageNet [27]. Fine-tuning the VGG16 network has been shown to be robust for pathology image classification [36, 12].

Methods	Datasets		
	Lym-region	Individual Lym	Nuclear Attr & Shape [20]
SUP-CNN	0.6985	0.4936	0.8487
DEN-CNN	0.8764	0.5576	0.8656
SP-CNN	0.9188	0.6262	0.8737
CSP-CNN	0.9526	0.7856	0.8788
CSP-CNN (5% data)	0.9215	0.7135	0.7128
Unsupervised features [37]	-	0.7132	-
Semi-supervised CNN [20]	-	-	0.8570
VGG16 [29]	0.9176	0.6925	0.8480

Table 2: Classification results measured by AUROC on there tasks described in Sec. 3.6, Sec. 3.7, Sec. 3.8. The proposed CSP-CNN outperforms the other methods significantly. Comparing the results of SP-CNN and our CSP-CNN, we can see that the proposed crosswise constraint boosts performance significantly. Even with only 5% labeled training data, our CSP-CNN (5% data) outperforms other methods on the first two datasets. The CSP-CNN (5% data) fails on the third dataset because when only using 5% of the training data, 5 out of 15 classes have less than 2 positive training instances which are too few for CNN training.

3.6. Classifying Lymphocyte-rich Regions

We use 20,876 images as the training set and the remaining 2,480 images as testing set. We use the Area Under ROC Curve (AUROC) as the evaluation metric. The results are shown in Tab. 2. Our proposed CSP-CNN achieves the best result on this dataset. Our CSP-CNN reduces the error of the best performing baseline SP-CNN by **42%**. Furthermore, with only 5% of the training data, our CSP-CNN (5% data) outperforms SP-CNN. The only difference between CSP-CNN and SP-CNN is the crosswise constraint, with

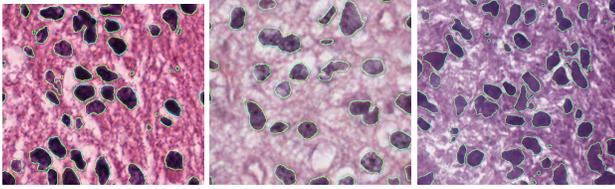


Figure 6: Randomly selected examples of nucleus segmentation using our CSP-CNN, on the MICCAI 2015 nucleus segmentation challenge dataset (best viewed in color). The segmentation boundaries are displayed in green.

which our CAE is capable of unsupervised nucleus detection. This supports our claim that our crosswise sparsity is essential to high performance.

3.7. Classifying Individual Lymphocytes

We compare our method with an existing unsupervised nucleus detection and feature extraction method [37]. We split training and test images 4 times and average the results. As the baseline method we carefully tuned a recently published unsupervised nucleus detection and feature extraction method [37], which is based on level sets, and applied a multi-layer neural network on top of the extracted features. We should note that the feature extraction step and the classification step have to be tuned separately in the baseline method, whereas our CSP-CNN method (Sec. 3.6) can be trained end-to-end. We show results in Tab. 2. Our CSP-CNN reduced the error of the SP-CNN by **25%**.

3.8. Nuclear Shape and Attribute Classification

We adopt the same 5-fold training and testing data separation protocol and report the results in Tab. 2. On this dataset the improvement of our method over the state-of-the-art is less significant than the improvement on other datasets because the images of nuclei are results of a fixed nucleus detection method which we cannot fine-tune with our proposed method.

3.9. Nucleus Segmentation

We use a sliding window approach to train and test our CNNs. A CNN outputs a feature map of the same size as its input. For evaluation, we followed the standard metric used in the MICCAI challenge: the DICE-average (average of two different versions of the DICE coefficient). We show results in Tab. 3. The proposed method achieves a significantly higher score than that of the challenge winner [6] and U-net [26]. Because the size of nuclei are only around 20×20 pixels, we eliminate our network’s pooling layers completely to preserve spatial information. We believe this is an important reason our method outperforms U-net. We show randomly selected segmentation examples in Fig. 6.

Methods	DICE-average
SUP-CNN	0.8216
DEN-CNN	0.8235
SP-CNN	0.8338
CSP-CNN	0.8362
CSP-CNN (5% data)	0.8205
Challenge winner [6]	0.80
U-net [26]	0.7942

Table 3: Nucleus segmentation results on the MICCAI 2015 nucleus segmentation challenge dataset. Our CSP-CNN outperforms the highest challenge score which is a DICE-average of 0.80, even with only 5% of the sliding windows during training. On this dataset we do not use pooling layers in the networks, because we find that pooling layers discard important spatial information, since the size of nuclei are only around 20×20 pixels.

4. Conclusions

We propose a crosswise sparse Convolutional Autoencoder (CAE) that for the first time, is capable of unsupervised nucleus detection and feature extraction simultaneously. We advocate that this CAE should be used to initialize classification or segmentation Convolutional Neural Networks (CNN) for supervised training. In this manner, the nucleus detection, feature extraction, and classification or segmentation steps are trained end-to-end. Our experimental results with four challenging datasets show that our proposed crosswise sparsity is essential to state-of-the-art results. In a future work we plan to perform unsupervised nucleus segmentation with the proposed CAE.

References

- [1] Pathologists’ Hourly Wages. <http://ww1.salary.com/Physician-Pathology-hourly-wages.html>. 1
- [2] The Cancer Genome Atlas. <https://cancergenome.nih.gov/>. 4, 5
- [3] Miccai 2015 workshop and challenges in imaging and digital pathology: The computational brain tumor cluster of event. <https://wiki.cancerimagingarchive.net/pages/viewpage.action?pageId=20644646>, 2015. 4, 5
- [4] N. Bayramoglu and J. Heikkilä. Transfer learning for cell nuclei classification in histopathology images. In *ECCV Workshops*, 2016. 1
- [5] Y.-I. Boureau, Y. L. Cun, et al. Marc’aurelio ranzato, y and boureau, lan and lecu, yann. In *NIPS*, 2008. 1, 3
- [6] H. Chen, X. Qi, L. Yu, Q. Dou, J. Qin, and P.-A. Heng. DCAN: Deep contour-aware networks for object instance segmentation from histology images. *Medical Image Analysis*, 2017. 8

- [7] L. Deng, M. L. Seltzer, D. Yu, A. Acero, A.-r. Mohamed, and G. E. Hinton. Binary coding of speech spectrograms using a deep auto-encoder. In *Interspeech*, 2010. 3
- [8] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 1
- [9] J. Galon, A. Costes, F. Sanchez-Cabo, A. Kirilovsky, B. Mlecnik, C. Lagorce-Pagès, M. Tosolini, M. Camus, A. Berger, P. Wind, et al. Type, density, and location of immune cells within human colorectal tumors predict clinical outcome. *Science*, 313, 2006. 5
- [10] B. Graham. Spatially-sparse convolutional neural networks. *arXiv preprint arXiv:1409.6070*, 2014. 4
- [11] A. Graves and N. Jaitly. Towards end-to-end speech recognition with recurrent neural networks. In *ICML*, volume 14, 2014. 1
- [12] L. Hou, K. Singh, D. Samaras, T. M. Kurc, Y. Gao, R. J. Seidman, and J. H. Saltz. Automatic histopathology image analysis with cnns. In *New York Scientific Data Summit*, 2016. 1, 3, 7
- [13] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 4, 5
- [14] R. Johnson and T. Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In *NIPS*, 2015. 3
- [15] I. Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *CVPR*, 2017. 1
- [16] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML*, 2013. 6
- [17] A. Makhzani and B. Frey. k-sparse autoencoders. 2014. 4
- [18] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *International Conference on Artificial Neural Networks (ICANN)*, 2011. 1, 3
- [19] C. Murdock, Z. Li, H. Zhou, and T. Duerig. Blockout: Dynamic model selection for hierarchical deep networks. In *CVPR*, 2016. 4
- [20] V. Murthy, L. Hou, D. Samaras, T. M. Kurc, and J. H. Saltz. Center-focusing multi-task CNN with injected features for classification of glioma nuclear images. In *Winter Conference on Applications of Computer Vision (WACV)*, 2017. 1, 5, 7
- [21] A. Ng. Sparse autoencoder. *CS294A Lecture notes*, 72, 2011. 3, 4
- [22] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2016. 1
- [23] M. Ranzato, C. Poultney, S. Chopra, and Y. LeCun. Efficient learning of sparse representations with an energy-based model. In *NIPS*, 2006. 1, 3
- [24] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 1
- [25] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 1
- [26] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 3, 6, 7, 8
- [27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 7
- [28] R. Salgado, C. Denkert, S. Demaria, N. Sirtaine, F. Klauschen, G. Pruneri, S. Wienert, G. Van den Eynden, F. L. Baehner, F. Penault-Llorca, et al. The evaluation of tumor-infiltrating lymphocytes (TILs) in breast cancer: recommendations by an international TILs working group 2014. *Annals of oncology*, 26, 2015. 5
- [29] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2014. 7
- [30] K. Sirinukunwattana, S. E. A. Raza, Y.-W. Tsang, D. R. Snead, I. A. Cree, and N. M. Rajpoot. Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images. *Medical Imaging*, 35, 2016. 1
- [31] Theano Development Team. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688, 2016. 7
- [32] A. S. Tripathi, A. Mathur, M. Daga, M. Kuse, and O. C. Au. 2-simdom: a 2-sieve model for detection of mitosis in multi-spectral breast cancer imagery. In *International Conference on Image Processing (ICIP)*.
- [33] R. Turkki, N. Linder, P. E. Kovanen, T. Pellinen, and J. Lundin. Antibody-supervised deep learning for quantification of tumor-infiltrating immune cells in hematoxylin and eosin stained breast cancer samples. *Journal of Pathology Informatics (JPI)*, 7, 2016. 5
- [34] T. F. Y. Vicente, L. Hou, C.-P. Yu, M. Hoai, and D. Samaras. Large-scale training of shadow detectors with noisily-annotated shadow examples. In *ECCV*, 2016. 6
- [35] J. Xu, L. Xiang, Q. Liu, H. Gilmore, J. Wu, J. Tang, and A. Madabhushi. Stacked sparse autoencoder (ssae) for nuclei detection on breast cancer histopathology images. *Medical Imaging*, 2016. 1, 3
- [36] Y. Xu, Z. Jia, Y. Ai, F. Zhang, M. Lai, I. Eric, and C. Chang. Deep convolutional activation features for large scale brain tumor histopathology image classification and segmentation. In *ICASSP*, 2015. 7
- [37] N. Zhou, X. Yu, T. Zhao, S. Wen, F. Wang, W. Zhu, T. Kurc, A. Tannenbaum, J. Saltz, and Y. Gao. Evaluation of nucleus segmentation in digital pathology images through large scale image synthesis. *SPIE Medical Imaging*, 2017. 7, 8