

# WBC-Profiler: an Unsupervised Feature Learning System for Leukocytes Characterization and Classification

Hong Yan<sup>1,\*</sup>, Xuanyu Mao<sup>2,3,\*</sup>, Yongquan Xia<sup>4</sup>, Zhiyang Li<sup>4</sup>, Chengbin Wang<sup>5</sup>, Rui Xia<sup>1</sup>, Xuejing Xu<sup>4</sup>, Zhiqiang Wang<sup>1</sup>, Xie Zhao<sup>4</sup>, Yan Li<sup>1</sup>, Han Shen<sup>4,†</sup>, and Hang Chang<sup>2†</sup>

<sup>1</sup>Department of Laboratory Medicine, Nanjing Chest Hospital, Nanjing 210029, China

<sup>2</sup>Synihealth Research, Wuhan 430074, China

<sup>3</sup>Department of Biological Sciences, University of the Pacific, Stockton, CA 95211, United States

<sup>4</sup>Department of Laboratory Medicine, Affiliated Drum Tower Hospital of Nanjing University Medical School, Nanjing 210029, China

<sup>5</sup>Department of Laboratory Medicine, Chinese PLA General Hospital, Beijing 100853, China

\*Co-first author

†Co-corresponding author. Email: shenhan10366@sina.com, hchang.synihealth@gmail.com

December 29, 2018

## Abstract

The characterization and classification of white blood cells (WBC) is critical for the diagnosis of anemia, leukemia and many other hematologic diseases. We developed WBC-Profiler, an unsupervised feature learning system for quantitative analysis of leukocytes. We demonstrate that WBC-Profiler enables automatic extraction of complex signatures from microscopic images without human-intervention and thereafter effective construction of leukocytes profiles, which decouples large scale complex leukocytes characterization from limitations in both human-based feature engineering/optimization and the end-to-end solutions provided by modern deep neural networks, and therefore has the potential to provide new opportunities towards meaningful studies/applications with scientific and/or clinical impact

## 1 Introduction

2 Hematology tests provide laboratory assessments of blood formation and blood disorders, and play critical roles in the  
3 indication, diagnosis and evaluation of many conditions, including infection, inflammation and anemia. Among vari-  
4 ous hematology tests, the differential count of while blood cells (WBC) provides an important tool in diagnosing and  
5 monitoring infection and leukemic disorders, and the ratio of various kinds of leukocytes are commonly used as impor-  
6 tant markers. For example, the neutrophil to lymphocyte ratio (NLR) is used as a marker of subclinical inflammation,  
7 and recent studies suggest that increased NLR is independent predictor of mortality in patients undergoing angiogra-  
8 phy or cardiac revascularization [1], meanwhile it is also associated with poor prognosis of various cancers [2], such  
9 as esophageal cancer [3] or advanced pancreatic cancer [4].

10 Accurate characterization, detection and classification of while blood cells into several categories, including Mono-  
11 cytes, Lymphocytes, Basophils, Eosinophils, Atypical lymphocytes and Neutrophilic granulocytes, is critical for the  
12 ratio (proportion) assessment of leukocytes in blood cell slides. However, such a high-demanding task in clinical labo-  
13 ratory heavily relies on the manual annotation by pathologies, which is not only labor-intensive but also challenging  
14 given the shortage of experienced medical experts in many clinical laboratories. To overcome these obstacles, many  
15 efforts have been made towards the automatic blood cell classification. Among which, early development of com-  
16 mercial systems in 1970s failed to revolutionize the field due to the high price and low accuracy [5] of the products,  
17 and Leuko, another commercial leukemia diagnosis system, was later on developed to reveal an improved accuracy  
18 on cell classification based on naive Bayes classifiers. Meanwhile, research efforts on cell classification have also

19 been evolving from fuzzy logic techniques [6], support vector machines (SVMs) [7] to cellular neural networks [8].  
20 However, these works were mostly developed with a limited amount of data, and/or focused on images taken from  
21 specified instruments, which leave their generalization capability insufficiently justified.

22 Motivated by recent neuroscience findings [9, 10], unsupervised learning and deep learning techniques [11–13]  
23 have gained momentum during the past decade for object representation and recognition (e.g., face representation  
24 and recognition) [14–17]. And their applications in various biomedical tasks [18–22] have demonstrated success with  
25 the potential to provide a new avenue to data-intensive clinical studies, among which, the leukocytes classification  
26 accuracy has been significantly improved due to the employment of deep neural networks (DNNs) [23]. However,  
27 systems of such kind typically only provide the end-to-end (i.e., from data to classification) solution, which leaves  
28 the characterization of white blood cells inaccessible, and thus impede the construction of leukocytes profile for many  
29 potential needs, including profile interpretation, profile optimization, profile differentiation among cell types as well  
30 as profile association with other meaningful endpoints.

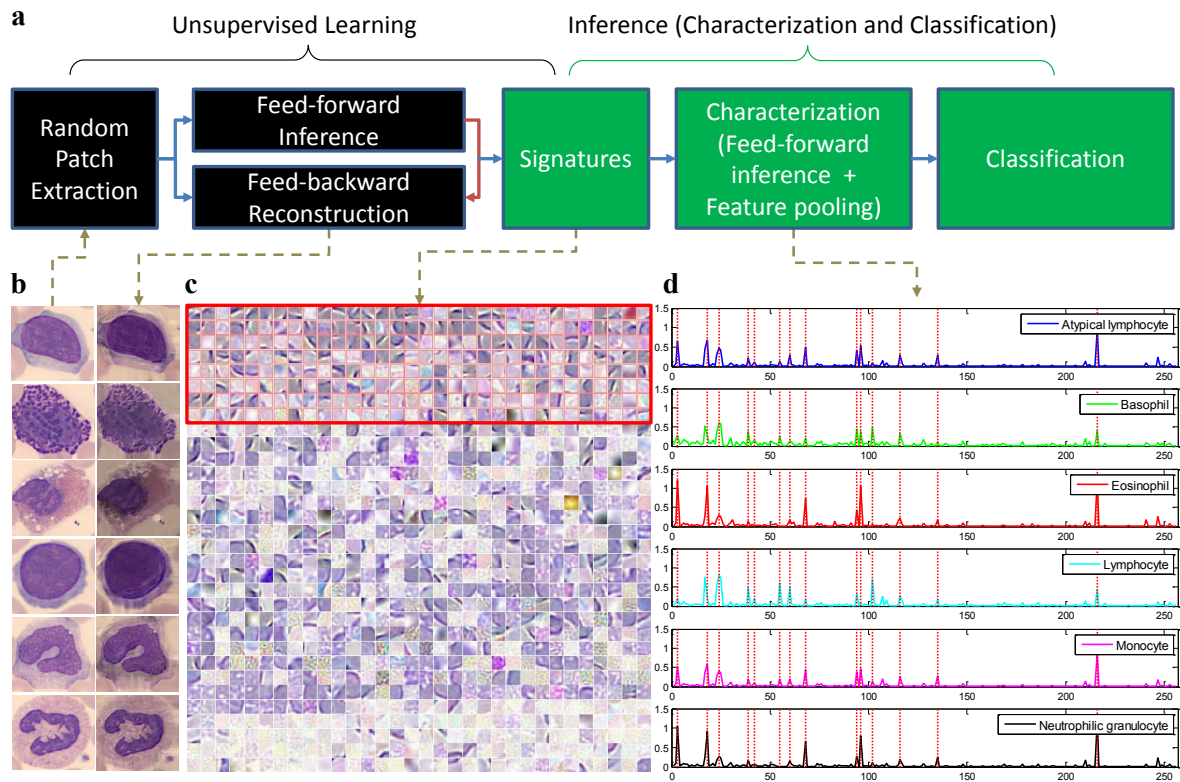
## 31 Results

32 We applied unsupervised feature learning for the automatic acquisition of intrinsic signatures directly from raw data  
33 (microscopic images) that facilitates the efficient and effective characterization and classification of leukocytes, skip-  
34 ping conventional steps such as feature engineering (i.e., manually design and optimize features) while providing  
35 meaningful cell profile that is typically inaccessible to many existing deep learning systems. We show that a well-  
36 designed unsupervised learning system is capable of automatic and efficient extraction of signatures (i.e., patterns  
37 encode both color and texture information) directly from microscopic images, which provides the foundation for both  
38 meaningful representation as well as effective classification of white blood cells.

39 The general principle of our leukocytes characterization and classification system (referred to as WBC-Profiler)  
40 is illustrated in Fig. 1. WBC-Profiler is a hierarchical unsupervised feature learning framework with feed-backward  
41 reconstruction and feed-forward feature inference. Unlike many unsupervised feature learning algorithms [24–27],  
42 WBC-Profiler involves only element-wise nonlinearity and matrix multiplication. Therefore, it provides an highly  
43 efficient and effective solution for Leukocytes characterization and classification.

44 To evaluate the capability of WBC-Profiler on leukocytes characterization, we visualized the signatures acquired  
45 by WBC-Profiler from our leukocytes image database in Fig. 1c. It contains 1024 elements, and captures information  
46 from both color and texture domains, which cannot be easily achieved manually. With the derived signatures, WBC-  
47 Profiler extracts patch-level features for each image patch through a feed-forward fashion, where an image patch refers  
48 to a sub-image with fixed size (i.e., 20 pixel by 20 pixel in our study) cropped from the original microscopic image,  
49 and the patch-level feature of a single image patch refers to the sparse code (coefficients) for the reconstruction of this  
50 specific image patch with derived signatures. The final profile (image-level representation) for the entire microscopic  
51 image is constructed with specific pooling operation over all patch-level features from the microscopic image. The  
52 selection of pooling strategy is typically guided by the effectiveness of the corresponding profile during classification,  
53 and in our case, mean-pooling was selected, which constructs the profile as the average sparse code (i.e., average re-  
54 construction coefficients) over all the image patches from the microscopic image. Heatmap of both selected signatures  
55 (see Fig. 2a) clearly indicates the differential expression of signatures across cell types, and provides an intuitive way  
56 to examine/evaluate the contribution of each signature for leukocyte image construction/composition. Furthermore,  
57 feature embedding of high-dimensional profile into 2-Dimensional (2D) space leads to cell-type-specific clusters (as  
58 illustrated in Fig. 2b-c), which demonstrates the effectiveness of WBC-Profiler in leukocytes characterization, espe-  
59 cially for the task of leukocytes classification.

60 To evaluate the capability of WBC-Profiler on leukocytes classification, we have compared it with one of the  
61 state-of-the-art techniques in the field of leukocytes classification (we refer it to as DeepVote, which combines the  
62 classification results from different deep neural networks to vote for the final decision) [23] and with one of the most  
63 successful deep learning techniques in the field of object detection and recognition, i.e., Faster R-CNN [28] with both  
64 vgg16 and res101 as the network architectures. During evaluation, half-half cross-validation was employed with 10  
65 iterations, where, at each iteration, we randomly selected 50% of the data per cell type for training, and used the rest  
66 for testing, and the performance in terms of average F1-measure and confusion matrix was illustrated in Figure 3,  
67 which demonstrates the effectiveness of WBC-Profiler for leukocytes classification.



**Figure 1: The concept of training and inference with WBC-Profiler for leukocytes characterization and classification.** **a**, basic structure of WBC-Profiler, where the *random patch extraction unit* selects a set of vectorized image patches randomly from input cell images; the *feed-backward reconstruction unit* reconstructs the original input signal (e.g., image patch) from a set of signatures learned from the input data; the *feed-forward inference unit* predicts the reconstruction coefficient (i.e., sparse code to be used as the feature representation for each individual image patch) for input data reconstruction; the *characterization unit* calculates the sparse codes for all image patches of a target cell image based on both the signatures and inference function derived from unsupervised learning, and summarize them into a single profile as the representation of the target cell image; and finally, the *classification unit* labels each cell image with different cell types. **b**, Examples of microscopic images of different types of white blood cells and the corresponding reconstruction results from the derived signatures. From top row to bottom row: Atypical lymphocytes, Basophils, Eosinophils, Lymphocytes, Monocytes and Neutrophilic granulocytes. **c**, Signatures (1024 in total) automatically learned from our cell image dataset by WBC-Profiler, where the top 256 signatures (ranked by random-forest based on their contribution to leukocytes classification) were highlighted within red bounding box. **d**, Mean profile of the top 256 signatures per cell type, where it is represented as the class-average contribution of each signature for the reconstruction of cell images in the same category.

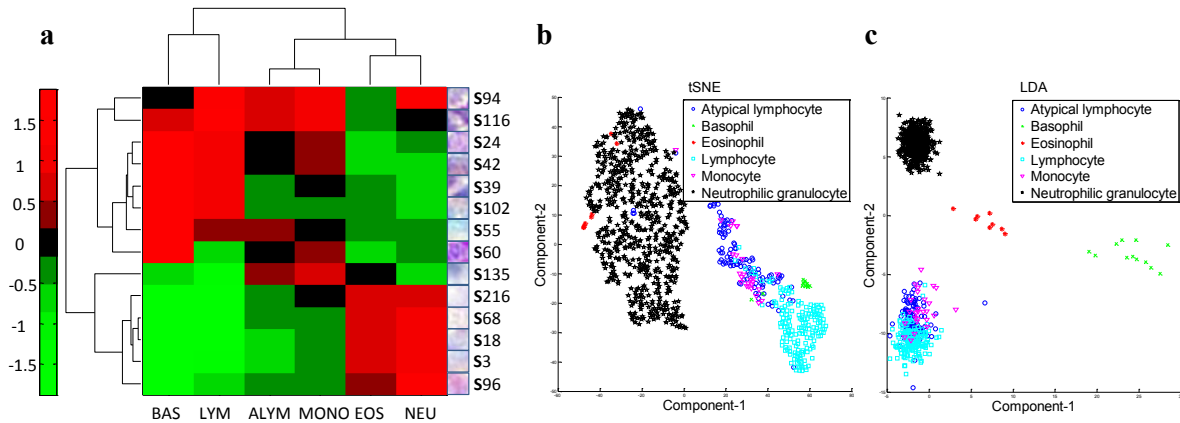


Figure 2: **WBC-Profiler provides effective leukocytes characterization.** **a**, Heatmap based on signatures selected from the mean profile (see Fig 1(d)), where zero-mean normalization is applied for better visualization. **b**, Feature embedding based on the top 256 signatures demonstrates the effectiveness of WBC-Profiler in unsupervised signature learning for leukocytes characterization, where t-Distributed Stochastic Neighbor Embedding (tSNE, unsupervised) reveals highly separable clusters. **c**, Separability of each clusters can be further improved through linear discriminant analysis (LDA, supervised).

## Conclusions

In summary, we developed WBC-Profiler, an unsupervised feature learning system for efficient and effective leukocytes characterization and classification. As demonstrated through a well-curated dataset, complex signatures, capturing both color and texture information, can be automatically and directly learned from the microscopic images. The designed architecture and the feed-forward nature of feature inference ensures WBC-Profiler's performance in precision, accuracy, and speed. Furthermore, WBC-Profiler decouples large scale complex leukocytes characterization from limitations in both human-based feature engineering/optimization and DNNs-based end-to-end solution, and therefore could further allow the extraction of highly multiplexed intrinsic properties and information from large scale leukocytes dataset towards meaningful endpoints with scientific and/or clinical impact.

## Acknowledgements

This work was supported by the Medical Key Science and Technology Development Projects of Nanjing (ZKX18016), the Medical Science and Technology Development Projects of Nanjing (YKK18167), and Synihealth Research.

## Author contributions

X.M. and H.C. wrote the software, performed the experiments, and analyzed the data. H.C., H.S., H.Y., and C.W. conceived and supervised the study. H.C., H.S., and H.Y. wrote the manuscript. Y.X., X.X., R.X., Z.W., Z.L., Y.L, and X.Z created the leukocytes image database.

## Competing interests

The authors declare no competing interests.

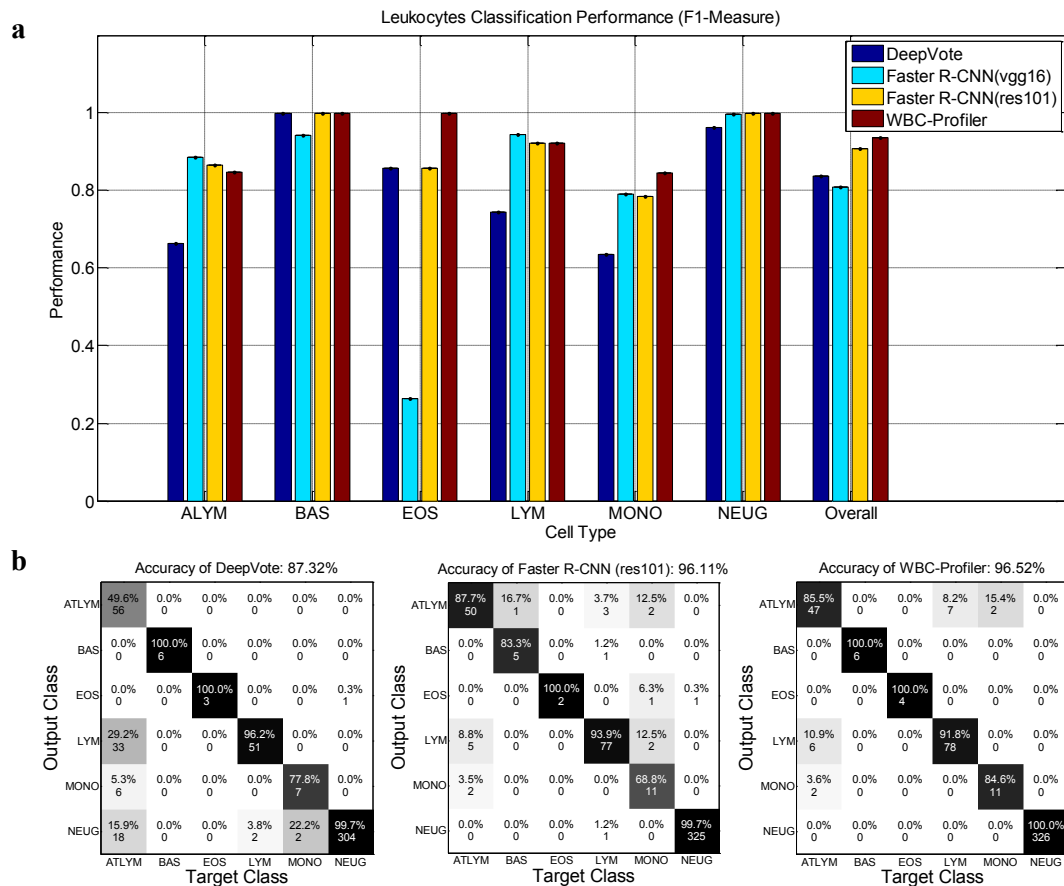


Figure 3: **WBC-Profiler enables accurate leukocytes classification.** **a**, Comparison of leukocytes classification performance in terms of F1-Measure for both individual class as well as overall class-average. **b**, Confusion matrix from DeepVote, Faster R-CNN (res101) and Stacked PSD for leukocytes classification with absolute cell image numbers (mean and rounded across all iterations) and normalized precision. Abbreviation: ATLYM - Atypical lymphocyte; BAS-Basophil; EOS - Eosinophil; LYM - Lymphocyte; MONO - Monocyte; NEUG - Neutrophilic granulocyte.

## 86 Methods

87 **Leukocytes Image Database.** A new dataset containing approximately 1000 microscopic images of 6 types of white  
88 blood cells has been collected at the Affiliated Drum Tower Hospital of Nanjing University Medical School to evaluate  
89 WBC-Profiler with detailed protocol listed as follows,

- 90 1. Sample selection. Abnormal blood test samples were selected by clinical experts for our study, where all samples  
91 were under daily intra-lab quality control as well as inter-lab quality control by the Clinical Lab Center, Jiangsu  
92 Province as well as MOPH (Ministry of Public Health) Clinical Lab Center.
- 93 2. Staining. Selected peripheral blood samples were filmed and stained with Wright-Gimsa through Sysmex SP 1000i.
- 94 3. Scanning. After staining, digital scanning was performed with OLYMPUS CX31RTSF at 1000X.
- 95 4. Labeling. The cell type labeling process were carried out by two clinical experts with more than 10 years of  
96 experience in the field.
- 97 5. Quality control. To ensure the data quality, only images meeting following criteria are selected:  
98 (a) technical correctness. No artifacts were introduced during staining and scanning process.  
99 (b) label consistency. The cell types of the target image independently labeled by the above two clinical experts  
100 were consistent with each other.
- 101 6. De-identification. To protect the privacy of patients, all images were de-identified by the clinical experts before any  
102 further quantitative study by WBC-Profiler.

103 **WBC-Profiler architecture.** WBC-Profiler is composed of unsupervised signature learning module, inference mod-  
104 ule and visualization module to realize efficient and effective leukocytes characterization and classification, as well as  
105 convenient interpretation of derived results.

106 *The unsupervised signature learning module* of WBC-Profiler is build upon the Stacked Predictive Sparse Decom-  
107 position (Stacked PSD) technique [29] for the construction of hierarchical unsupervised learning framework, which  
108 is suggested to be able to capture higher-level dependencies of input variables, thereby improving the ability of the  
109 system to capture underlying regularities in the data. Unlike many unsupervised feature learning algorithms [24–27],  
110 the feed-forward feature inference of PSD is very efficient, as it involves only element-wise nonlinearity and matrix  
111 multiplication. Therefore, it provides an highly efficient and effective solution for Leukocytes characterization. In  
112 our study, we used only one PSD layer with 1024 kernels (signatures) at a fixed size of 20-by-20 pixels. The training  
113 process was set to be 100 iterations, and it converged with around 20 iterations.

114 *The inference module* of WBC-Profiler consists of characterization and classification units, where the former ex-  
115 tracts patch-level features with a sliding window from the leukocytes image and construct the image-level profile  
116 through specific pooling operation, and the latter utilizes the pre-built profiles for leukocytes classification. In our  
117 study, we extracted patch-level features with a sliding window at the step size fixed to be 5 pixels, and adopted mean-  
118 pooling strategy among different popular pooling strategies through the evaluation on classification performance via  
119 support vector machine (SVM) classifier.

120 *The visualization module* of WBC-Profiler provides intuitive and convenient means for data visualization and  
121 interpretation, including signature visualization (as illustrated in Fig.1c), profile visualization (as illustrated in Fig.1d),  
122 heatmap (as illustrated in Fig. 2a ) and feature embedding (as illustrated in Fig.2b-c).

123 **Code availability.** Matlab scripts for WBC-Profiler, including unsupervised learning module, inference module  
124 and visualization module, are available as Supplementary Software, and further updates will be make available at  
125 <http://bmihub.org/project/wbc-profiler>.

126 **Data availability.** The data that support the findings of this study are available from the corresponding author upon  
127 request. Example data are available in the Supplementary Data and Supplementary Software packages.

## 128 References

- 129 [1] X. Wang, G. Zhang, X. Jiang, H. Zhu, Z. Lu, and L. Xu, “Neutrophil to lymphocyte ratio in relation to risk of all-cause  
130 mortality and cardiovascular events among patients undergoing angiography or cardiac revascularization: A meta-analysis of  
131 observational studies,” vol. 234, no. 1, pp. 206–213, 2014. Exported from <https://app.dimensions.ai> on 2018/11/26. 1

- 132 [2] A. J. Templeton, M. G. McNamara, B. eruga, F. E. Vera-Badillo, P. Aneja, A. Ocaa, R. Leibowitz-Amit, G. Sonpavde, J. J.  
133 Knox, B. Tran, I. F. Tannock, and E. Amir, “Prognostic role of neutrophil-to-lymphocyte ratio in solid tumors: A systematic  
134 review and meta-analysis,” *JNCI: Journal of the National Cancer Institute*, vol. 106, no. 6, p. dju124, 2014. 1
- 135 [3] J. Wang, Y. Jia, N. Wang, X. Zhang, B. Tan, G. Zhang, and Y. Cheng, “The clinical significance of tumor-infiltrating neu-  
136 trophils and neutrophil-to-cd8+ lymphocyte ratio in patients with resectable esophageal squamous cell carcinoma,” *Journal*  
137 *of Translational Medicine*, vol. 12, p. 7, Jan 2014. 1
- 138 [4] P. Xue, M. Kanai, Y. Mori, T. Nishimura, N. Uza, Y. Kodama, Y. Kawaguchi, K. Takaori, S. Matsumoto, S. Uemoto, and  
139 T. Chiba, “Neutrophil-to-lymphocyte ratio for predicting palliative chemotherapy outcomes in advanced pancreatic cancer  
140 patients,” *Cancer Medicine*, vol. 3, no. 2, pp. 406–415. 1
- 141 [5] K. Preston, “High-resolution leukocyte analyzers: retrospective and prospective,” *Appl. Opt.*, vol. 26, pp. 3258–3265, Aug  
142 1987. 1
- 143 [6] P. Sobrevilla, E. Montseny, and J. Keller, “White blood cell detection in bone marrow images,” in *18th International Confer-*  
144 *ence of the North American Fuzzy Information Processing Society - NAFIPS (Cat. No.99TH8397)*, pp. 403–407, June 1999.  
145 2
- 146 [7] D. M. Ushizima, A. C. Lorena, and A. C. P. L. F. de Carvalho, “Support vector machines applied to white blood cell recogni-  
147 tion,” in *Fifth International Conference on Hybrid Intelligent Systems (HIS’05)*, pp. 6 pp.–, Nov 2005. 2
- 148 [8] S. Wang and M. Wang, “A new detection algorithm (nda) based on fuzzy cellular neural networks for white blood cell  
149 detection,” *IEEE Transactions on Information Technology in Biomedicine*, vol. 10, pp. 5–10, Jan 2006. 2
- 150 [9] T. S. Lee and D. Mumford, “Hierarchical bayesian inference in the visual cortex,” *Journal of the Optical Society of America*,  
151 vol. 20, pp. 1434–1448, 2003. 2
- 152 [10] T. S. Lee, D. Mumford, R. Romero, and V. A. Lamme, “The role of the primary visual cortex in higher level vision,” *Vision*  
153 *Research*, vol. 38, no. 15/16, pp. 2429–2454, 1998. 2
- 154 [11] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” in *Proceedings*  
155 *of the IEEE*, pp. 2278–2324, 1998. 2
- 156 [12] F. J. Huang and Y. LeCun, “Large-scale learning with svm and convolutional for generic object categorization,” in *Proceed-*  
157 *ings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 1, CVPR ’06*,  
158 (Washington, DC, USA), pp. 284–291, IEEE Computer Society, 2006. 2
- 159 [13] G. E. Hinton and S. Osindero, “A fast learning algorithm for deep belief nets,” *Neural Computation*, vol. 18, pp. 1527–1554,  
160 2006. 2
- 161 [14] P. Y. Simard, D. Steinkraus, and J. C. Platt, “Best practices for convolutional neural networks applied to visual document  
162 analysis,” in *Proceedings of the Seventh International Conference on Document Analysis and Recognition - Volume 2, ICDAR*  
163 *’03*, (Washington, DC, USA), pp. 958–, IEEE Computer Society, 2003. 2
- 164 [15] M. Osadchy, Y. L. Cun, M. L. Miller, and P. Perona, “Synergistic face detection and pose estimation with energy-based  
165 model,” in *In Advances in Neural Information Processing Systems (NIPS)*, pp. 1017–1024, MIT Press, 2005. 2
- 166 [16] B. Kwolek, “Face detection using convolutional neural networks and gabor filters,” in *Proceedings of the 15th international*  
167 *conference on Artificial Neural Networks: biological Inspirations - Volume Part I, ICANN’05*, (Berlin, Heidelberg), pp. 551–  
168 556, Springer-Verlag, 2005. 2
- 169 [17] S. Sukittanon, A. C. Surendran, J. C. Platt, and C. J. C. Burges, “Convolutional networks for speech detection.,” in *INTER-*  
170 *SPEECH*, ISCA, 2004. 2
- 171 [18] S. J. McKenna, I. W. Ricketts, A. Y. Cairns, and K. A. Hussein, “A comparison of neural network architectures for cervical  
172 cell classification,” in *1993 Third International Conference on Artificial Neural Networks*, pp. 105–109, May 1993. 2
- 173 [19] H. A. Elsalamony, “Detection of some anaemia types in human blood smears using neural networks,” *Measurement Science*  
174 *and Technology*, vol. 27, no. 8, p. 085401, 2016. 2
- 175 [20] S. Manik, L. M. Saini, and N. Vadera, “Counting and classification of white blood cell using artificial neural network (ann),” in  
176 *2016 IEEE 1st International Conference on Power Electronics, Intelligent Control and Energy Systems (ICPEICES)*, pp. 1–5,  
177 July 2016. 2
- 178 [21] J. Zhang, H. Hu, S. Chen, Y. Huang, and Q. Guan, “Cancer cells detection in phase-contrast microscopy images based on  
179 faster r-cnn,” in *2016 9th International Symposium on Computational Intelligence and Design (ISCID)*, vol. 1, pp. 363–367,  
180 Dec 2016. 2
- 181 [22] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, “Dermatologist-level classification of skin  
182 cancer with deep neural networks,” *Nature*, vol. 542, no. 7639, pp. 115–118, 2017. 2

- 183 [23] W. Yu, J. Chang, C. Yang, L. Zhang, H. Shen, Y. Xia, and J. Sha, "Automatic classification of leukocytes using deep neural  
184 network," in *2017 IEEE 12th International Conference on ASIC (ASICON)*, pp. 1041–1044, Oct 2017. 2
- 185 [24] H. Lee, A. Battle, R. Raina, and A. Y. Ng, "Efficient sparse coding algorithms," in *In NIPS*, pp. 801–808, NIPS, 2007. 2, 6
- 186 [25] H. Lee, C. Ekanadham, and A. Y. Ng, "Sparse deep belief net model for visual area v2," in *Advances in Neural Information  
187 Processing Systems 20*, MIT Press, 2008. 2, 6
- 188 [26] C. Poultney, S. Chopra, and Y. Lecun, "Efficient learning of sparse representations with an energy-based model," in *Advances  
189 in Neural Information Processing Systems (NIPS 2006)*, MIT Press, 2006. 2, 6
- 190 [27] K. Yu, T. Zhang, and Y. Gong, "Nonlinear learning using local coordinate coding," in *Advances in Neural Information Pro-  
191 cessing Systems 22* (Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, eds.), pp. 2223–2231, 2009. 2,  
192 6
- 193 [28] C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, eds., *Advances in Neural Information Processing  
194 Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec,  
195 Canada*, 2015. 2
- 196 [29] H. Chang, Y. Zhou, A. Borowsky, K. E. Barner, P. T. Spellman, and B. Parvin, "Stacked predictive sparse decomposition for  
197 classification of histology sections," *International Journal of Computer Vision*, vol. 113, no. 1, pp. 3–18, 2015. 6