

Automatic Microscopic Cell Counting by Use of Deeply-Supervised Density Regression Model

Shenghua He¹, Kyaw Thu Minn^{2,3}, Lilianna Solnica-Krezel³, Mark Anastasio^{1,2,4}, and Hua Li⁴

¹Department of Computer Science and Engineering

²Department of Biomedical Engineering

³Department of Developmental Biology

⁴Department of Radiation Oncology

Washington University in St. Louis, St. Louis, MO, USA

ABSTRACT

Accurately counting cells in microscopic images is important for medical diagnoses and biological studies, but manual cell counting is very tedious, time-consuming, and prone to subjective errors, and automatic counting can be less accurate than desired. To improve the accuracy of automatic cell counting, we propose here a novel method that employs deeply-supervised density regression. A fully convolutional neural network (FCNN) serves as the primary FCNN for density map regression. Innovatively, a set of auxiliary FCNNs are employed to provide additional supervision for learning the intermediate layers of the primary CNN to improve network performance. In addition, the primary CNN is designed as a concatenating framework to integrate multi-scale features through shortcut connections in the network, which improves the granularity of the features extracted from the intermediate CNN layers and further supports the final density map estimation. The experimental results on immunofluorescent images of human embryonic stem cells demonstrate the superior performance of the proposed method over other state-of-the-art methods.

Keywords: Automatic cell counting, microscopic images, density regression, deeply-supervised learning, concatenating network

1. INTRODUCTION

Accurately counting the number of cells in microscopic images greatly aids medical diagnoses and biological studies.¹ However, manual cell counting, which is slow, expensive, and prone to subjective errors, is not practically feasible for high-throughput processes. An automatic and efficient solution with improved counting accuracy is highly desirable, but automatic cell counting is challenging due to the low image contrast, strong tissue background, and significant inter-cell occlusions in 2D microscopic images.²⁻⁵

To address these challenges,^{6,7} density-based counting methods have received increasing attention, due to their superior performance to those traditional detection-based methods.^{5,8-10} Generally speaking, density-based methods employ machine-learning tools to learn a density regression model (DRM) that estimates the cell density distribution from the characteristics/features of a given image. The number of cells can be subsequently estimated by integrating the estimated density map. For example, Lempitsky et al.⁶ proposed a supervised learning framework to learn a linear DRM and employ it for visual object counting tasks. Differently, Xie et al.⁷ utilized a fully convolutional regression network (FCRN) to learn a DRM for regressing cell spatial densities over the image. The FCRN, as a specific fully convolutional neural network (FCNN), integrates informative feature extraction and powerful function learning to estimate the density map of a given image. It demonstrates promising performance for cell counting tasks, especially for counting overlapped cells.

Generally, the layers in the FCRN are hierarchically constructed, and the output of each layer relies on the outputs of its previous layers. There are two potential shortcomings of this traditional network design. First, the intermediate layers are optimized based on the gradient information back-propagated only from the final layer of

Further author information: (Send correspondence to Hua Li & Mark Anastasio)
E-mail: {li.hua, anastasio}@wustl.edu, Telephone: 1 314 537 7145 & 1 314 935 3637

the network, not directly from the adjacent layer. Second, this design allows only adjacent layers to be connected, while limiting the integration of multi-scale features (or information) and overall network performance. These two shortcomings might lead to sub-optimized intermediate layers, and can eventually affect the overall cell counting accuracy. Thus, the need to improve the accuracy of automated cell counting methods remains.

Recently, in order to improve the effectiveness of learning intermediate layers in a designed deep neural networks, deeply-supervised learning (or deep supervision) has been proposed and has shown promising performance for addressing various computer vision tasks, such as classification¹¹ and segmentation.^{12,13} In addition, concatenating CNN frameworks have also attracted great attention. These networks can concatenate multi-scale features by shortcut connections of non-adjacent layers within the network, and so achieve better results than the traditional networks in such computer vision tasks as segmentation¹⁴ and detection.¹⁵

Motivated by these works, this study proposes a novel density regression-based automatic cell counting method. A FCNN is used as a primary FCNN (PriCNN) to learn the density regression model (DRM) that performs an end-to-end mapping from a cell image to the corresponding density map. A set of auxiliary CNN (AuxCNNs) are built to assess the features at the intermediate layers in the PriCNN and to directly supervise the training of these layers. In addition, by use of concatenation layers, the multi-scale features from non-adjacent layers are integrated to improve the granularity of the features extracted from the intermediate layers for further supporting final density map estimation. Experimental results, evaluated on a set of immunofluorescent images of human embryonic stem cells (hESC), have demonstrated the superior performance of the proposed deep supervision-based DRM method compared to other state-of-the-art methods.

2. METHODOLOGY

2.1 Background: Density-Based Automatic Cell Counting

The goal of density regression-based cell counting methods is to learn a density regression function F , which can be employed to estimate the density map of a given image.^{6,7} Given an image $X \in \mathbb{R}^{M \times N}$ which includes N_c cells, the density map $Y \in \mathbb{R}^{M \times N}$ of X can be considered as the superposition of a set of N_c normalized 2D discrete Gaussian kernels that are placed at the centroids of the N_c cells. Therefore, the number of cells can be counted by integrating the density map over the image.

Let $S = \{(s_{k_x}, s_{k_y}) \in \mathbb{N}^2 : k = 1, 2, \dots, N_c\}$ represent the cell centroid positions in X . Each pixel $Y_{i,j}$ on the density map Y can be expressed as:

$$Y_{i,j} = \sum_{k=1}^{N_c} G_\sigma(i - s_{k_x}, j - s_{k_y}), \quad \forall \quad i \in M, j \in N, \quad (1)$$

where $G_\sigma(n_x, n_y) = C \cdot e^{-\frac{n_x^2 + n_y^2}{2\sigma^2}} \in \mathbb{R}^{(2K_G+1) \times (2K_G+1)}$, $n_x, n_y = -K_G, -K_G + 1, \dots, K_G$, is a normalized 2D Gaussian kernel that satisfies $\sum_{n_x=-K_G}^{K_G} \sum_{n_y=-K_G}^{K_G} G_\sigma(n_x, n_y) = 1$. Here, σ^2 is the isotropic covariance, $(2K_G+1)$ is the kernel size, and C is a normalization constant.

The density regression-based cell counting process generally includes three steps: (1) map an image to a feature map, (2) estimate a cell density map from the feature map, and (3) integrate the density map for automatic cell counting. In the first step, each pixel $X_{i,j}$ in X can be assumed to be associated with a real-valued feature vector $\phi(X_{i,j}) \in \mathbb{R}^Z$. The feature map $P \in \mathbb{R}^{M \times N \times Z}$ of X can be generated using specific feature extraction methods, such as the dense scale invariant feature transform (SIFT) descriptor,¹⁶ ordinary filter banks,¹⁷ or codebook learning.¹⁸ In the second step, the estimated density $\hat{Y}_{i,j}$ of each pixel $X_{i,j}$ in X can be obtained by applying a pre-trained density regression function F on the given $\phi(X_{i,j})$:

$$\hat{Y}_{i,j} = F(\phi(X_{i,j}); \Theta), \quad (2)$$

where Θ is a parameter vector that determines the function F . Finally, in the third step, the number of cells in X , N_c , can be counted by integrating the estimated densities \hat{Y} over the image region:

$$N_c \approx \hat{N}_c = \sum_{i=1}^M \sum_{j=1}^N \hat{Y}_{i,j}. \quad (3)$$

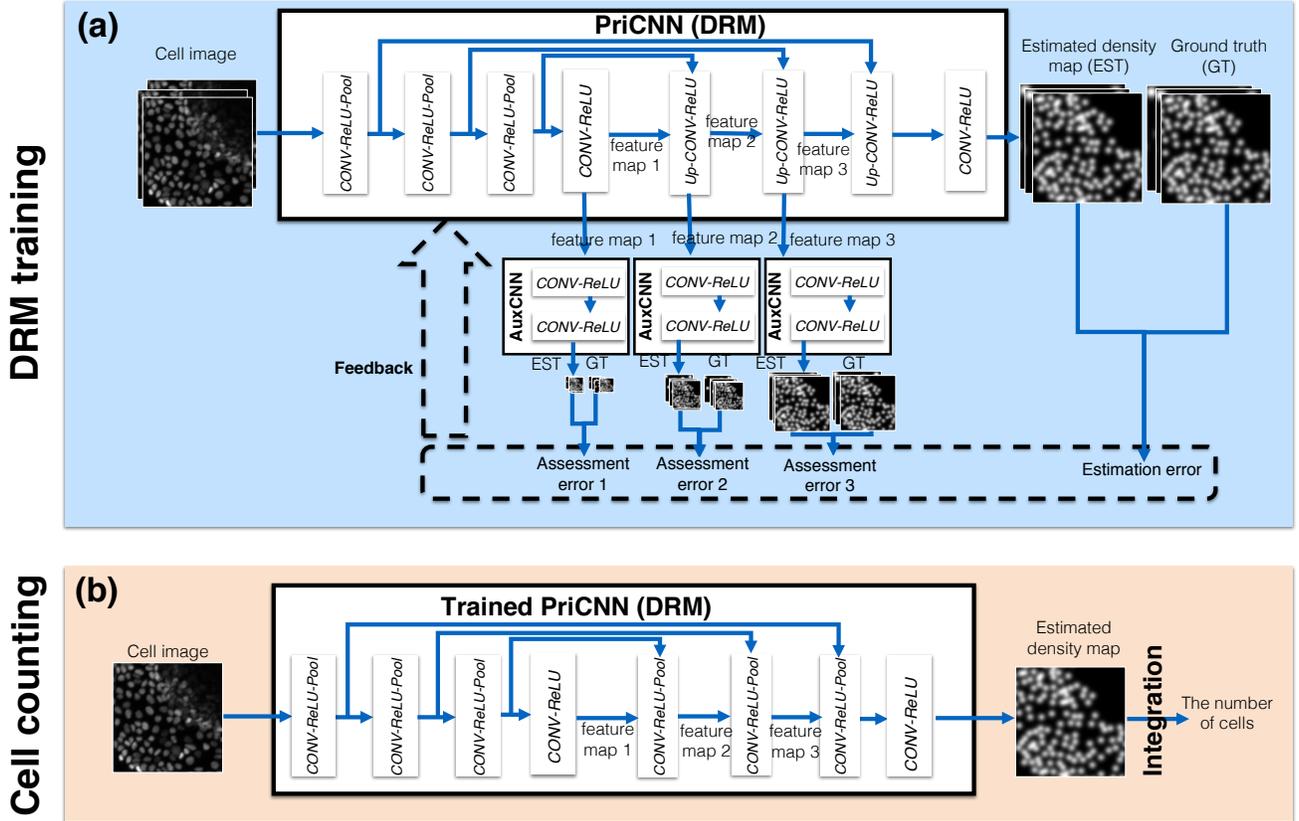


Figure 1: The framework of the proposed automatic cell counting method. The first phase (a) is DRM training, and the second phase (b) is cell counting by use of the trained DRM.

A key task in density regression-based cell counting methods is learning the function F by use of training datasets. The learning of F and the related cell counting method proposed in this study are described below.

2.2 The Proposed Automatic Cell Counting Framework

We propose a novel automatic cell counting method that employs deeply-supervised density regression model in this study. The framework, shown in Figure 1, includes two phases: 1) DRM training and 2) cell counting by use of the trained DRM. The network architecture of the proposed DRM is described in Section 2.2.1. The two phases in the framework are described in Sections 2.2.2 and 2.2.3, respectively.

2.2.1 The DRM Network Architecture

The DRM is built as a primary FCNN (PriCNN) with the purpose of estimating the density map \hat{Y} of an image X , such that:

$$Y \approx \hat{Y} = F(X; \Theta), \quad (4)$$

where $F(X; \Theta)$ is a density regression function, and Θ is a parameter vector that determines F .

Motivated by the network architecture of FCRN,⁷ the designed PriCNN (DRM) includes 8 chained blocks. As shown in the upper row of Figure 1, each of the first three blocks includes a convolutional (CONV) layer, a ReLU layer, and a max-pooling (Pool) layer; the fourth block in the PriCNN includes a CONV layer and a ReLU layer; each of the fifth to seventh blocks includes a up-sampling (UP) layer, a CONV layer, and a ReLU layer; and the last block includes a chain of a CONV layer and a ReLU layer.

In addition, concatenation layers are employed in the PriCNN to integrate multi-scale features and thus improve the granularity of the features, which assists in the final density map estimation. This design is motivated

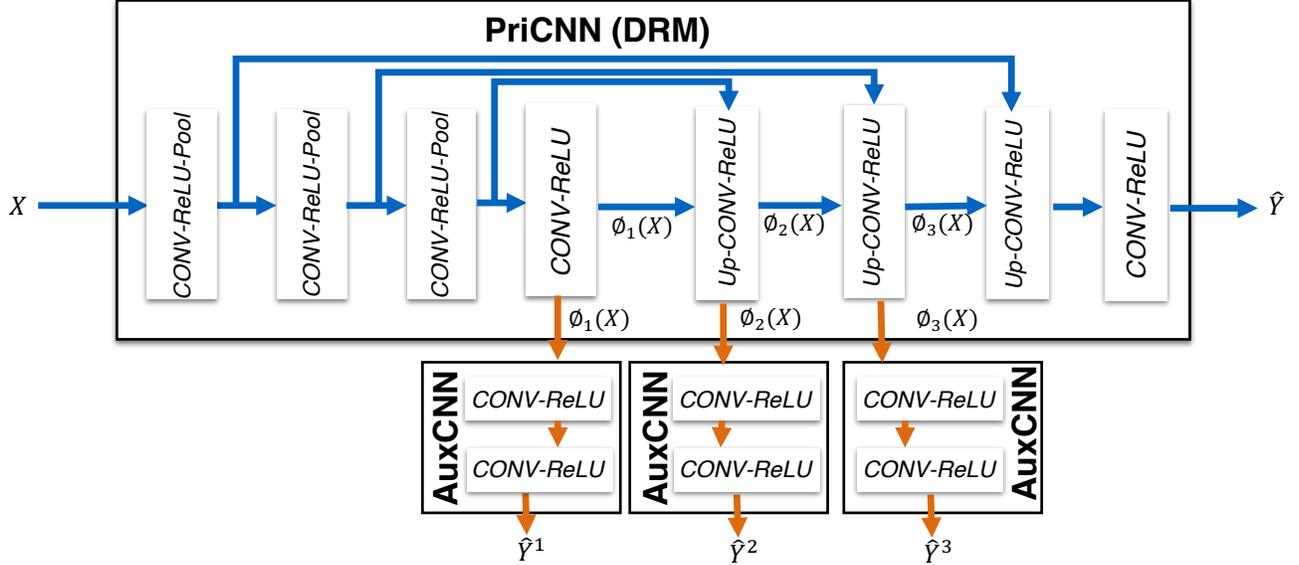


Figure 2: The network architectures of the PriCNN and 3 AuxCNNs in the proposed approach.

by a network architecture described by Ronneberger et. al.¹⁴ As shown in Figure 1, the outputs from each of the first three blocks are multi-resolution, low-dimension, and highly-representative feature maps. Three shortcut connections are established to connect the first and seventh blocks, the second and sixth blocks, and the third and fifth blocks, respectively. With these shortcut connections, multi-resolution features can be concatenated between non-adjacent layers. The integration of multi-scale features can further improve the performance of the network compared to the traditional FCRN, which allows only adjacent layers to be connected.

2.2.2 DRM Training Process

1. AuxCNN-supported DRM Training

Training the designed DRM (or PriCNN) with such a hierarchical structure is a challenging task. As described in Section 1, all the layers in the original FCRN⁷ are learned based on the feedback only from the final layer. Therefore, the intermediate layers might be sub-optimized, which can significantly affect the accuracy of the final estimated density map.

Innovatively, three auxiliary FCNNs (AuxCNNs) are employed to provide additional supervision for learning the intermediate layers of the PriCNN. As shown in Figure 2, each AuxCNN contains two CONV-ReLU blocks for estimating a low-resolution density map from the feature map generated at an intermediate layer of the PriCNN. By jointly minimizing the errors between the estimated density maps and the corresponding ground truth density maps at different resolution levels, the optimization of the intermediate layers in the PriCNN can be improved, which eventually improves the overall performance of the PriCNN.

2. Jointly Training the PriCNN and AuxCNN

The parameters of the PriCNN can be learned by jointly training the PriCNN and the AuxCNNs with a set of given training data $D = \{(X_i, Y_i)\}_{i=1,2,\dots,B}$, where X_i and Y_i represent the i -th image and its associated ground-truth density map, respectively. The training is completed by minimizing the differences between the estimated density maps at different resolution levels and the ground truth density maps.

As shown in Figure 2, the PriCNN can be denoted mathematically as $F(X; \Theta)$, where Θ is the parameter vector of F and X is an input image. All the trainable parameters in the first four blocks, the 5th, the 6th, and the last two blocks can be denoted as Θ_1 , Θ_2 , Θ_3 , and Θ_4 , respectively. Therefore, $\Theta = (\Theta_1, \Theta_2, \Theta_3, \Theta_4)$ and $F(X; \Theta) = F(X; \Theta_1, \Theta_2, \Theta_3, \Theta_4)$. Also, the output feature maps of the 4-th, 5-th, and 6-th blocks can be denoted as $\phi_1(X; \Theta_1)$, $\phi_2(X; \Theta_1, \Theta_2)$, and $\phi_3(X; \Theta_1, \Theta_2, \Theta_3)$, respectively. Similarly, the three AuxCNNs can be denoted as $A_k(\phi_k; \theta_k)$ ($k = 1, 2, 3$), where θ_k is the parameter vector of the k -th AuxCNN, and A_k is the density

map estimated by use of the k -th AuxCNN. Therefore, the cooperative training of the PriCNN and AuxCNNs is performed by jointly minimizing four loss functions, defined below:

$$\begin{cases} L_k(\Theta_1, \dots, \Theta_k, \theta_k) &= \frac{1}{B} \sum_{i=1}^B \|A_k(\phi_k(X_i; \Theta_1, \dots, \Theta_k); \theta_k) - Y_i^k\|^2, k = 1, 2, 3, \\ L(\Theta) &= \frac{1}{B} \sum_{i=1}^B \|F(X_i, \Theta) - Y_i\|^2, \end{cases} \quad (5)$$

where Y_i^k is the ground truth low-resolution density map (GTLR) generated from Y_i . $Y_i^k \in \mathbb{R}^{M_k \times N_k}$ is generated from the original ground-truth density map $Y_i \in \mathbb{R}^{M \times N}$ by summing every adjunct $a \times b$ in Y , with $a_k = \frac{M}{M_k}$ and $b_k = \frac{N}{N_k}$. $L(\Theta)$ is the average mean square error (MSE) between the estimated density maps and their ground truths. $L_k(\Theta_1, \dots, \Theta_k, \theta_k)$ is the average MSE between the low-resolution density maps estimated by k -th AuxCNN and their corresponding GTLR density maps.

To improve the computational efficiency of the optimization of the PriCNN and AuxCNNs, we construct a combined loss function, defined as below:

$$L_{overall}(\Theta, \theta_1, \theta_2, \theta_3) = L(\Theta) + \sum_{k=1}^3 \lambda_k L_k(\Theta_1, \dots, \Theta_k, \theta_k), \quad (6)$$

where $\lambda_k \in [0, 1]$ is a parameter that controls the relative strength of the supervision under the k -th AuxCNN for learning the intermediate layers in the PriCNN. Eqn.(6) is numerically minimized via stochastic gradient descent (SGD) methods.¹⁹

2.2.3 Density estimation and cell counting

During the cell counting phase of the framework (Figure 1), the number of cells in a to-be-tested image $X^t \in \mathbb{R}^{M' \times N'}$ can be estimated by use of the trained DRM represented by $F(X; \Theta^*)$:

$$\hat{N}_c = \sum_{i=1}^{M'} \sum_{j=1}^{N'} [F(X^t; \Theta^*)]_{i,j}, \quad (7)$$

where $[F(X^t; \Theta^*)]_{i,j}$ is the estimated density at pixel (i, j) . In this step, the dimensions of the to-be-tested image can be different because arbitrary input image sizes are allowed by the trained PriCNN.

3. EXPERIMENTAL RESULTS

3.1 Dataset

A set of 49 immunofluorescent images of human embryonic stem cells (hESC) was employed in this study to test the performance of the proposed method. Each image was 512×512 pixels, and the 49 hESC images were manually annotated by identifying the centroid of each cell within each image. Statistically, the cell number among these images is about 518 ± 316 .

For each annotated image in the training dataset, the corresponding ground truth density maps were generated by placing a normalized 2D discrete Gaussian kernel with isotropic covariance, σ^2 , at each annotated cell centroid in the image (details shown in Section 2.2). The values of σ and $2K_G + 1$ were set to 3 pixels and 21 pixels, respectively. A pair consisting of an image and a density map was considered as a training sample for training the PriCNN and AuxCNNs. In this study, 5-fold cross validation was employed to evaluate the cell counting performance.

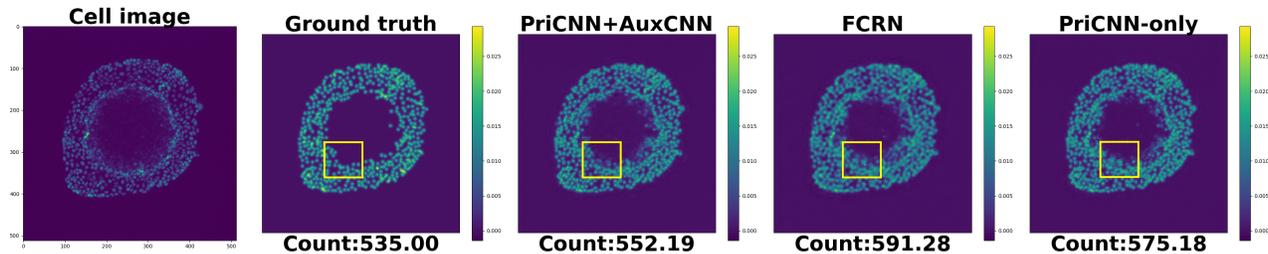


Figure 3: Density estimation for one hESC image example.

3.2 Method Implementation

We compared the performance of the proposed method (denoted as PriCNN+AuxCNN) with a state-of-the-art method, FCRN.⁷ The FCRN used the same network architecture as the PriCNN, but without concatenation layers. In addition, a PriCNN without AuxCNNs (or a FCRN with concatenation layers) was also compared to illustrate the performance improvement by use of AuxCNNs.

In the PriCNN, the convolution kernel size in the first 7 blocks was set to 3×3 , while that in the last block was set to 1×1 . The numbers of kernels in the first to 8th CONV layers are set to 32, 64, 128, 512, 128, 64, 32, and 1, respectively. The pooling size in each pool layer was set to 2×2 , and the Up layers performed bi-linear interpolation.

In the first block of the AuxCNN, the kernel size was set to 3×3 and the number of kernels was 32, while the comparable values in the second block were 1×1 and 1, respectively. In addition, the ground truth low-resolution density map (GTLR) $Y_i^k \in \mathbb{R}^{M_k \times N_k}$ was generated from the original ground-truth density map $Y_i \in \mathbb{R}^{M \times N}$ by summing local regions with size of $(a_1, b_1) = (8, 8)$, $(a_2, b_2) = (4, 4)$, and $(a_3, b_3) = (2, 2)$, respectively.

All the three methods, PriCNN+AuxCNN, PriCNN-only, and FCRN were trained under the same hyperparameter configurations, including a learning rate of 0.0001 and a batch size of 100. In addition, all the parameters were orthogonally initialized.²⁰

3.3 Results

In this study, mean absolute error (MAE) and standard deviation of absolute errors (STD) were employed to evaluate the cell counting performance. MAE measures the mean of the absolute errors (MAE) between the estimated cell counts and their ground truths for all images in the validation set. The STD measures the standard deviation of the absolute errors. Table 1 shows that the proposed method yields superior cell counting performance to the other two methods in terms of MAE and STD. In addition, Figure 3 presents the estimated density maps of one hESC image example estimated by the three methods. The numbers of cells counted from density maps are indicated below each density map. From the figure, we can see that the proposed method (PriCNN+AuxCNN) can estimate a density map that is more similar to the ground truth, compared to the other two methods. Also, our estimated cell count is closer to the ground truth.

Table 1: Performance of the proposed cell counting method

Performance	PriCNN+AuxCNN	PriCNN-only	FCRN ⁷
MAE	32.89	42.17	44.90
STD	26.35	30.97	35.30

4. DISCUSSION

Convolutional neural networks (CNN) have succeeded in computer vision tasks, including image classification,²¹ segmentation,^{14,22} and object detection.²³ The success is because that CNNs can integrate informative feature extraction and powerful nonlinear function learning. Furthermore, fully convolutional neural networks (FCNN),

such as FCN²⁴ and U-Net,¹⁴ allow flexible input image sizes, and have been employed to perform an efficient end-to-end mapping from an image (one domain) to a probability map (another domain). Both the PriCNN (the proposed DRM in the study) and FCRN⁷ are some specific FCNNs, which explain the descent cell counting accuracy they have achieved in this study.

In this study, only a set of experimental immunofluorescent hESC images was employed for DRM training and validation. However, the generalization of the proposed method should not be limited to only the experimental immunofluorescent hESC images. In future, we will evaluate the proposed method on image sets of other modalities. In addition, other competing general object counting methods will also be compared with our proposed method.

5. CONCLUSION

In this study, for the first time, a deeply-supervised density regression framework is proposed for automatic cell counting. The results obtained on experimental hESC images demonstrate the superior cell counting performance of the proposed method, compared with the state of the art.

ACKNOWLEDGMENTS

This work was supported in part by award NIH R01EB020604, R01EB023045, R01NS102213, and R21CA223799.

REFERENCES

- [1] Coates, A. S., Winer, E. P., Goldhirsch, A., Gelber, R. D., Gnant, M., Piccart-Gebhart, M., Thürlimann, B., Senn, H.-J., Members, P., André, F., et al., “Tailoring therapiesimproving the management of early breast cancer: St gallen international expert consensus on the primary therapy of early breast cancer 2015,” *Annals of oncology* **26**(8), 1533–1546 (2015).
- [2] Matas, J., Chum, O., Urban, M., and Pajdla, T., “Robust wide-baseline stereo from maximally stable extremal regions,” *Image and vision computing* **22**(10), 761–767 (2004).
- [3] Barinova, O., Lempitsky, V., and Kholi, P., “On detection of multiple object instances using hough transforms,” *IEEE Transactions on Pattern Analysis and Machine Intelligence* **34**(9), 1773–1784 (2012).
- [4] Arteta, C., Lempitsky, V., Noble, J. A., and Zisserman, A., “Learning to detect cells using non-overlapping extremal regions,” in [*International Conference on Medical Image Computing and Computer-Assisted Intervention*], 348–356, Springer (2012).
- [5] Xing, F., Su, H., Neltner, J., and Yang, L., “Automatic ki-67 counting using robust cell detection and online dictionary learning,” *IEEE Transactions on Biomedical Engineering* **61**(3), 859–870 (2014).
- [6] Lempitsky, V. and Zisserman, A., “Learning to count objects in images,” in [*Advances in neural information processing systems*], 1324–1332 (2010).
- [7] Xie, W., Noble, J. A., and Zisserman, A., “Microscopy cell counting and detection with fully convolutional regression networks,” *Computer methods in biomechanics and biomedical engineering: Imaging & Visualization* **6**(3), 283–292 (2018).
- [8] Arteta, C., Lempitsky, V., Noble, J. A., and Zisserman, A., “Detecting overlapping instances in microscopy images using extremal region trees,” *Medical image analysis* **27**, 3–16 (2016).
- [9] Cireşan, D. C., Giusti, A., Gambardella, L. M., and Schmidhuber, J., “Mitosis detection in breast cancer histology images with deep neural networks,” in [*International Conference on Medical Image Computing and Computer-assisted Intervention*], 411–418, Springer (2013).
- [10] Liu, F. and Yang, L., “A novel cell detection method using deep convolutional neural network and maximum-weight independent set,” in [*Deep Learning and Convolutional Neural Networks for Medical Image Computing*], 63–72, Springer (2017).
- [11] Lee, C.-Y., Xie, S., Gallagher, P., Zhang, Z., and Tu, Z., “Deeply-supervised nets,” in [*Artificial Intelligence and Statistics*], 562–570 (2015).
- [12] Zeng, G., Yang, X., Li, J., Yu, L., Heng, P.-A., and Zheng, G., “3d u-net with multi-level deep supervision: fully automatic segmentation of proximal femur in 3d mr images,” in [*International Workshop on Machine Learning in Medical Imaging*], 274–282, Springer (2017).

- [13] Dou, Q., Yu, L., Chen, H., Jin, Y., Yang, X., Qin, J., and Heng, P.-A., “3d deeply supervised network for automated segmentation of volumetric medical images,” *Medical image analysis* **41**, 40–54 (2017).
- [14] Ronneberger, O., Fischer, P., and Brox, T., “U-net: Convolutional networks for biomedical image segmentation,” in [*International Conference on Medical image computing and computer-assisted intervention*], 234–241, Springer (2015).
- [15] Dong, H., Yang, G., Liu, F., Mo, Y., and Guo, Y., “Automatic brain tumor detection and segmentation using u-net based fully convolutional networks,” in [*Annual Conference on Medical Image Understanding and Analysis*], 506–517, Springer (2017).
- [16] Vedaldi, A. and Fulkerson, B., “Vlfeat: An open and portable library of computer vision algorithms,” in [*Proceedings of the 18th ACM international conference on Multimedia*], 1469–1472, ACM (2010).
- [17] Fiaschi, L., Köthe, U., Nair, R., and Hamprecht, F. A., “Learning to count with regression forest and structured labels,” in [*Pattern Recognition (ICPR), 2012 21st International Conference on*], 2685–2688, IEEE (2012).
- [18] Sommer, C., Straehle, C. N., Koethe, U., Hamprecht, F. A., et al., “Ilastik: Interactive learning and segmentation toolkit.,” in [*ISBI*], **2**(5), 8 (2011).
- [19] Bottou, L., “Large-scale machine learning with stochastic gradient descent,” in [*Proceedings of COMP-STAT’2010*], 177–186, Springer (2010).
- [20] Saxe, A. M., McClelland, J. L., and Ganguli, S., “Exact solutions to the nonlinear dynamics of learning in deep linear neural networks,” *arXiv preprint arXiv:1312.6120* (2013).
- [21] He, K., Zhang, X., Ren, S., and Sun, J., “Deep residual learning for image recognition,” in [*Proceedings of the IEEE conference on computer vision and pattern recognition*], 770–778 (2016).
- [22] He, S., Zheng, J., Maehara, A., Mintz, G., Tang, D., Anastasio, M., and Li, H., “Convolutional neural network based automatic plaque characterization for intracoronary optical coherence tomography images,” in [*Medical Imaging 2018: Image Processing*], **10574**, 1057432, International Society for Optics and Photonics (2018).
- [23] Ren, S., He, K., Girshick, R., and Sun, J., “Faster r-cnn: Towards real-time object detection with region proposal networks,” in [*Advances in neural information processing systems*], 91–99 (2015).
- [24] Long, J., Shelhamer, E., and Darrell, T., “Fully convolutional networks for semantic segmentation,” in [*Proceedings of the IEEE conference on computer vision and pattern recognition*], 3431–3440 (2015).