

Detecting and Classifying Nuclei on a Budget

Joseph G. Jacobs^{1,2}, Gabriel J. Brostow², Alex Freeman³, Daniel C. Alexander^{1,2}, and Eleftheria Panagiotaki^{1,2}

¹ Centre for Medical Image Computing, University College London, London, UK

² Department of Computer Science, University College London, London, UK

³ Department of Histopathology, University College London Hospitals NHS Foundation Trust, University College London, London, UK
j.jacobs@cs.ucl.ac.uk

Abstract. The benefits of deep neural networks can be hard to realise in medical imaging tasks because training sample sizes are often modest. Pre-training on large data sets and subsequent transfer learning to specific tasks with limited labelled training data has proved a successful strategy in other domains. Here, we implement and test this idea for detecting and classifying nuclei in histology, important tasks that enable quantifiable characterisation of prostate cancer. We pre-train a convolutional neural network for nucleus detection on a large *colon* histology dataset, and examine the effects of fine-tuning this network with different amounts of *prostate* histology data. Results show promise for clinical translation. However, we find that transfer learning is not always a viable option when training deep neural networks for nucleus *classification*. As such, we also demonstrate that semi-supervised ladder networks are a suitable alternative for learning a nucleus classifier with limited data.

1 Introduction

Measures of cell nuclei show increasing promise for improving cancer characterization, providing useful diagnostic and prognostic information for different pathologies. For instance, the amount of different types of cells in prostate tissue (epithelial, fibroblast, etc.) strongly correlates to the Gleason grade of prostate cancer[3,5]. Lee et al.[7] show that nuclei orientation entropy in prostatectomies is a predictor of biochemical recurrence in cancer patients. However, such quantitative histological analysis is rarely used in clinical practice: manual nucleus detection and classification is *extremely* time consuming due to the high resolution of histological images, and the requisite expertise is also expensive. There is a critical need for computer-aided diagnosis tools for nuclei in histology.

Automatic nucleus detection is a well studied problem[1,6]. The current state-of-the-art systems use convolutional neural networks (CNNs) to perform spatial regression for predicting the location of nuclei[12,15]. Similarly, the best nucleus type classification methods also use CNNs. Sirinukunwattana et al.[12] use an ensemble of CNN predictions to classify images containing previously detected nuclei. Wang et al.[14] go a step further and train CNNs for simultaneous nucleus

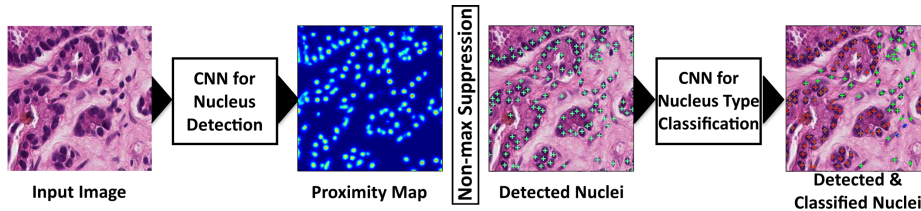


Fig. 1: Pipeline for detecting and classifying nuclei in histology.

detection and classification in lung biopsies while Bayramoglu and Heikkilä[2] show that transfer learning from natural images is useful for improving both the training time and performance of nucleus classification CNNs. While these perform impressively, a limitation of CNNs is that they typically require fully supervised training with thousands of labelled nuclei to prevent overfitting. This can be a major barrier for entry within the medical community as the data needs to be labelled by expert clinicians, which makes producing large labelled datasets expensive, both in time and cost.

This paper examines methods for learning a prostate nucleus detector and classifier given modest amounts of labelled prostate nuclei data. Specifically, we explore the viability of transfer learning for prostate nucleus detection by fine-tuning CNNs pre-trained on colon data and semi-supervised learning with Γ -ladder networks for prostate nucleus classification. These methods attempt to exploit the availability of large amounts of data from other domains (transfer learning) and large amounts of unlabelled prostate data (semi-supervised learning) respectively. The following sections describe the detail of these methods (§ 2) and our experiments (§ 3).

2 Methods

This section describes the CNN models (subsection 2.1), the fine-tuning procedure used to perform transfer learning (subsection 2.2) and the ladder network architecture used for semi-supervised learning (subsection 2.3).

2.1 Detecting and Classifying Nuclei with CNNs

Nucleus Detection Like [6,12,15] we formulate nucleus detection as a regression task as seen in Fig. 1. Given an input histology image, the nucleus detector predicts a function $d(x)$ that expresses the proximity of each pixel x to the nearest nucleus centroid. The local maxima in the resulting proximity map correspond to predicted nucleus centroids. We use the proximity function from [6]:

$$d(x) = \mathbb{I}\left[D(x) \leq d_{\max}\right] \left(e^{\alpha \left(1 - \frac{D(x)}{d_{\max}}\right)} - 1 \right) \quad (1)$$

Table 1: The fully convolutional network for nucleus detection.

#	Type	Filter Size	Stride	Padding
1	Convolution	$7 \times 7 \times 3 \times 16$	1×1	3×3
2	ReLU			
3	Max Pooling	3×3	1×1	1×1
4	Convolution	$5 \times 5 \times 16 \times 16$	1×1	2×2
5	ReLU			
6	Max Pooling	3×3	1×1	1×1
7	Convolution	$5 \times 5 \times 16 \times 16$	1×1	2×2
8	ReLU			
9	Max Pooling	3×3	1×1	1×1
10	Convolution	$11 \times 11 \times 16 \times 128$	1×1	5×5
11	ReLU			
12	Convolution	$1 \times 1 \times 128 \times 128$	1×1	0×0
13	ReLU			
14	Convolution	$1 \times 1 \times 128 \times 1$	1×1	0×0

where $\mathbb{I}[a]$ is an indicator function, $D(x)$ is the Euclidean distance from x to the nearest nucleus centroid while α and d_{\max} control the height and radius of peaks in $d(x)$. We introduce a novel fully convolutional network (FCN) architecture (Table 1) that performs inference on an entire image in a single pass. This significantly speeds up both training and test time compared to sliding window methods⁴.

Nucleus Type Classifier The nucleus type classifier is a standard multi-class classification CNN that takes as input a 27×27 px nucleus patch and classifies it as either epithelial, inflammatory or miscellaneous. The network structure is the standard single patch predictor model described in [12].

2.2 Transfer Learning with CNNs

Unlike other machine learning methods, CNNs do not require hand-engineered input features. The convolutional layers in a CNN act as feature extractors that are learnt directly from data. However, this can be a major limitation as these convolutional filters need to be trained on large datasets to prevent overfitting. One way to avoid re-learning the convolutional filters for every task is by transfer learning. Instead of training a CNN from scratch, we begin with a model that is pre-trained on a separate, large dataset. This ensures that the model has useful convolutional filters when we begin training. The training procedure then fine-tunes these CNN weights for a particular task.

⁴ It takes approximately 80 minutes to process a $107\,250 \times 103\,168$ whole-slide prostatectomy image on an NVIDIA GTX980 (incl. disk I/O), comparable to [16].

We examine the suitability of using transfer learning for reducing the amount of training data required to train both nucleus classifiers and detectors. For both tasks, we pre-train CNNs on a publicly available dataset of labelled colon nuclei[12] and fine-tune the entire CNN with varying amounts of prostate data. This allows investigation of the trade-off between the amount of labelled prostate training data and performance of nucleus detection/classification CNNs.

2.3 Semi-supervised Learning with Ladder Networks

Another method for reducing the amount of labelled training data required is by using a semi-supervised learning framework. Given a dataset of N labelled images and M unlabelled images where often $M \gg N$, semi-supervised learning frameworks attempt to learn classification models that exploit both the labelled and unlabelled images. In this instance, we explore the suitability of the ladder network architecture[11] for learning a nucleus classifier. Ladder networks turn a standard neural network into a semi-supervised model by treating it as the encoder in a denoising autoencoder.

A standard neural network is turned into a ladder network by (i) adding a decoder network to turn the network into an autoencoder and (ii) adding skip connections from every layer in the encoder to the corresponding layer in the decoder. During training, noise is added to the outputs of each layer in the encoder and the training objective is to minimise the weighted sum of the supervised cost function and the unsupervised cost functions⁵. A special case of ladder networks is the Γ -ladder network where we only consider the denoising cost in the top-most layer of decoder network (i.e. we set the denoising cost weights of all other decoder layers to zero). In our experiments, we use the Γ -ladder CNNs to perform semi-supervised learning of nucleus classifiers as they are faster to train and have fewer hyperparameters to adjust.

3 Results & Discussion

3.1 Experimental Setup

Dataset All experiments were run on a dataset of H&E stained prostate biopsies collected from 34 cancer patients. The biopsies were digitised at $20\times$ magnification ($0.55\mu\text{m}$ per pixel) with a Leica SCN400 slide scanner and we extracted $400 \times 250 \times 250\text{px}$ images. A histopathologist with over 10 years experience in genitourinary pathology (AF) dot annotated 16,562 nuclei in these images. 10,062 of these were also labelled with one of three nuclei type labels: 4212 epithelial, 1866 inflammatory (lymphocytes, plasma cells and macrophages) and 3984 miscellaneous (fibroblasts, blood vessel walls, nerves, etc.). We divided this data

⁵ The supervised objective is the cross entropy cost at the top of the encoder while the unsupervised objectives are the denoising mean squared errors at each decoder layer. We refer the reader to [11] for a more detailed description of ladder networks.

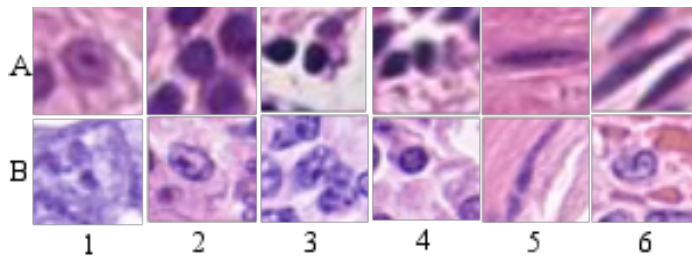


Fig. 2: Sample nuclei from the prostate (row A) and colon (row B) datasets. Columns 1 & 2 are epithelial nuclei, 3 & 4 are inflammatory nuclei and 5 & 6 are other miscellaneous nuclei.

into three patient-stratified sets: 40% for training, 10% for validation and 50% for testing.

For the fine-tuned models we pre-train the CNNs on a publicly available colon biopsy dataset[12]. The images are also H&E stained and digitised at $20\times$ magnification. The dataset contains 29,756 dot annotated nuclei, with type labels for 22,444 nuclei: 7,722 epithelial, 6,971 inflammatory and 7,751 miscellaneous.

CNN Training The CNN weights were randomly initialised from a normal distribution with mean 0 and standard deviation of 10^{-2} . The data was augmented with 90° , 180° and 270° rotations as well as flips along the horizontal and vertical axes. We trained the networks using stochastic gradient descent with Nesterov momentum[13] with a learning rate of 10^{-4} and minibatch sizes of 2 250×250 px image patches for the nucleus detectors⁶ and 128 27×27 nucleus patches for the nucleus classifiers. To prevent overfitting, the number of training epochs was determined independently for each network based on the value of the cost function on the held-out validation set. The optimal number of epochs ranged from as few as 20 training epochs for classification networks pre-trained on colon data to 5,000 training epochs for fully supervised classification CNNs trained using just 1% of labelled prostate training data. Other hyperparameters such as the unsupervised cost weight were similarly optimised on the held-out validation set independently for each network to prevent overfitting.

Evaluation Metrics We quantify the performance of a model by measuring the number of true positives (TP), false positives (FP) and false negatives (FN) produced by the model. TPs, FPs and FNs are well defined for classification problems. For nucleus detection, we define a TP as a predicted centroid that falls within a 6px radius of the ground truth annotation. FPs are predictions that

⁶ We note that since the FCN performs dense prediction on an input image, training/testing with a single 250×250 px image patch is equivalent to training/testing with 62,500 neighbouring patches using a patch-based method.

do not meet this criterion and FNs are ground truth annotations not associated with predictions. Based on these, we report four metrics for nucleus detection: (i) precision, $P = \frac{TP}{TP+FP}$ (ii) recall, $R = \frac{TP}{TP+FN}$, (iii) F₁ score, $F_1 = \frac{2PR}{P+R}$ and (iv) the area under the precision-recall curve (AUPR). For nucleus classification, we report the overall accuracy, the individual class F₁ scores, unweighted average of class F₁ scores (macro F₁) and the weighted average of class F₁ scores.

3.2 Nucleus Detection

Table 2: Precision, recall, F₁ score and AUPR for nucleus detectors trained with different amounts of labelled prostate images.

Labelled	Metrics	Baseline CNN	Fine-tuned CNN
1% 2 images	Precision	0.805 ± 0.002	0.827 ± 0.007
	Recall	0.862 ± 0.012	0.873 ± 0.004
	F ₁ Score	0.833 ± 0.006	0.849 ± 0.003
	AUPR	0.866 ± 0.004	0.896 ± 0.001
3% 6 images	Precision	0.825 ± 0.005	0.836 ± 0.012
	Recall	0.863 ± 0.009	0.872 ± 0.015
	F ₁ Score	0.844 ± 0.003	0.853 ± 0.001
	AUPR	0.877 ± 0.006	0.899 ± 0.005
5% 10 images	Precision	0.824 ± 0.007	0.845 ± 0.002
	Recall	0.875 ± 0.005	0.865 ± 0.003
	F ₁ Score	0.849 ± 0.003	0.855 ± 0.002
	AUPR	0.885 ± 0.003	0.901 ± 0.001
100% 200 images	Precision	0.843 ± 0.013	0.846 ± 0.004
	Recall	0.885 ± 0.007	0.882 ± 0.004
	F ₁ Score	0.864 ± 0.005	0.864 ± 0.003
	AUPR	0.910 ± 0.003	0.911 ± 0.004

Table 2 compares the baseline method (a CNN trained from scratch with the given labelled images) against a fine-tuned CNN pre-trained with colon data. The precision, recall and F₁ scores reported on the table are for the point on the precision-recall curve with the highest F₁ score. The results show that fine-tuned CNNs consistently outperform the baseline method. Although the precision, recall and F₁ scores of the baseline methods are similar to that of the fine-tuned models, the more revealing metric is the AUPR. The fine-tuned CNNs have much higher AUPR than baseline CNNs. Using just 1% of labelled prostate data, the fine-tuned CNN AUPR is comparable to that of the baseline method that uses 100% of labelled data (Fig. 3). This indicates that the fine-tuned CNNs are more robust to the choice of the threshold parameter used to discard false positives. The results also suggest that the convolutional filters learnt with the colon data are generalisable to prostate data for this task.

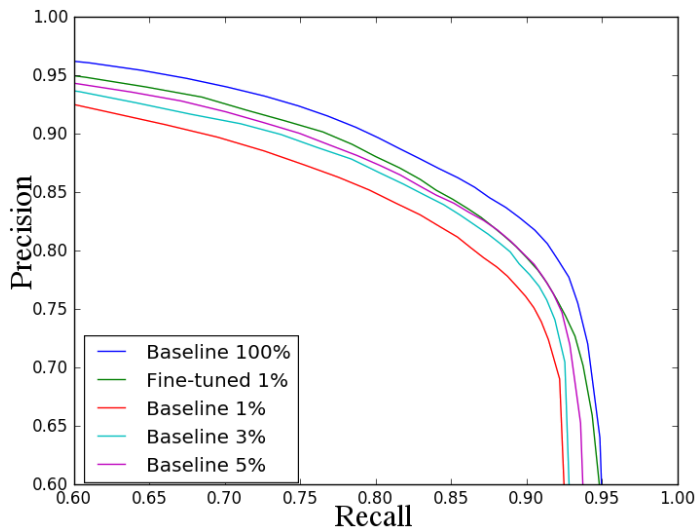


Fig. 3: Precision-recall curves for the baseline models and the 1% fine-tuned model.

3.3 Nucleus Classification

For nucleus classification, we compare our baseline method (a fully supervised CNN) against a fine-tuned CNN and a Γ -ladder CNN (Table 3). The results show that fine-tuning does not work very well when using 1% of labelled training data. Despite a 4-5% increase across the mean F_1 scores, we note that the scores have larger standard deviations compared to the baseline and Γ -ladder CNN, especially for inflammatory nuclei. The Γ -ladder CNN performs substantially better than the other two models across the different metrics using just 1% of the labelled data. The Γ -ladder CNN trained with 1% of labelled data even outperforms the baseline trained on 3% of labelled data on five of the six metrics.

We see a considerable jump in performance for the baseline and fine-tuned models when 3% of labelled data is used for training. While the Γ -ladder CNN improves as well and is still the best performing of the three models, the increase in performance is less substantial. Similarly, there is a marginal improvement in performance of all the three models as we increase the amount of labelled data to 5%. When using 100% of labelled data, we see identical performance for all models, with Γ -ladder CNNs performing marginally better than the other two.

The experiments indicate that Γ -ladder CNNs are the most robust of the three models. They perform well even when given a very small amount of labelled data and either matches or improves the performance of fully supervised and fine-tuned models when we increase the amount of labelled data. The large variation in inflammatory nucleus classification performance at 1% of labelled data could be explained by the small fraction of inflammatory nuclei in the prostate dataset

Table 3: F_1 metrics for nucleus classifiers trained with different amounts of labelled prostate nuclei patches.

Labelled	F_1 Scores	Baseline CNN	Fine-tuned CNN	Γ -ladder CNN
1% 40 nuclei	Weighted F_1	0.672 ± 0.007	0.721 ± 0.043	0.757 ± 0.017
	Macro F_1	0.639 ± 0.021	0.695 ± 0.061	0.738 ± 0.019
	Epithelial F_1	0.719 ± 0.031	0.765 ± 0.027	0.806 ± 0.015
	Inflammation F_1	0.486 ± 0.089	0.574 ± 0.145	0.654 ± 0.034
	Other F_1	0.713 ± 0.018	0.746 ± 0.017	0.755 ± 0.014
3% 120 nuclei	Weighted F_1	0.739 ± 0.008	0.763 ± 0.010	0.774 ± 0.008
	Macro F_1	0.725 ± 0.012	0.752 ± 0.008	0.762 ± 0.010
	Epithelial F_1	0.778 ± 0.009	0.810 ± 0.007	0.810 ± 0.009
	Inflammation F_1	0.664 ± 0.033	0.703 ± 0.006	0.707 ± 0.018
	Other F_1	0.733 ± 0.013	0.743 ± 0.018	0.769 ± 0.006
5% 200 nuclei	Weighted F_1	0.772 ± 0.002	0.780 ± 0.009	0.779 ± 0.007
	Macro F_1	0.761 ± 0.003	0.769 ± 0.007	0.765 ± 0.008
	Epithelial F_1	0.811 ± 0.006	0.821 ± 0.007	0.822 ± 0.008
	Inflammation F_1	0.713 ± 0.012	0.720 ± 0.012	0.704 ± 0.019
	Other F_1	0.758 ± 0.003	0.765 ± 0.017	0.770 ± 0.009
100% ~ 4000 nuclei	Weighted F_1	0.831 ± 0.004	0.828 ± 0.002	0.835 ± 0.004
	Macro F_1	0.820 ± 0.005	0.819 ± 0.002	0.825 ± 0.004
	Epithelial F_1	0.868 ± 0.003	0.863 ± 0.002	0.872 ± 0.003
	Inflammation F_1	0.772 ± 0.006	0.775 ± 0.003	0.778 ± 0.011
	Other F_1	0.821 ± 0.006	0.818 ± 0.001	0.824 ± 0.003

compared to the other classes. However, as previously noted there is an even larger variation in inflammatory nucleus classification performance for the fine-tuned CNN compared to the other models. This could potentially be explained by differences between the colon and prostate datasets. Inflammatory cells in the prostate dataset are mainly lymphocytes (Fig. 2, A3) while inflammatory cells in the colon are mainly macrophages (Fig. 2, B3) which are active and therefore look very similar to abnormal epithelial cells (Fig. 2, B2) with visible nucleoli.

4 Conclusions & Future Work

This paper adapts the general principles of transfer learning and semi-supervised learning for detecting and classifying cell nuclei on a budget. We demonstrate that transfer learning is suitable for learning nucleus detectors and classifiers given limited labelled data. However, it could potentially cause problems if there are biological differences in the tissue characteristics between the dataset used for pre-training and the dataset used for fine-tuning, as seen when attempting to learn a nucleus classifier with transfer learning. In this instance, we demonstrate that semi-supervised learning with Γ -ladder networks is a suitable alternative.

In future work, we will explore methods for including the full ladder network architecture, as well as histology from different organs and pathologies that could

benefit from this application (e.g. breast cancer). Additionally, a limitation of ladder networks is that they have more hyperparameters to optimise compared to standard neural networks. As such, future work will explore adapting other semi-supervised learning for neural networks [8], possibly adding query selection [4].

Acknowledgments We thank the EPSRC for funding EP’s (EP/N021967/1), DA’s (EP/M020533) and GB’s (EP/K015664/1, EP/K503745/1) work on this topic, the UCL Department of Computer Science for JJ’s studentship and the UCL Computer Science Cluster team.

References

1. Arteta, C., Lempitsky, V., Noble, A.J., Zisserman, A.: Learning to Detect Cells Using Non-overlapping Extremal Regions. In: Ayache, N., Delingette, H., Golland, P., Mori, K. (eds.) MICCAI 2012. LNCS, vol. 7510, pp. 348–356. Springer, Berlin, Heidelberg (2012)
2. Bayramoglu, N., Heikkilä, J.: Transfer Learning for Cell Nuclei Classification in Histopathology Images. In: Hua, G., Jégou, H. (eds.) ECCV 2016 Workshops. LNCS, vol. 9915, pp. 532–539. Springer, Cham (2012)
3. Chatterjee, A., Watson, G., Myint, E., Sved, P., McEntee, M., Bourne, R.M.: Changes in Epithelium, Stroma, and Lumen Space Correlate More Strongly with Gleason Pattern and Are Stronger Predictors of Prostate ADC Changes than Cellularity Metrics. *Radiology* 277(3), 751–762 (2015)
4. Gal, Y., Islam, R., Ghahramani, Z.: Deep Bayesian Active Learning with Image Data. In: NIPS Bayesian Deep Learning Workshop (2016)
5. Gorelick, L., Veksler, O., Gaed, M., Gómez, J.A., Moussa, M., Bauman, G.S., Fenster, A., Ward, A.D.: Prostate Histopathology: Learning Tissue Component Histograms for Cancer Detection and Classification. *IEEE Transactions on Medical Imaging* 32(10), 1804–1818 (2013)
6. Kainz, P., Urschler, M., Schuster, S., Wohlhart, P., Lepetit, V.: You Should Use Regression to Detect Cells. In: Navab et al. [9], pp. 276–283
7. Lee, G., Ali, S., Veltri, R., Epstein, J.I., Christudass, C., Madabhushi, A.: Cell Orientation Entropy (CORE): Predicting Biochemical Recurrence from Prostate Cancer Tissue Microarrays. In: Mori, K., Sakuma, I., Sato, Y., Barillot, C., Navab, N. (eds.) MICCAI 2013. LNCS, vol. 8151, pp. 396–403. Springer, Berlin, Heidelberg (2013)
8. Maaløe, L., Sønderby, C.K., Sønderby, S.K., Winther, O.: Auxiliary deep generative models. In: Balcan, M., Weinberger, K.Q. (eds.) ICML 2016. pp. 1445–1453. Proceedings of Machine Learning Research, PMLR (2016)
9. Navab, N., Hornegger, J., III, W.M.W., Frangi, A.F. (eds.): MICCAI 2015, LNCS, vol. 9351. Springer, Cham (2015)
10. Ourselin, S., Joskowicz, L., Sabuncu, M.R., Unal, G., Wells, W. (eds.): MICCAI 2016, LNCS, vol. 9901. Springer, Cham (2016)
11. Rasmus, A., Berglund, M., Honkala, M., Valpola, H., Raiko, T.: Semi-supervised Learning with Ladder Networks. In: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (eds.) NIPS 2015. pp. 3546–3554. Curran Associates, Inc. (2015)

12. Sirinukunwattana, K., Raza, S.E.A., Tsang, Y.W., Snead, D.R.J., Cree, I.A., Rajpoot, N.M.: Locality Sensitive Deep Learning for Detection and Classification of Nuclei in Routing Colon Cancer Histology Images. *IEEE Transactions in Medical Imaging* 35(5), 1196–1206 (2016)
13. Sutskever, I., Martens, J., Dahl, G., Hinton, G.: On the Importance of Initialization and Momentum in Deep Learning. In: Dasgupta, S., McAllester, D. (eds.) *ICML 2013*. pp. 1139–1147. *Proceedings of Machine Learning Research*, PMLR (2013)
14. Wang, S., Yao, J., Xu, Z., Huang, J.: Subtype Cell Detection with an Accelerated Deep Convolution Neural Network. In: Ourselin et al. [10], pp. 640–648
15. Xie, Y., Xing, F., Kong, X., Su, H., Yang, L.: Beyond Classification: Structured Regression for Robust Cell Detection Using Convolutional Neural Network. In: Navab et al. [9], pp. 358–365
16. Xu, Z., Huang, J.: Detecting 10,000 Cells in One Second. In: Ourselin et al. [10], pp. 676–684