

Deeply-Supervised Density Regression for Automatic Cell Counting in Microscopy Images

Shenghua He^a, Kyaw Thu Minn^{b,c}, Lilianna Solnica-Krezel^{c,d}, Mark A. Anastasio^{e,*}, and Hua Li^{e,f,g,*}

^a*Department of Computer Science and Engineering,
Washington University in St. Louis, St. Louis, MO 63110 USA*

^b*Department of Biomedical Engineering,
Washington University in St. Louis, St. Louis, MO 63110 USA*

^c*Department of Developmental Biology,
Washington University School of Medicine in St. Louis, St. Louis, MO 63110 USA*

^d*Center of Regenerative Medicine,
Washington University School of Medicine in St. Louis, St. Louis, MO 63110 USA*

^e*Department of Bioengineering,
University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA*

^f*Cancer Center at Illinois,*

University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA

^g*Carle Cancer Center, Carle Foundation Hospital, Urbana, IL 61801 USA*

**Corresponding Author*

Abstract

Accurately counting the number of cells in microscopy images is required in many medical diagnosis and biological studies. This task is tedious, time-consuming, and prone to subjective errors. However, designing automatic counting methods remains challenging due to low image contrast, complex background, large variance in cell shapes and counts, and significant cell occlusions in two-dimensional microscopy images. In this study, we proposed a new density regression-based method for automatically counting cells in microscopy images. The proposed method processes two innovations compared to other state-of-the-art density regression-based methods. First, the density regression model (DRM) is designed as a concatenated fully convolutional regression network (C-FCRN) to employ multi-scale image features for the estimation of cell density maps from given images. Second, auxiliary convolutional neural networks (AuxCNNs) are employed to assist in the training of intermediate layers of the designed C-FCRN to improve the DRM performance on unseen datasets. Experimental studies evaluated on four datasets demonstrate the superior performance of the proposed method.

Keywords: Automatic cell counting, Microscopy images, Fully convolutional neural network, Deeply-supervised learning

1. Introduction

Numerous microscopy image analysis methods have been proposed for various medical diagnoses and biological studies that include counting the number of cells [1, 2, 3], locating cell positions [4, 5, 6], acquiring cell shapes [7, 8, 9, 10], and classifying cell categories [11, 12, 13]. Especially, the number of cells in a microscopy image can indicate the presence of diseases [14], help differentiate tumor types [15], assist in understanding cellular and molecular genetic mechanisms [16, 17], and provide useful information to many other applications [18, 19]. Manually counting cells in microscopy images is tedious, time-consuming, prone to subjective errors, and not feasible for high-throughput process in real-world biomedical applications. During the past decades, many automatic cell counting methods have been proposed [20, 21, 22, 23, 24, 25, 26]. However, designing efficient automatic methods with sufficient counting accuracy still remains a challenging task due to various image acquisition techniques, low image contrast, complex tissue background, large variations in cell sizes, shapes and counts, and significant inter-cell occlusions in two-dimensional (2D) microscopy images.

The reported automatic cell counting methods can be categorized into *detection-based* and *regression-based* methods. Generally, detection-based methods first determine the cell centroid locations and subsequently count them to estimate the number of cells [24, 25, 23, 27]. Therefore, the performance of these methods highly relies on the accuracy of cell centroid detection results. Traditional detection-based methods have been designed based on feature extraction [28], morphological processing [29], H-minima/maxima transform [29], Laplacian of Gaussian filtering [30], maximally stable extremal region detection [24], radial symmetry-based voting [31], or conventional supervised learning strategies [4]. Recently, deep learning strategies have shown superior ability of extracting informative image features and generating inferences in all kinds of medical image analysis tasks [32, 25, 33]. A bunch of deep learning-based detection methods have been proposed [27, 34, 35, 10, 9, 36, 37, 6, 5]. For example, Falk et al. [5] trained a fully convolutional neural network (U-Net) to compute a probability map of cell existing in a given image. The number of cells can then be determined by searching for the local maxima on the probability map with a non-maxima suppression method. Xie et al. [36] applied the non-maxima suppression process to a dense proximity map for cell detection. The proximity map was produced by a fully residual convolutional network-based structural regression model (StructRegNet), and exhibits higher responses at locations near cell centroids to benefit for local maximum searching. Tofighi et al. [35] used a *a priori*-guided deep neural network for cell nuclei detection. In their method, nuclei shape *a priori* is employed as a regularizer in a model learning process to improve the cell detection accuracy. Liu et al. [27] trained a CNN model to determine the final cell detection result from the results generated by several traditional cell counting methods. The selection process was formulated as a maximum-weight independent set (MWIS) problem, a combinatorial optimization problem that has been studied in many applications of clustering,

segmentation, and tracking. Paulauskaite et al. [38] recently performed an experimental investigation of the Mask R-CNN method, which was proposed by He et al. [39], to detect overlapping cells with a two-stage procedure of determining potential cell regions and jointly classifying and predicating cell masks. The method was validated on fluorescence and histology images and showed promising results on detecting overlapping cells. However, it still remains difficult to detect cells that are highly occluded, densely concentrated, and surrounded by histopathological structures.

Compared to detection-based methods, regression-based cell counting methods have received more and more attention due to their superior performance on counting occluded cells [40, 41, 42, 1, 43, 2, 44, 3, 45, 46]. Some regression-based methods learn a cell counter through a regression process directly without requiring cell detection. In these methods, the number of cells is the direct and only output, and no cell location information can be provided. For example, Khan et al. [40] and Xue et al. [41] learned a convolutional neural network-based cell counter from small image patches which can increase the amount of training samples. The total number of cells across the whole image can then be obtained by summing those on image patches. These methods might suffer from redundant estimation issues across the patch boundaries, and might not be efficient since they have to infer for each image patch separately before cell counting. Differently, Cohen et al. [42] learned a cell counter with a fully convolutional neural network (FCNN). They utilized the “sliding window” mechanism associated with the convolutional layers of the FCNN to address the redundant counting issues across the overlapped regions among image patches. Their method counts the number of cells by directly inferring a count map for the whole image. The method performance might be affected by the sizes of sliding windows.

Other regression-based methods learn a spatial cell density regression model (DRM) across a full-size image instead of learning direct cell counters [1, 43, 3, 47]. In these methods, the number of cells can be obtained by integrating the regressed density map, and the local maxima in the density map can be considered as cell centroid locations. Therefore, both the number and the centroid locations of cells can be obtained. Conventional density regression-based methods learn DRMs from extracted handcrafted image features, in which the feature extraction is independent of the DRM learning. For example, Lempit-sky et al. [1] used local features (e.g. scale-invariant feature transform (SIFT) features) to learn a linear DRM by use of a regularized risk regression-based learning framework. Differently, Fiaschi et al. [43] learned a nonlinear DRM based on regression random forest methods. In their method, image features computed by ordinary filter banks were employed as the model input. The performance of these methods relies on the effectiveness of feature extraction methods, that of the DRM learning algorithms, and the match between them.

Instead of using handcrafted image features to learn a DRM, some methods were proposed to integrate the feature learning into end-to-end nonlinear DRM learning by use of deep convolutional neural networks. The learned end-to-end DRMs use images as their direct inputs to compute the corresponding density

maps [48, 3, 47, 49]. As one of the pioneering work using this strategy, Xie et al. [3] proposed a fully convolutional regression network (FCRN) to learn such a DRM integrating image feature extraction and density map estimation for arbitrary-sized input images. By use of CNNs in feature extraction and model learning, their method demonstrated superior cell counting performance than conventional density regression-based methods, especially on microscopy images containing severely overlapped cell regions. Following Xie et al.’s work, Zheng et al. [49] trained a FCRN by incorporating a manifold regularization based on the graph Laplacian of the estimated density maps to reduce the risk of overfitting. Liu et al. [50] employed a post-processing CNN to further regress the estimated density map to improve the accuracy of cell counting.

However, in the original FCRN work, the network layers of a FCRN are structured hierarchically and the output of each layer relies merely on the output of its direct adjacent layer. This restricts the FCRN to produce a more authentic density map for cell counting. In addition, the training of original FCRN is based on a single loss that is measured at the final output layer, and all its intermediate layers are optimized based on the gradients back-propagated from this single loss only. The decreased gradients potentially trap the optimization of intermediate layers into unsatisfying local minima and jeopardize the overall network performance.

Recently, CNNs that concatenate multi-scale features by shortcut connections of non-adjacent layers have been reported and demonstrated promising performance than conventional hierarchical networks for many applications [51, 52]. In these concatenated network architectures, the multi-scale image features extracted by all the layers along the down-sampling path can be integrated into the input of the layers along the up-sampling path to further improve the model performance. Also, deeply-supervised (or deep supervision) learning strategies, aiming at enhancing the training of intermediate layers of designed neural networks by providing direct supervisions for them, have been proposed and have yielded promising performance for several computer vision tasks including image classification [53] and segmentation [54, 55]. To the best of our knowledge, deeply-supervised learning has not been employed in learning a density regression model for cell counting task except our preliminary work [56].

In this study, a novel density regression-based method for automatically counting cells in microscopy images is proposed. It addresses the two shortcomings that exist in the original FCRN by integrating the concatenation design and deeply-supervised learning strategy into the FCRN. Specifically, the density regression model (DRM) is designed as a concatenated FCRN (C-FCRN) to employ multi-scale image features for the estimation of cell density maps from given images. The C-FCRN can fuse multi-scale features and improve the granularity of the extracted features to benefit the density map regression. It also facilitates the learning of intermediate layers in the down-sampling path by back-propagating the gradients conveyed via the shortcut connections. In addition, auxiliary convolutional neural networks (AuxCNNs) were employed to assist in training the C-FCRN by providing direct and deep supervision on learning its intermediate layers to improve the cell counting performance.

The remainder of the manuscript is organized as follows. The proposed automatic cell counting method is described in Section 2. Section 3 describes the testing datasets and the implementation details of the proposed method. Section 4 contains the experimental results. A discussion and conclusion are provided in Section 5 and Section 6, respectively.

2. The Proposed Cell Counting Method

2.1. Background: Density regression-based cell counting

The salient mathematical aspects of the density regression-based counting process can be described as below. For a given two-dimensional microscopy image $X \in \mathbb{R}^{M \times N}$ that includes N_c cells, the density map corresponding to X can be represented as $Y \in \mathbb{R}^{M \times N}$. Each value in Y represents the number of cells at the corresponding pixel of X . Let $\phi(X)$ be a feature map extracted from X , a density regression function $F_\phi(\phi(X), \Theta)$ can be defined as a mapping function from X to Y :

$$Y = F_\phi(\phi(X); \Theta), \quad (1)$$

where the vector Θ parameterizes F_ϕ . The number of cells in X can be subsequently computed by:

$$N_c = \sum_{i=1}^M \sum_{j=1}^N Y_{i,j} = \sum_{i=1}^M \sum_{j=1}^N [F_\phi(\phi(X); \Theta)]_{i,j}, \quad (2)$$

where $[F_\phi(\phi(X); \Theta)]_{i,j}$ is the computed density associated with the pixel $X_{i,j}$. The key component of density regression-based methods is to learn $F_\phi(\phi(X), \Theta)$ from $\phi(X)$ and the corresponding Θ [1, 43]. In the fully convolutional regression network (FCRN) [3], $F_\phi(\phi(X), \Theta)$ can be simplified to $F(X, \Theta)$ because it can be learned directly from X .

2.2. Concatenated FCRN-based cell counting method

The proposed concatenated FCRN (C-FCRN) is shown in Figure 1, which integrates a concatenated neural network design and deeply-supervised learning strategy into the original FCRN. The C-FCRN network includes 8 blocks. Three concatenation layers (red lines in Figure 1) are established to connect the intermediate outputs along the down-sampling path to the input of the fifth to seventh blocks along the up-sampling path, respectively. This C-FCRN design integrates multi-scale features from non-adjacent layers to improve the granularity of the extracted features for density map regression, and subsequently improve the model performance on cell counting. The first three blocks in the C-FCRN are employed to extract low-dimension feature maps. Each of them includes a convolutional (CONV) layer, a ReLU layer, and a max-pooling (Pool) layer. The fourth block, including a CONV layer and a ReLU layer, is used to further extract highly-representative features. The fifth to seventh blocks are employed to gradually restore the resolutions of feature maps while refining the

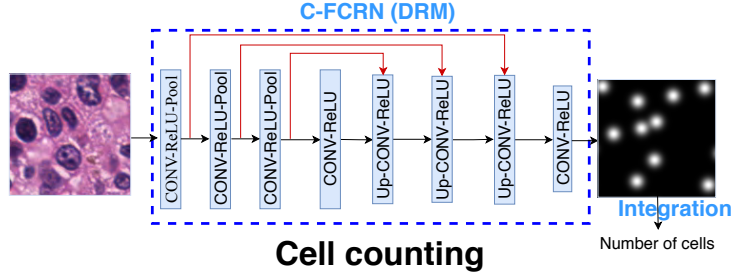


Figure 1: Framework of the proposed C-FCRN based automatic cell counting method. Different from the original FCRN, three shortcut connections (red lines) are established to connect the first three intermediate blocks to the fifth to seventh blocks, respectively.

extracted feature maps. Each of these blocks includes an up-sampling (UP) layer, a CONV layer, and a ReLU layer. The last block, including a chain of a CONV layer and a ReLU layer, is employed to estimate the final density map.

In C-FCRN, the CONV layer in each block is associated with a set of learnable kernels and is employed to extract local features from the output of its previous layer. The ReLU layer in each block is employed to increase the nonlinear properties of the network without affecting the receptive fields of the CONV layer by setting negative responses from its previous layer to zero while keeping the positive ones unchanged. Each Pool layer in the first three blocks performs a down-sampling operation on an input feature map by outputting only the maximum value in every down-sampled region in the feature map. Therefore, multi-scale informative features are extracted progressively along with the decrease of the spatial size of an input feature map. In contrast, each Up layer in the fifth to seventh block performs an up-sampling operation to gradually restore the resolution of the final estimated density map. This network design permits integration of feature extraction into the density regression process. Therefore, no additional feature extraction methods are required.

Given a to-be-tested image $X \in \mathbb{R}^{M \times N}$ and the trained density regression function $F(X; \Theta)$, the density map corresponding to X can be estimated as $\hat{Y} = F(X; \Theta)$. Therefore, the number of cells in X can be conveniently estimated based on the equation below:

$$\hat{N}_c = \sum_{i=1}^M \sum_{j=1}^N \hat{Y}_{i,j} = \sum_{i=1}^M \sum_{j=1}^N [F(X; \Theta)]_{i,j}, \quad (3)$$

where $[F(X; \Theta)]_{i,j}$ represents the estimated density of pixel (i, j) in the X .

2.3. Deeply-supervised C-FCRN training with auxiliary CNNs

The task of training the C-FCRN corresponds to learning a nonlinear density regression function $F(X, \Theta)$ with parameters Θ . However, training such a hierarchical and concatenated deep neural network by solving the corresponding highly non-convex optimization problem is a challenging task. Motivated by

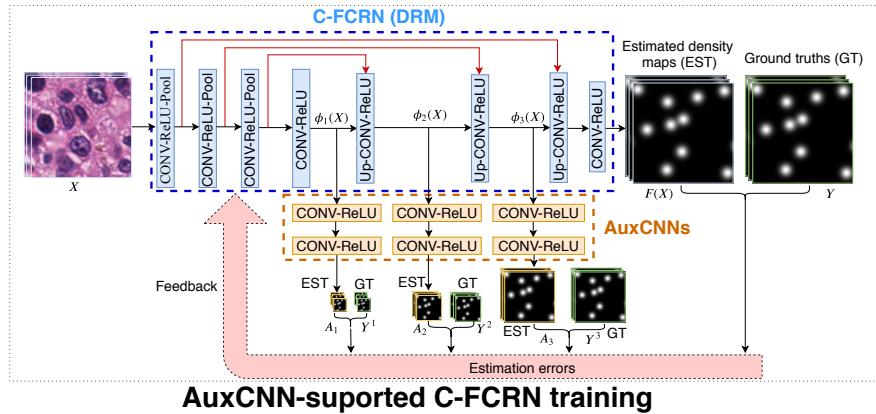


Figure 2: Framework of the AuxCNN-supported C-FCRN training process. The blue dash-line region indicates the C-FCRN. The orange dash-line region indicates the three AuxCNNs. EST and GT represents the estimated and ground truth density maps with varied resolutions, respectively.

the deeply-supervised learning strategies [53, 54, 55], we employed three auxiliary convolutional neural networks (AuxCNNs) to provide direct supervision for learning the intermediate layers of the C-FCRN. The AuxCNN-supported C-FCRN training process is shown in Figure 2. Each AuxCNN contains two CONV-ReLU blocks, which estimate a low-resolution density map from each input feature map, respectively. The difference between the estimated density maps and the related ground truth are employed to support the C-FCRN training.

The Θ in the density regression function $F(X, \Theta)$ can be re-defined as $\Theta = (\Theta_1, \Theta_2, \Theta_3, \Theta_4)$, in which Θ_1 represents the trainable parameters in the first four blocks, Θ_2 represents the parameters in the 5-th block, Θ_3 represents the parameters in the 6-th block, and Θ_4 represents the parameters in the last 7-th and 8-th blocks, respectively. The outputs of the 4-th, 5-th, and 6-th blocks can then be denoted as $\phi_1(X; \Theta_1)$, $\phi_2(X; \Theta_1, \Theta_2)$, and $\phi_3(X; \Theta_1, \Theta_2, \Theta_3)$. They are also the inputs of the 1-st, 2-nd, and 3-rd AuxCNNs, respectively. Given each input ϕ_k ($k = 1, 2, 3$), the output of each AuxCNN is a low-resolution density map $A_k(\phi_k; \theta_k)$, where θ_k represents the parameter vector of the k -th AuxCNN.

$F(X; \Theta)$ and $A_k(\phi_k; \theta_k)$ are jointly trained through the minimization of a combined loss function [53],

$$\begin{aligned}
 L_{cmb}(\Theta, \theta_1, \theta_2, \theta_3) = & L(\Theta) + \sum_{k=1}^3 \alpha_k L_k(\Theta_1, \dots, \Theta_k, \theta_k) \\
 & + \lambda(\|\Theta\|^2 + \sum_{k=1}^3 \|\theta_k\|^2), \quad k = 1, 2, 3,
 \end{aligned} \tag{4}$$

where $L(\Theta)$ represents a loss function that measures the average mean square

errors (MSE) between the estimated density map from the C-FCRN and the corresponding ground truth density map. $L_k(\Theta_1, \dots, \Theta_k, \theta_k)$ represents a loss function that measures the average MSE between a low-resolution density map estimated by the k -th AuxCNN and the corresponding low-resolution ground-truth (LRGT) density map. The parameter $\alpha_k \in [0, 1]$ controls the supervision strength under the k -th AuxCNN. The parameter λ controls the strength of l_2 penalty to reduce overfitting and $L_k(\Theta_1, \dots, \Theta_k, \theta_k)$ ($k = 1, 2, 3$) and $L(\Theta)$ are defined as:

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

where Y_b represents the full-size ground truth density map of the b -th training data X_b of B training images. Here, Y_b^k represents the low-resolution ground-truth (LRGT) density map, which is generated from Y_b by summing local regions in the original ground truth density map. An example of the summing process is shown in Figure 3.

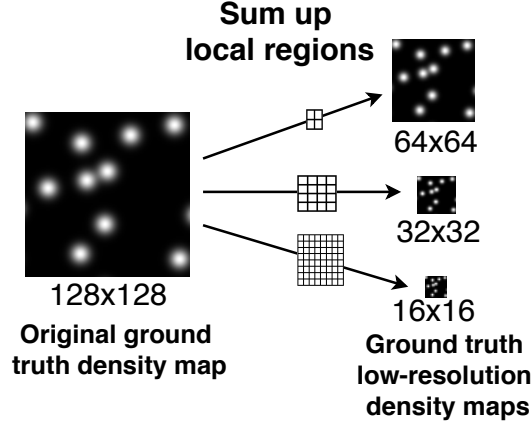


Figure 3: Example of constructing ground truth low-resolution density maps from an original ground truth of 128×128 pixels by summing up every local regions with size 2×2 , 4×4 and 8×8 pixels, respectively.

The loss L_{cmb} can be numerically minimized via momentum stochastic gradient descent (SGD) methods [57] based on the Eqn. (6) shown below:

$$\begin{cases} \Delta \Theta_k^{(t+1)} = \beta \Delta \Theta_k^{(t)} - (1 - \beta) \left(\eta \frac{\partial L_{cmb}^{(t)}}{\partial \Theta_k^{(t)}} \right), \\ \Theta_k^{(t+1)} = \Theta_k^{(t)} - \Delta \Theta_k^{(t+1)}, \end{cases} \quad (6)$$

where $\Theta_k^{(t)}$ is the updated parameters Θ_k at the t -th iteration; β is a momentum parameter that controls the contribution of the result from the previous iteration; and η is a learning rate that determines the parameter updating speed.

Since $L_k(\Theta_1, \dots, \Theta_k, \theta_k)$ only relates to θ_k and Θ_m ($m = 1, 2, \dots, k$), the gradient w.r.t the model parameters Θ_k can be computed by:

$$\frac{\partial L_{cmb}^{(t)}}{\partial \Theta_k^{(t)}} = \frac{\partial L^{(t)}}{\partial \Theta_k^{(t)}} + \sum_{m=k}^3 \alpha_m \frac{\partial L_m^{(t)}}{\partial \Theta_k^{(t)}} + 2\lambda \Theta_k^{(t)}, \quad (7)$$

with the back-propagation algorithm [58]. The learned $F(X; \Theta)$, represented by the trained C-FRCN model, can be used to estimate density maps for arbitrary-sized images because fully convolutional layers are employed in the C-FRCN.

In the rest of this paper, the proposed C-FRCN deeply-supervised by auxiliary CNNs during the training process is denoted as **C-FRCN+Aux**.

3. Datasets and method implementation

3.1. Datasets

Four microscopy image datasets were considered in this study, which are synthetic images of bacterial cells, experimental images of bone marrow cells, colorectal cancer cells, and human embryonic stem cells (hESCs), respectively. Table 1 illustrates the data details. Sample images from the four datasets are shown in Figure 4.

Table 1: Four datasets employed in this study

Dataset	Bacterial cells	Bone marrow cells	Colorectal cancer cells	hESCs
# of images	200	40	100	49
Image size	256×256	600×600	500×500	512×512
Count statistics	174 ± 64	126 ± 33	310 ± 216	518 ± 316

Image size is represented by pixel, and count statistics is represented by mean and standard variations of cell numbers in all the images in each dataset.

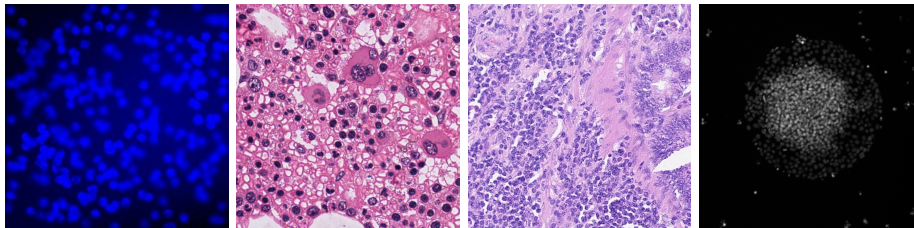


Figure 4: Example images of the four datasets used in this study. From left to right: Synthetic bacterial cells, Bone marrow cells, Colorectal cancer cells, and Human embryonic stem cells.

3.1.1. Synthetic bacterial cells

This is a public synthetic dataset generated by Lempitsky et al. [1] by use of the method proposed by Lehmussola et al. [59]. This dataset contains 200 RGB synthetic fluorescent microscopy images of bacterial cells. The size of each image is $256 \times 256 \times 3$ pixels. The cells in these images are designed to be clustered and occluded with each other. This dataset is appropriate for testing the performance of the proposed method.

3.1.2. Bone marrow cells

This dataset includes 40 Hematoxylin-Eosin (H&E) stained bright-field RGB microscopy images, which were created from 10 images acquired from the human bone marrow tissues of 8 different patients [60]. The original image size of each H&E image is $1200 \times 1200 \times 3$ pixels. Each of the 10 original image was split into 4 images with the size of 600×600 pixels, following the process in Ception-Count [42]. The images in this dataset have inhomogeneous tissue background, and large cell shape variance.

3.1.3. Colorectal cancer cells

This dataset includes 100 H&E stained histology RGB images of colorectal adenocarcinomas acquired from 9 patients [61]. Knowing the number of colorectal adenocarcinomas can help with better understanding of colorectal cancer tumor for exploring various treatment strategies. Images in this dataset yield high inhomogeneous tissue region, noisy background, and large variance in the numbers of cells. This dataset is suitable to test the robustness and accuracy of given cell counting methods.

3.1.4. Human embryonic stem cells

This dataset contains 49 immunofluorescent images of human embryonic stem cells (hESC) that are differentiated into varied cell types [62]. The differentiation efficiency of the hESC population can be potentially observed based on the counted number of cells from each differentiation type in the images. The images in this dataset yield low image contrast and severe cell occlusion and clusters. In addition, high background noise exists in images.

3.2. Ground truth density map generation

Both the full-size and low-resolution ground truth (LRGT) density maps of the training images need to be constructed in order to train the C-FCRN and three AuxCNNs simultaneously. The full-size ground truth density map Y of an image X in the four data sets (described in Section 3.1) is defined as the superposition of a set of normalized 2D discrete Gaussian kernels [3]. The number of Gaussian kernels in Y is identical to the number of cells in X , and each kernel is centered at a cell centroid in X (as shown in Figure 5). Intuitively, the density map design can be interpreted in the perspective of microscopy imaging. Due to the limitation of imaging system and the point spread function (PSF), the intensity of each single pixel in image X is affected by the PSF, and can be considered as a combination of the PSF-affected intensities of the pixel itself and its surrounding pixels. Accordingly, the density map is generated by simulating the imaging system and setting PSF as a Gaussian function. Integrating the density over Y gives an estimate of the counts of cells in image X . This definition is also the same as the definition described in Lempitsky et al. [1], one of the compared methods in this study. This process would allow density regression-based methods to solve the problem of counting the overlapping cells. In the synthetic bacterial cell dataset, the ground truth cell centroids and numbers

were automatically annotated during the image generation [1], while they are manually annotated on images in the other three experimental datasets. The manual annotations for bone marrow cell images and colorectal cell images were provided by [60] and [61], respectively. The hESC annotation was performed by a graduate student under the supervision of and validation of a biologist expert [62].

Let $S = \{(s_{x_k}, s_{y_k}) \in \mathbb{N}^2\}$ represent N_c cell centroid positions in X , where $k = 1, 2, \dots, N_c$. Each $Y_{i,j}$ in Y can be expressed as:

$$\begin{cases} Y_{i,j} = \sum_{k=1}^{N_c} G_\sigma(i - s_{x_k}, j - s_{y_k}), \\ G_\sigma(n_x, n_y) = C \cdot e^{-\frac{n_x^2 + n_y^2}{2\sigma^2}}, n_x, n_y \in \{-K_G, \dots, 0, \dots, K_G\}, \end{cases} \quad (8)$$

where $G_\sigma(n_x, n_y) \in \mathbb{R}^{(2K_G+1) \times (2K_G+1)}$ is a normalized 2D Gaussian kernel, and $\sum_{n_x=-K_G}^{K_G} \sum_{n_y=-K_G}^{K_G} G_\sigma(n_x, n_y) = 1$. σ^2 is the isotropic covariance, K_G is an integer that determines the kernel size $(2K_G + 1) \times (2K_G + 1)$ pixels, and C is a normalization constant. In light of the different sizes of cells in these four different datasets, the parameter σ was set to 5 pixels for bone marrow images and 3 pixels for images in the other three datasets, respectively. The parameter K_G was set to 10 pixels for all four image datasets.

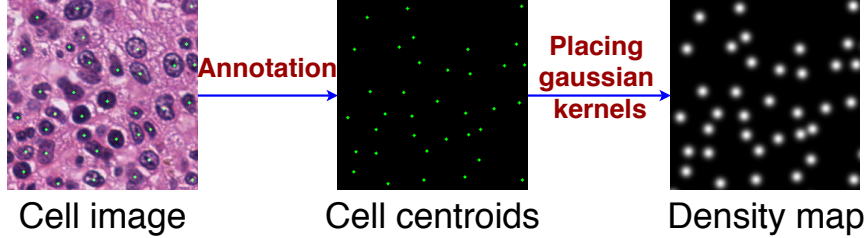


Figure 5: Example of generating density map from a given cell centroid set.

Corresponding to the bi-linear interpolation performed by the Up layers in C-FCRN, the three low-resolution ground truth (LRGT) density maps $Y^k \in \mathbb{R}^{M_k \times N_k}$ ($k = 1, 2, 3$) were generated from the original full-size ground-truth density map $Y \in \mathbb{R}^{M \times N}$ by summing local regions with size of 8×8 , 4×4 , and 2×2 pixels, respectively. Examples of ground truth of the images from the marrow bone cell dataset are shown in Figure 5, and the corresponding LRGT density map construction process is shown in Figure 3.

All images employed in the study were preprocessed by normalizing pixel values to a uniform range $[0, 1]$ in order to accelerate and stabilize the model training process [63]. The normalized images were subsequently employed as the inputs of the networks for both training and testing purpose. Random rotation with an arbitrary angle within $[0, 40^\circ]$ and random flipping on the training images was performed as a data augmentation operation to mitigate overfitting. During the training process, the ground truth density maps were amplified by 100 in order to force the C-FCRN and AuxCNNs to fit cell area

rather than the background [3]. Correspondingly, the estimated density maps estimated from the testing image were scaled back by a factor of 0.01 before counting cell numbers.

3.3. C-FCRN and AuxCNN network parameter settings

The convolution kernel size in the first 7 blocks of C-FCRN was set to 3×3 , while that in the last block was set to 1×1 . The numbers of kernels in the first to 8-th CONV layers were 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. The size of the C-FCRN input image was set to 128×128 pixels, so did the output density map. Three AuxCNNs yield the similar network structures, in which the kernel size of the first block in AuxCNN was set to 3×3 and the number of kernels was set to 32, while that in the second block were set to 1×1 and 1, respectively.

3.4. C-FCRN+Aux training and testing

Six thousand epochs were employed for model training, and that can permit the convergence of the training process in this study. In each training epoch, 100 image patches of 128×128 pixels were randomly cropped from each image for training. All the cropped image patches and their corresponding density maps were employed for training DRMs in the following epoch. The weight vector in the combined loss function $L_{cmb}(\Theta, \theta_1, \theta_2, \theta_3)$ in Eqn. 4 was set to $(\alpha_1, \alpha_2, \alpha_3) = (\frac{1}{64}, \frac{1}{16}, \frac{1}{4})$, considering that the task of a higher-resolution density estimation is more correlated to the task of original density estimation task. A momentum SGD method was used to minimize the combined loss function for jointly training the FCRN and AuxCNNs. The learning rates for training the C-FCRN+Aux were determined by operating a line search in a set of values $\{0.05, 0.01, 0.005, 0.0001, 0.0005, 0.001\}$ and selecting the one that results in the lowest validation loss. Other hyper-parameters were set to the fixed values of $\beta = 0.99$, $\lambda = 0.01$, and batch size = 100 considering the variations of these hyper-parameter values did not significantly improve the training performance based our trials. All the network parameters in the C-FCRN+Aux were orthogonally initialized [64].

The model performance was investigated by use of 5-fold cross validation on all four image datasets. When conducting cross validation on one of the four image datasets, the image dataset was randomly split into 5 folds of images for model training and validation. Specifically, every time, 4 of them were employed as the training dataset and the rest one as the validation dataset. Repeat the process for 5 times until each fold of data was used as validation dataset once. The average validation performance over the five times were measured as the evaluation result.

The proposed C-FCRN+Aux was implemented by use of python programming language with libraries including Python 3.5, NumPy 1.14, Keras 2.0, and Tensorflow 1.3.1. Model training and validation were performed on a Nvidia Titan X GPU with 12 GB of VRAM and several Intel(R) Xeon(R) CPUs with E5-2620 v4 @ 2.10GHz.

3.5. Other methods for comparison

The proposed method (denoted as C-FCRN+Aux) was compared to other eight state-of-the-art methods, which include four regression-based counting methods [3, 42, 1], and four detection-based counting methods [24, 39, 36, 5].

Those four to-be-compared regression-based counting methods include the original FCRN method [3], the C-FCRN without AuxCNNs-supporting training (denoted as C-FCRN-only), the Count-Ception [42] method, and the Lempitsky’s method [1]. The original FCRN and the C-FCRN-only methods are nonlinear density regression-based methods. The Count-Ception [42] method is a nonlinear counter regression-based method, which employs a fully convolutional neural network (FCNN) to perform redundant cell counting in each overlapped local region and average out the estimated results to obtain the final cell count. The Lempitsky’s method is a linear density regression-based method, which learns the DRM by use of a regularized risk linear regression. Its hyper-parameter settings can be found in [1].

The loss functions for training the FCRN and C-FCRN were defined as the MSE between the ground truth density maps and the estimated density maps measured in a batch of training data. Differently, the loss function in the Count-Ception method was specified as the mean absolute error between the ground truth and the estimated count maps. The ground truth count map was generated according to its definition in the literature [42]. A momentum SGD method was used to minimize the loss functions in all these three methods. The learning rates and other hyper-parameters for model training in these methods were determined in the same way as those were described in Section 3.4. All the network parameters in FCRN and C-FCRN-only were orthogonally initialized [64]; while those in the Count-Ception model were initialized by Glorot weight initialization [42]. The local region size in the Count-Ception was set to 32×32 as suggested in the literature [42].

The four referred detection-based counting methods include three deep-learning methods, StructRegNet [36], U-Net [5] and Mask R-CNN [39], and the Arteta’s method [24]. In the detection-based cell counting methods, the number of cells is determined by the number of detected cell centroids or cell regions. The StructRegNet used a fully residual convolutional network to regress a dense proximity map that exhibits higher responses at locations near cell centroids. Then the thresholding and non-maximum post-processes were employed to count the number of cell centroids. Differently, the U-Net method employed a U-Net to predict a cell probability map, and count cell centroids from it. The mask R-CNN method detects the cells by first detecting possible cell regions and then jointly predicting and segmenting these regions to get cells. The thresholds for the post-processes were tuned by visually checking detection results for two random validation images. The to-be-compared Arteta’s method [24] aims to segment an image into non-overlapped cell regions by use of a conventional machine learning technique. The results related to Arteta’s method on the bacterial dataset was referred to the literature [24].

The experiment settings related to the three deep learning detection-based counting methods are described as below. The StructRegNet model was built

up based on the instructions presented by Xie et al. [36]. The ground truth proximity map was generated by an exponential function defined as:

$$\mathcal{M}(u, v) = \begin{cases} \frac{e^{\alpha(1-\frac{D(i,j)}{d})}-1}{e^\alpha-1}, & D(i, j) \leq d, \\ 0, & D(i, j) > d, \end{cases} \quad (9)$$

where $D(i, j)$ is the Euclidean distance from a pixel (i, j) to its closest annotated cell centroid; d is a distance threshold and α is the decay ration, and both of them are used to control the shape of this exponential function. As suggested in literature [36], $\alpha = 3, d = 15$ was set in this study; the loss function for model training was a weighted MSE between the ground truth and estimated proximity map measured in a training batch. In this loss function, pixels closer to cell centroids were assigned to higher weights than those far-away pixels, and obtained more attention in the model training.

Although the task in this study is to annotate cell centroids, considering that the original U-Net method [51] requires fully annotation of complete cell masks, we reformulated the cell counting task as a segmentation problem in order to adapt the U-Net model to infer a segmentation map containing a small 2D disk at each cell centroid for each image, as suggested by Falk et al. [5]. When generating the ground truth segmentation maps, the radii of the 2D disks were set to 4 pixels, 8 pixels, 5 pixels and 3 pixels for the bacterial cell, bone marrow cell, colorectal cancer cell and hESC datasets, respectively, based on the average cell size of each dataset. The U-Net was trained by minimizing a binary cross-entropy loss with a momentum SGD method. The learning rates were determined by operating a line search in a set of values $\{0.05, 0.01, 0.005, 0.0001, 0.0005, 0.001\}$ and selecting the one that results in the lowest validation loss. Other hyper-parameters were set to the fixed values of $\beta = 0.99$, $\lambda = 0.01$, and batch size = 100. All the network parameters in the U-Net were orthogonally initialized. The same adaptation was performed for the Mask R-CNN method, except that a separate segmentation map was generated for each cell. For example, a set of N_c separate segmentation maps were prepared as the ground truth for an image containing N_c cells. ResNet-101 was chosen as feature extraction network in the Mask R-CNN model, since it yields better performance than the ResNet-50. The image scaling factor parameter was set to 2. The model was trained with image patches of $512 \times 512 \times 3$ pixels that were randomly cropped from the scaled images in the training mode, and then tested on the whole scaled images. The sizes of anchors related to the region proposal networks for the bacterial cell dataset and the bone marrow cell dataset were set to $\{8, 16, 32, 64\}$ and $\{8, 16, 32, 64, 128\}$, respectively. The Mask R-CNN model was trained by jointly minimizing the bounding box loss, classification loss, and segmentation loss. A stochastic gradient descent method was employed to minimize the losses. The batch size and learning rate were set to 4 and 0.001, respectively. The other parameter settings can be found in the repository [65].

The implementations of the six to-be-compared deep learning-based methods, including the FCRN, C-FCRN-only, Count-Ception, U-Net, Mask R-CNN, and StructRegNet, were based on the same Python, Tensorflow and Keras li-

baries as described in Secion 3.4. In addition, the buildup of Mask R-CNN model was based on an open-sourced repository [65]. A Matlab implementation of Lempitsky’s method provided by Lempitsky et al. [1] was used to evaluate the Lempitsky’s method. The results related to Arteta’s method on the bacterial dataset was directly referred to the literature [24].

3.6. Performance evaluation metrics

Mean absolute count error (MAE), mean relative count error (MRE), and their related standard deviations (denoted by STDa and STD_r) were employed as the evaluation metrics:

$$\begin{aligned}
 \text{MAE} &= \frac{1}{T} \sum_{t=1}^T |N_{c_t} - \hat{N}_{c_t}|, \\
 \text{STDa} &= \sqrt{\frac{1}{T-1} \sum_{t=1}^T (|N_{c_t} - \hat{N}_{c_t}| - \text{MAE})^2}, \\
 \text{MRE} &= \frac{1}{T} \sum_{t=1}^T \frac{|N_{c_t} - \hat{N}_{c_t}|}{N_{c_t}}, \\
 \text{STD}_r &= \sqrt{\frac{1}{T-1} \sum_{t=1}^T \left(\frac{|N_{c_t} - \hat{N}_{c_t}|}{N_{c_t}} - \text{MRE} \right)^2}.
 \end{aligned} \tag{10}$$

where T is the number of validation images, N_{c_t} and \hat{N}_{c_t} are the ground truth cell count and the estimated cell count in the t -th image respectively. MAE measures the mean of the absolute errors between the estimated cell counts and their ground truths for all the validation images. Considering the large variance in the numbers of cells in colorectal images and hESC images, MRE was also considered for method evaluation because they measure the relative errors between the ground-truth counts and the estimated counts. STDa and STD_r indicate the stability of the cell counting process. A lower MAE or MRE indicates a better cell counting accuracy, and a lower STDa or STD_r means a more stable counting performance.

4. Experimental results

4.1. Cell counting performance

Cell counting performance of the proposed “C-FCRN+Aux” method and the other eight methods on the four datasets are reported in the Figure 6 and Table 2. The proposed method demonstrates superior cell counting performance to the other eight methods in terms of MAE and MRE. Compared to the regression-based methods, all four detection-based methods achieve worse counting performance in terms of MAE and MRE. Also, all three non-linear density regression-based methods (the proposed method, FCRN, C-FCRN-only) demonstrate superior counting performance compared to Lempitsky’s method, one of the conventional linear methods.

A paired t -test was performed on the absolute counting errors related to the proposed method (C-FCRN+Aux) and its closest counterpart C-FCRN-only. In this test, the null hypothesis H_0 was defined as the population means of absolute errors related to the C-FCRN+Aux is higher than that of C-FCRN, and vise

Table 2: MAE±STD performance evaluated on the four data sets.

MAE ± STD	Bacterial cells	Bone marrow cells	Colorectal cancer cells	hESC
Lempitsky’s method	3.52 ± 2.99	—	—	—
Altera’s method	5.06*	—	—	—
Mask R-CNN	36.92 ± 19.73	44.4 ± 14.17	—	—
U-Net	27.77 ± 25.48	48.00 ± 18.98	—	—
StructRegNet	9.80 ± 8.68	12.75 ± 8.62	45.97 ± 47.97	189.14 ± 231.86
FCRN	2.75 ± 2.47	8.46 ± 7.63	42.58 ± 33.51	44.90 ± 35.39
C-FCRN-Only	2.58 ± 2.28	8.68 ± 7.37	39.55 ± 35.80	42.17 ± 30.97
Count-ception	2.79 ± 2.68	7.89 ± 6.83	34.14 ± 29.04	35.87 ± 35.77
C-FCRN+Aux	2.37 ± 2.27	6.55 ± 5.26	29.34 ± 25.4	32.89 ± 26.35

* indicates the result reported in the literature [24], in which the method was tested on a set of 100 testing bacterial cell images. Differently, the results from other methods related to this dataset were evaluated on a complete set of 200 bacterial cell images in this study, since the cross validation-based evaluation allows each image to be considered as a testing image for once. In addition, the Lempitsky’s method was only validated on the bacterial cell dataset because this dataset provides handcrafted image features for validation purpose. The results from the U-Net and Mask R-CNN were not reported on colorectal cancer cell and hESC datasets, due to their failure in providing reasonable detection results on the two datasets.

versa for hypothesis H_1 . The p -values for the tests on the synthetic cell, bone marrow cell, colorectal cancer cell, and hESC datasets are 6.19×10^{-4} , 0.042, 5×10^{-7} and 2.8×10^{-3} , respectively. A similar paired t -test was performed on the absolute counting errors related to C-FCRN+Aux and original FCRN, and the corresponding p -values related to the four datasets are 0.024, 0.012, 7.35×10^{-5} and 0.017, respectively. The paired t -test results show that the MAEs related to the proposed method were lower than its two counterparts: C-FCRN and FCRN-only with statistical significance.

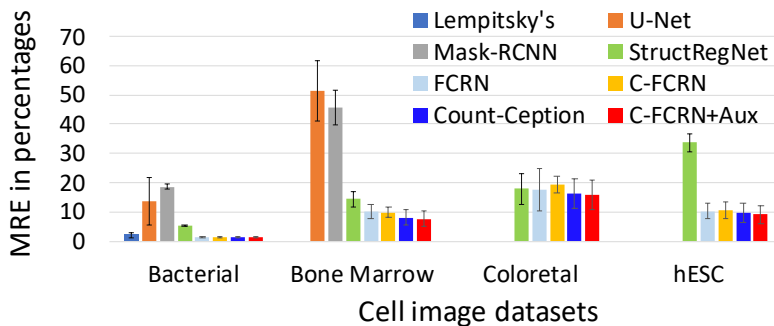
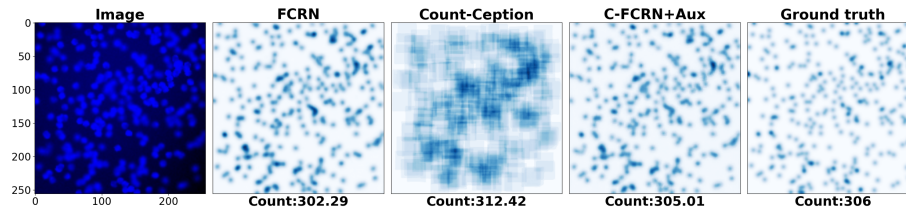
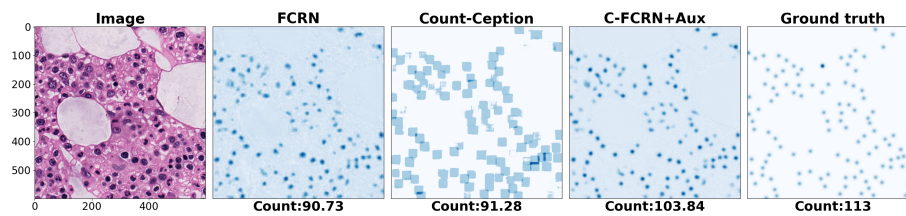


Figure 6: The MRE performance evaluated on four different datasets. “C-FCRN+Aux” represents the proposed method in this study. No MRE results were reported for Arteta’s method.

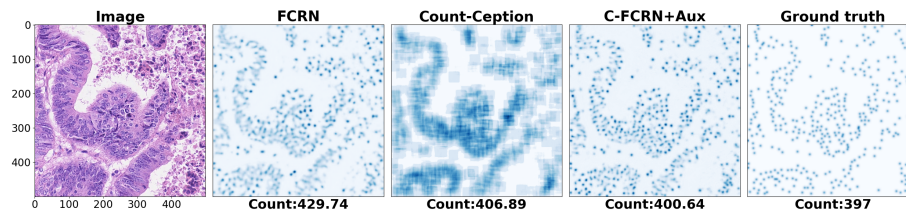
Figure 7 shows the estimated density/count map of a testing example in each of the four datasets. The density maps estimated by the C-FCRN+Aux



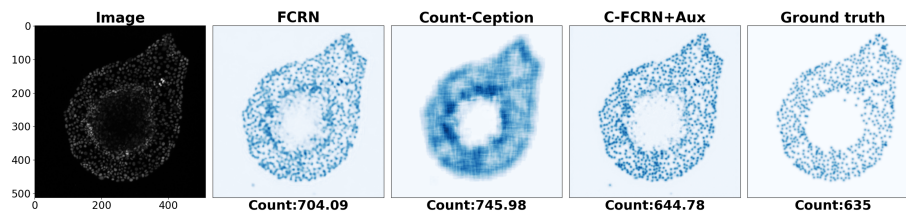
(a) Bacterial cells



(b) Bone marrow cells



(c) Colorectal cancer cells



(d) human embryonic stem cells (hESC)

Figure 7: Estimated density or count maps from a sample image in each of the four datasets. The panels from left to right on each row show the cell images and the density/count maps estimated by the FCRN, the Count-Ception, and the proposed method (C-FCRN+Aux), and the associated ground truth density maps, respectively.

appear visually closer to the ground truth density maps compared to the FCRN method. It is noted that the Count-Ception method predicts a count map directly without providing cell centroid locations, which is different from the other density regression-based methods.

Figure 8 shows the result of a testing example in each of the bacterial and bone marrow cell datasets by use of three detection-based methods (Mask R-CNN, U-Net and StructRegNet). The StructRegNet achieves more accurate results than the other two. One of the possible reasons is that the StructRegNet model is trained to regress a dense proximity map, in which the pixels closer to cell centroids can get more attention than those far-away pixels; this is different from the U-Net and Mask R-CNN model. This can benefit more for local maximum searching in the non-maximum post-process and yield better cell detection performance. It was also observed that the three detection-based methods commonly failed in detecting clustered and occluded cells in the bacterial image example. Also, they either under-detect or over-detect cells in the bone marrow image example. These images contain strongly heterogenous backgrounds and the shapes of cells vary largely. The inaccuracy of cell detection with these detection-based methods confirms their lower cell counting accuracy shown in the Table 2 and Figure 6.

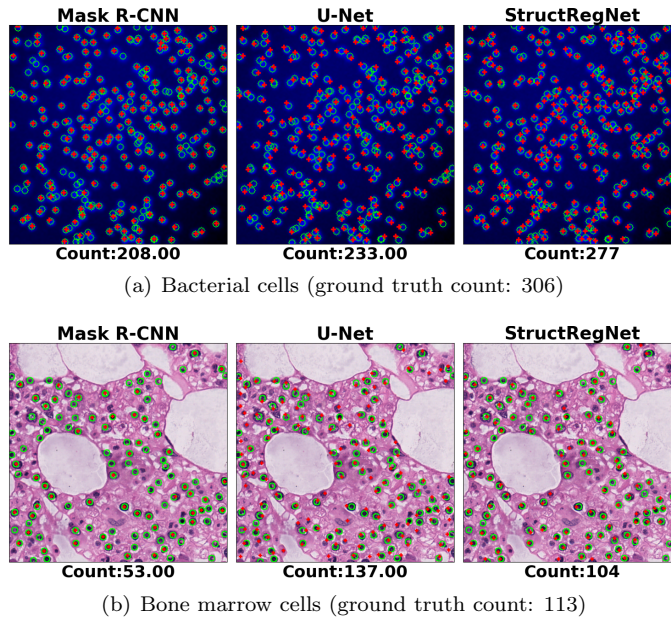
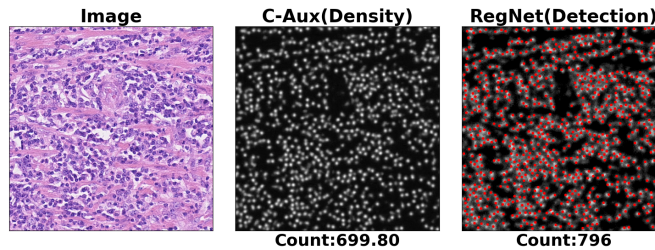


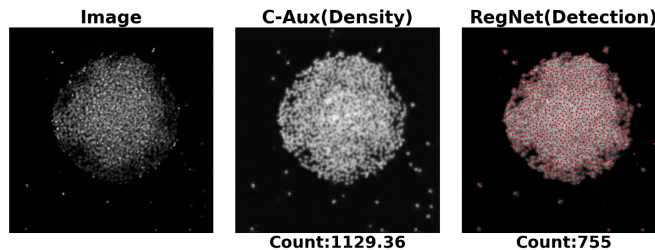
Figure 8: Example results of three deep-learning detection-based cell counting methods (Mask R-CNN, U-Net, and StructRegNet). Panels (a) and (b) show the prediction results on the bacterial and bone marrow datasets, respectively. The green cycles and red dots in each image represent the ground truth annotations and the detected cell centroids, respectively.

Figure 9 shows the result on an example in each of the colorectal and hESC

datasets by use of the proposed method and the StructRegNet method, which are the best-performing regression-based method and detection-based method tested in this study, respectively. The cells are commonly concentrated in colorectal cell images and seriously clustered and occluded in the hESC images. Cell detection in these two scenarios is extremely challenging. The StructRegNet method shows much worse counting performance compared to the proposed method.



(a) Colorectal cancer cells (ground truth count: 712)



(b) hESC (ground truth count: 1100)

Figure 9: Example prediction results based on the proposed C-FCRN+Aux method and the detection-based method (StructRegNet). Here, “image”, “C-Aux” and “RegNet” represent the processed image and the estimated density map using the “C-FCRN+Aux” method and the computed proximity map using the “StructRegNet” methods. The red dots represent the detected cell centroids based on the computed proximity map, respectively.

4.2. Benefits of using AuxCNNs to support C-FCRN training

The accuracy of the estimated density map along the training process was investigated to demonstrate that AuxCNNs supports the overall model training process. Figure 10 shows the curves of validation losses vs. the number of epochs for the proposed method and the other two nonlinear density regression methods (C-FCRN-only and FCRN) on four datasets. One of the five validation curves generated during the 5-fold cross validation procedure is presented for each method as an example. The curves generated for the first 500 epochs are shown because the validation losses keep stable after the 500-th epoch. As shown in Figure 10, the curves from all three methods converge when the number of epochs increases, which reflects the stability of training process. In

addition, the curves of the proposed C-FCRN+Aux method are significantly lower compared to the other two for all four datasets, which demonstrate that the proposed method allows to train a model that yields better model-fitting with the deep supervisions from the AuxCNNs. This analysis of validation loss over the training process is consistent with the results shown in Tables 2 and Figure 6, and reflects the better model fitting and generalization of our DRM to the validation data.

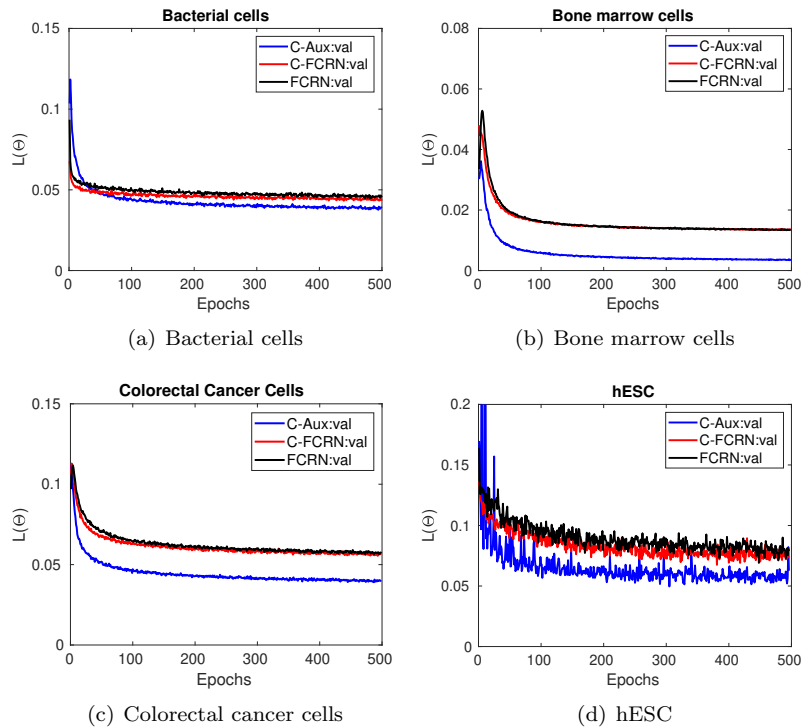


Figure 10: Validation losses as the functions of epochs are plotted for the DRM training on the four datasets. C-Aux and C-FCRN are abbreviations of C-FCRN+Aux and C-FCRN-Only methods, respectively.

4.3. Computation efficiency comparison

The computation efficiencies of the seven deep convolutional neural network-based methods, including the proposed method, FCRN, C-FCRN-Only, CountCeption, StructRegNet, U-Net and Mask R-CNN, were compared. The average processing time on testing images with the same GPU settings was employed as the comparison metric. Table 3 shows that the proposed method costs comparable counting time compared to the FCRN and the C-FCRN-Only methods. The CountCeption method is the more time-consuming one in comparison to

the other three regression-based methods. In the Count-Ception method, max-pooling layers are not employed in the network, and filters with large spatial size (5×5 pixels) are employed for extracting multi-scale features from images. These two reasons induce a large amount of convolution operations between high-dimension feature maps and large-sized filters, therefore, leading to the high computation workload in the Count-Ception method.

Density regression-based methods are less time-consuming than the three detection-based methods (StructRegNet, U-Net, and Mask R-CNN). The main reason is that the non-maximum suppression post-processing for cell detection costs a considerable amount of time. Mask R-CNN takes particularly longer time because of its superior larger network size and its aim at predicting separate masks for each cell, which is a much more complex task compared to the cell counting task.

Table 3: Computational Efficiency Comparison

Seconds/image	Bacterial cells	Bone marrow cells	Colorectal cancer cells	hESC
Mask R-CNN	0.55279	0.89527	—	—
U-Net	0.07646	0.16125	—	—
StructRegNet	0.06648	0.08035	0.18690	0.36167
FCRN	0.00468	0.02568	0.017	0.01901
C-FCRN-Only	0.00511	0.02846	0.01925	0.02134
Count-Ception	0.25111	0.18185	0.16308	0.19208
C-FCRN+Aux	0.00554	0.03113	0.02233	0.02275

Seconds/image represents the processing time for one image.

5. Discussion

The method proposed in this study combines the advantage of FCRN design with concatenation layers and a deeply-supervised learning strategy. It solves the two shortcomings that exist in the original FCRN. The concatenated layers integrates multi-scale features extracted from non-adjacent layers to improve the granularity of the extracted features and further support the density map estimation. The deeply-supervised learning strategy permits a direct supervision from AuxCNNs on learning its intermediate layers to mitigate the potential vanishing gradient issue and improve the cell counting performance. The results on four image datasets show superior cell counting performance of the proposed method compared to the other eight state-of-the-art methods. In addition, compared to the original FCRN, the proposed method improve the counting performance on four datasets ranging from 13% to 31% in terms of MAE. The computational efficiency of the proposed method is comparable to other density regression-based methods. The proposed method is capable of processing arbitrary-size images and estimating their density maps by use of fully convolutional layers in the C-FCRN. The proposed method could also be applied to heterogeneous cell assemblies, if cell types of interest are annotated in the training images. This deeply supervised learning framework will encourage

the trained DRM to focus on the cell types of interest but consider cells of other types as background.

The proposed method, other four regression-based and four detection-based methods were investigated on four challenging datasets. In general, the density regression-based methods yielded better performance and had three advantages over the detection-based methods. First, the regression-based methods count cells without cell detection, which can avoid challenging cell detection issues that commonly exist in microscopy images. Second, density regression-based methods are convenient for deployment, since they do not require trivial post-processings such as thresholding and non-maximum suppression. Thirdly, density regression-based methods can count cells more efficiently, i.e. the counting for an image of 512×512 pixels takes about $20ms$. The three advantages enable the density-regression based methods to be potentially applied to real-time clinical applications. In addition, it should be noted that even though the detection-based methods yield lower performance on this cell counting task, they are more suitable for the segmentation of cells of other types for other applications [66, 67]. Generally, for those cell types of interest, the cells in the acquired microscopy images are less overlapped and the cell masks can be fully annotated. In addition, the kernel sizes shown in Eq. 8 is determined by K_G , which is chosen according to the sizes of cells in the processed image to guarantee that the touching areas between occluded cells have been appropriately represented on the related density map. In this study, the radii of cells in the four datasets are less than 8 pixels. We then set the kernel size $(2K_G + 1) \times (2K_G + 1)$ to 21×21 pixels.

In the current study, all images are pre-processed by simply normalizing the intensities to the range of $[0, 1]$ to increase the stability of the model training process. In the future, we will investigate other image denoising and/or image enhancement methods to more accurately count cells for images that exhibit highly inhomogeneous tissue backgrounds and noises, or yield low image contrast. Also, the cell centroids used for generating ground truth density maps in the three experimental datasets were manually annotated by human experts, which may be subject to subjective errors. This might be one of the reasons that the MREs of these three experimental datasets (shown in Figure 6) were higher than that of the synthetic bacterial dataset. More accurate annotation strategies will be investigated to reduce the uncertainty in generating ground truth density maps. In this study, a uniform network architecture of C-FCRN+Aux was applied to learn DRMs separately on each of the four distinct datasets. We will adapt some other variants of FCRNs in the future that aim at crowd counting tasks [68, 69, 70] for varied datasets.

6. Conclusion

A deeply-supervised density regression model is proposed in this study for accurately and robustly counting the number of cells in microscopy images. The proposed method is capable of processing varied-size images containing dense cell clusters, large variations of cell morphology and inhomogeneous background

noise. Extensive experiments based on four datasets representing different image modalities and image acquisition techniques demonstrated the efficiency, robustness, and generality of the proposed method. The proposed method can be potentially to be applied to real-time clinical applications. It also holds the promise to be applied to a number of different problems, such as object counting (other than cells) in crowded scenes.

Acknowledgment

This work was supported in part by award NIH R01EB020604, R01EB023045, R01NS102213, R01CA233873, R21CA223799, and a grant from Children Discovery Institute (LSK). The dataset of human embryonic stem cells are provided by Solnica-Krezel group at Washington University School of Medicine.
5 The authors greatly appreciate the useful discussion with Dr. Su Ruan at The University at Rouen and Dr. Frank Brooks at The University of Illinois at Urbana-Champaign. The authors would like to thank the anonymous reviewers for valuable comments and suggestions.

References

10 References

- [1] V. Lempitsky, A. Zisserman, Learning to count objects in images, in: Advances in neural information processing systems, 2010, pp. 1324–1332.
- [2] C. Arteta, V. Lempitsky, J. A. Noble, A. Zisserman, Interactive object counting, in: European Conference on Computer Vision, Springer, 2014,
15 pp. 504–518.
- [3] W. Xie, J. A. Noble, A. Zisserman, Microscopy cell counting and detection with fully convolutional regression networks, Computer methods in biomechanics and biomedical engineering: Imaging & Visualization 6 (3) (2018) 283–292.
- [4] F. Xing, L. Yang, Robust nucleus/cell detection and segmentation in digital
20 pathology and microscopy images: a comprehensive review, IEEE reviews in biomedical engineering 9 (2016) 234–263.
- [5] T. Falk, D. Mai, R. Bensch, Ö. Çiçek, A. Abdulkadir, Y. MARRAKCHI, A. Böhm, J. Deubner, Z. Jäckel, K. Seiwald, et al., U-net: deep learning for cell counting, detection, and morphometry, Nature methods 16 (1)
25 (2019) 67–70.
- [6] C. F. Koyuncu, G. N. Gunesli, R. C. Atalay, C. Gunduz-Demir, Deepdistance: A multi-task deep regression model for cell detection in inverted microscopy images, Medical Image Analysis (2020) 101720.

- 30 [7] S. Zhang, D. Metaxas, Large-scale medical image analytics: Recent methodologies, applications and future directions (2016).
- [8] M. Wainberg, D. Merico, A. Delong, B. J. Frey, Deep learning in biomedicine, *Nature biotechnology* 36 (9) (2018) 829.
- [9] F. Xing, Y. Xie, H. Su, F. Liu, L. Yang, Deep learning in microscopy image analysis: A survey, *IEEE transactions on neural networks and learning systems* 29 (10) (2017) 4550–4568.
- 35 [10] G. Carneiro, Y. Zheng, F. Xing, L. Yang, Review of deep learning methods in mammography, cardiovascular, and microscopy image analysis, in: *Deep Learning and Convolutional Neural Networks for Medical Image Computing*, Springer, 2017, pp. 11–32.
- 40 [11] H. Irshad, A. Veillard, L. Roux, D. Racoceanu, Methods for nuclei detection, segmentation, and classification in digital histopathology: a review—current status and future potential, *IEEE reviews in biomedical engineering* 7 (2013) 97–114.
- 45 [12] C. L. Chen, A. Mahjoubfar, L.-C. Tai, I. K. Blaby, A. Huang, K. R. Niazi, B. Jalali, Deep learning in label-free cell classification, *Scientific reports* 6 (2016) 21471.
- [13] C. Hu, S. He, Y. J. Lee, Y. He, E. M. Kong, H. Li, M. A. Anastasio, G. Popescu, Label-free cell viability assay using phase imaging with computational specificity, *bioRxiv*.
- 50 [14] B. Venkatalakshmi, K. Thilagavathi, Automatic red blood cell counting using hough transform, in: *Information & Communication Technologies (ICT), 2013 IEEE Conference on*, IEEE, 2013, pp. 267–271.
- [15] A. S. Coates, E. P. Winer, A. Goldhirsch, R. D. Gelber, M. Gnant, M. Piccart-Gebhart, B. Thürlimann, H.-J. Senn, P. Members, F. André, et al., Tailoring therapies—improving 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) (2015) 1533–1546.
- 55 [16] L. Solnica-Krezel, Conserved patterns of cell movements during vertebrate gastrulation, *Current biology* 15 (6) (2005) R213–R228.
- [17] S.-C. Zhang, M. Wernig, I. D. Duncan, O. Brüstle, J. A. Thomson, In vitro differentiation of transplantable neural precursors from human embryonic stem cells, *Nature biotechnology* 19 (12) (2001) 1129.
- 60 [18] J. A. Thomson, J. Itskovitz-Eldor, S. S. Shapiro, M. A. Waknitz, J. J. Swiergiel, V. S. Marshall, J. M. Jones, Embryonic stem cell lines derived from human blastocysts, *science* 282 (5391) (1998) 1145–1147.
- 65

- [19] O. V. Lagutin, C. C. Zhu, D. Kobayashi, J. Topczewski, K. Shimamura, L. Puelles, H. R. Russell, P. J. McKinnon, L. Solnica-Krezel, G. Oliver, Six3 repression of wnt signaling in the anterior neuroectoderm is essential for vertebrate forebrain development, *Genes & development* 17 (3) (2003) 368–379.
- [20] J. Matas, O. Chum, M. Urban, T. Pajdla, Robust wide-baseline stereo from maximally stable extremal regions, *Image and vision computing* 22 (10) (2004) 761–767.
- [21] O. Barinova, V. Lempitsky, P. Kholi, On detection of multiple object instances using hough transforms, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34 (9) (2012) 1773–1784.
- [22] C. Arteta, V. Lempitsky, J. A. Noble, A. Zisserman, Learning to detect cells using non-overlapping extremal regions, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer, 2012, pp. 348–356.
- [23] F. Xing, H. Su, J. Neltner, L. Yang, Automatic ki-67 counting using robust cell detection and online dictionary learning, *IEEE Transactions on Biomedical Engineering* 61 (3) (2014) 859–870.
- [24] C. Arteta, V. Lempitsky, J. A. Noble, A. Zisserman, Detecting overlapping instances in microscopy images using extremal region trees, *Medical image analysis* 27 (2016) 3–16.
- [25] D. C. Cireşan, A. Giusti, L. M. Gambardella, J. Schmidhuber, Mitosis detection in breast cancer histology images with deep neural networks, in: *International Conference on Medical Image Computing and Computer-assisted Intervention*, Springer, 2013, pp. 411–418.
- [26] S. He, K. T. Minn, L. Solnica-Krezel, H. Li, M. Anastasio, Automatic microscopic cell counting by use of unsupervised adversarial domain adaptation and supervised density regression, in: *Medical Imaging 2019: Digital Pathology*, Vol. 10956, International Society for Optics and Photonics, 2019, p. 1095604.
- [27] F. Liu, L. Yang, 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*, Springer, 2017, pp. 63–72.
- [28] C. Sommer, L. Fiaschi, F. A. Hamprecht, D. W. Gerlich, Learning-based mitotic cell detection in histopathological images, in: *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, IEEE, 2012, pp. 2306–2309.
- [29] P. Soille, *Morphological image analysis: principles and applications*, Springer Science & Business Media, 2013.

- [30] H. Kong, H. C. Akakin, S. E. Sarma, A generalized laplacian of gaussian filter for blob detection and its applications, *IEEE transactions on cybernetics* 43 (6) (2013) 1719–1733.
- 110 [31] D. Reisfeld, H. Wolfson, Y. Yeshurun, Context-free attentional operators: The generalized symmetry transform, *International Journal of Computer Vision* 14 (2) (1995) 119–130.
- [32] S. He, J. Zheng, A. Maehara, G. Mintz, D. Tang, M. Anastasio, H. Li, Convolutional neural network based automatic plaque characterization for intracoronary optical coherence tomography images, in: *Medical Imaging 2018: Image Processing*, Vol. 10574, International Society for Optics and Photonics, 2018, p. 1057432.
- 115 [33] S. He, W. Zhou, H. Li, M. A. Anastasio, Learning numerical observers using unsupervised domain adaptation, in: *Medical Imaging 2020: Image Perception, Observer Performance, and Technology Assessment*, Vol. 11316, International Society for Optics and Photonics, 2020, p. 113160W.
- [34] R. Zhu, D. Sui, H. Qin, A. Hao, An extended type cell detection and counting method based on fcn, in: *2017 IEEE 17th International Conference on Bioinformatics and Bioengineering (BIBE)*, IEEE, 2017, pp. 51–56.
- 125 [35] M. Tofghi, T. Guo, J. K. Vanamala, V. Monga, Prior information guided regularized deep learning for cell nucleus detection, *IEEE transactions on medical imaging*.
- [36] Y. Xie, F. Xing, X. Shi, X. Kong, H. Su, L. Yang, Efficient and robust cell detection: A structured regression approach, *Medical image analysis* 44 (2018) 245–254.
- 130 [37] H. Bischof, C. Payer, D. Stern, M. Feiner, M. Urschler, Segmenting and tracking cell instances with cosine embeddings and recurrent hourglass networks.
- [38] A. Paulauskaite-Taraseviciene, K. Sutiene, J. Valotka, V. Raudonis, T. Iesmantas, Deep learning-based detection of overlapping cells, in: *Proceedings of the 2019 3rd International Conference on Advances in Artificial Intelligence*, 2019, pp. 217–220.
- 135 [39] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- 140 [40] A. Khan, S. Gould, M. Salzmann, Deep convolutional neural networks for human embryonic cell counting, in: *European Conference on Computer Vision*, Springer, 2016, pp. 339–348.

- 145 [41] Y. Xue, N. Ray, J. Hugh, G. Bigras, Cell counting by regression using convolutional neural network, in: European Conference on Computer Vision, Springer, 2016, pp. 274–290.
- [42] J. P. Cohen, G. Boucher, C. A. Glastonbury, H. Z. Lo, Y. Bengio, Countception: Counting by fully convolutional redundant counting, in: Computer Vision Workshop (ICCVW), 2017 IEEE International Conference on, IEEE, 2017, pp. 18–26.
- 150 [43] L. Fiaschi, U. Köthe, R. Nair, F. A. Hamprecht, Learning to count with regression forest and structured labels, in: Pattern Recognition (ICPR), 2012 21st International Conference on, IEEE, 2012, pp. 2685–2688.
- [44] E. Walach, L. Wolf, Learning to count with CNN boosting, in: European Conference on Computer Vision, Springer, 2016, pp. 660–676.
- 155 [45] N. Nitta, T. Sugimura, A. Isozaki, H. Mikami, K. Hiraki, S. Sakuma, T. Iino, F. Arai, T. Endo, Y. Fujiwaki, et al., Intelligent image-activated cell sorting, *Cell* 175 (1) (2018) 266–276.
- [46] S. S. Alahmari, D. Goldgof, L. Hall, H. A. Phoulady, R. H. Patel, P. R. Mouton, Automated cell counts on tissue sections by deep learning and unbiased stereology, *Journal of chemical neuroanatomy* 96 (2019) 94–101.
- 160 [47] Q. Liu, A. Junker, K. Murakami, P. Hu, Automated counting of cancer cells by ensembling deep features, *Cells* 8 (9) (2019) 1019.
- [48] J. Sierra, J. Pineda, E. Viteri, A. Tello, M. Millán, V. Galvis, L. Romero, A. Marrugo, Generating density maps for convolutional neural network-based cell counting in specular microscopy images, in: *Journal of Physics: Conference Series*, Vol. 1547, IOP Publishing, 2020, p. 012019.
- 165 [49] Y. Zheng, Z. Chen, Y. Zuo, X. Guan, Z. Wang, X. Mu, Manifold-regularized regression network: A novel end-to-end method for cell counting and localization, in: *Proceedings of the 2020 the 4th International Conference on Innovation in Artificial Intelligence*, 2020, pp. 121–124.
- 170 [50] Q. Liu, A. Junker, K. Murakami, P. Hu, A novel convolutional regression network for cell counting, in: *2019 IEEE 7th International Conference on Bioinformatics and Computational Biology (ICBCB)*, IEEE, 2019, pp. 44–49.
- 175 [51] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *International Conference on Medical image computing and computer-assisted intervention*, Springer, 2015, pp. 234–241.
- 180 [52] H. Dong, G. Yang, F. Liu, Y. Mo, Y. Guo, Automatic brain tumor detection and segmentation using U-Net based fully convolutional networks, in: *Annual Conference on Medical Image Understanding and Analysis*, Springer, 2017, pp. 506–517.

- 185 [53] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, Z. Tu, Deeply-supervised nets, in: *Artificial Intelligence and Statistics*, 2015, pp. 562–570.
- [54] G. Zeng, X. Yang, J. Li, L. Yu, P.-A. Heng, G. Zheng, 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*, Springer, 2017, pp. 274–282.
- 190 [55] Q. Dou, L. Yu, H. Chen, Y. Jin, X. Yang, J. Qin, P.-A. Heng, 3d deeply supervised network for automated segmentation of volumetric medical images, *Medical image analysis* 41 (2017) 40–54.
- [56] S. He, K. T. Minn, L. Solnica-Krezel, M. Anastasio, H. Li, Automatic microscopic cell counting by use of deeply-supervised density regression model, in: *Medical Imaging 2019: Digital Pathology*, Vol. 10956, International Society for Optics and Photonics, 2019, p. 109560L.
- 195 [57] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: *Proceedings of COMPSTAT'2010*, Springer, 2010, pp. 177–186.
- [58] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning representations by back-propagating errors, *nature* 323 (6088) (1986) 533.
- 200 [59] A. Lehmussola, P. Ruusuvuori, J. Selinummi, H. Huttunen, O. Yli-Harja, Computational framework for simulating fluorescence microscope images with cell populations, *IEEE Transactions on Medical Imaging* 26 (7) (2007) 1010–1016.
- 205 [60] P. Kainz, M. Urschler, S. Schuler, P. Wohlhart, V. Lepetit, You should use regression to detect cells, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer, 2015, pp. 276–283.
- [61] K. Sirinukunwattana, S. E. A. Raza, Y.-W. Tsang, D. R. Snead, I. A. Cree, N. M. Rajpoot, Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images, *IEEE transactions on medical imaging* 35 (5) (2016) 1196–1206.
- 210 [62] K. T. Minn, Y. C. Fu, S. He, S. C. George, M. A. Anastasio, S. A. Morris, L. Solnica-Krezel, High-resolution transcriptional and morphogenetic profiling of cells from micropatterned human embryonic stem cell gastruloid cultures, *DEVELOPMENTAL-CELL-D-20-00079*.
- 215 [63] J. Jin, M. Li, L. Jin, Data normalization to accelerate training for linear neural net to predict tropical cyclone tracks, *Mathematical Problems in Engineering* 2015.
- 220 [64] A. M. Saxe, J. L. McClelland, S. Ganguli, Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, arXiv preprint arXiv:1312.6120.

- [65] W. Abdulla, Mask r-cnn for object detection and instance segmentation on keras and tensorflow.
- 225 [66] D. Zhang, Y. Song, D. Liu, H. Jia, S. Liu, Y. Xia, H. Huang, W. Cai, Panoptic segmentation with an end-to-end cell r-cnn for pathology image analysis, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2018, pp. 237–244.
- 230 [67] J. W. Johnson, Adapting mask-rcnn for automatic nucleus segmentation, arXiv preprint arXiv:1805.00500.
- [68] C. Zhang, H. Li, X. Wang, X. Yang, Cross-scene crowd counting via deep convolutional neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 833–841.
- [69] E. Walach, L. Wolf, Learning to count with cnn boosting, in: European Conference on Computer Vision, Springer, 2016, pp. 660–676.
- 235 [70] V. A. Sindagi, V. M. Patel, CNN-based cascaded multi-task learning of high-level prior and density estimation for crowd counting, in: 2017 14th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), IEEE, 2017, pp. 1–6.