



OPEN

A large dataset of white blood cells containing cell locations and types, along with segmented nuclei and cytoplasm

Zahra Mousavi Kouzehkanan^{1,2,11}, Sepehr Saghari^{2,3,11}, Sajad Tavakoli^{2,4}, Peyman Rostami^{2,5}, Mohammadjavad Abaszadeh¹, Farzaneh Mirzadeh^{2,6}, Esmail Shahabi Satsar^{2,7}, Maryam Gheidishahran⁸, Fatemeh Gorgi⁹, Saeed Mohammadi¹⁰ & Reshad Hosseini¹✉

Accurate and early detection of anomalies in peripheral white blood cells plays a crucial role in the evaluation of well-being in individuals and the diagnosis and prognosis of hematologic diseases. For example, some blood disorders and immune system-related diseases are diagnosed by the differential count of white blood cells, which is one of the common laboratory tests. Data is one of the most important ingredients in the development and testing of many commercial and successful automatic or semi-automatic systems. To this end, this study introduces a free access dataset of normal peripheral white blood cells called Raabin-WBC containing about 40,000 images of white blood cells and color spots. For ensuring the validity of the data, a significant number of cells were labeled by two experts. Also, the ground truths of the nuclei and cytoplasm are extracted for 1145 selected cells. To provide the necessary diversity, various smears have been imaged, and two different cameras and two different microscopes were used. We did some preliminary deep learning experiments on Raabin-WBC to demonstrate how the generalization power of machine learning methods, especially deep neural networks, can be affected by the mentioned diversity. Raabin-WBC as a public data in the field of health can be used for the model development and testing in different machine learning tasks including classification, detection, segmentation, and localization.

The issue of precise and early diagnosis is the most vital step in the medical treatment process. According to the World Health Organization, about 2 billion people currently do not have access to primary medical and pharmaceutical services¹. In the meantime, laboratory tests play an essential role in the diagnosis and treatment of diseases. It is estimated that about 70% of the decisions related to the diagnosis and treatment of disease, as well as the discharge and admission of a patient, rely on the results of laboratory tests². In this regard, the differential count of white blood cells is one of the common laboratory tests necessary to be considered in the diagnosis of various diseases such as blood disorders (such as leukemia, anemia, polycythemia, etc.) and immune system-related diseases (such as autoimmune anemias, allergy, etc.)³.

White blood cells called leukocytes include two groups of phagocytes and lymphocytes. While phagocytes comprise cells of the innate immune system and function rapidly after infection, lymphocytes mediate the acquired immune response. Phagocytes, themselves, can be divided into granulocytes (neutrophils, basophils, and eosinophils) and monocytes. In Table 1 and Fig. 1, you can see the characteristics and images of the five categories of white blood cells. Table 2 shows some examples of the diseases that occur with an increase or

¹School of ECE, College of Engineering, University of Tehran, Tehran, Iran. ²Nimaad Health Equipment Development Company, Tehran, Iran. ³Graduated Bachelor of Laboratory of Sciences, Paramedical Faculty of Guilan, University of Medical of Sciences, Langarud, Gilan, Iran. ⁴Faculty of Electrical Engineering, K. N. Toosi University of Technology, Tehran, Iran. ⁵School of Mechanical Engineering, Sharif University of Technology, Tehran, Iran. ⁶School of Medicine, Tarbiat Modares University, Tehran, Iran. ⁷Flow Cytometry Department, Takhte Tavous Patobiology Lab, Tehran, Iran. ⁸Department of Hematology and Blood Transfusion, School of Allied Medical Sciences, Iran University of Medical Sciences, Tehran, Iran. ⁹Bachelor of Laboratory of Sciences, Faculty of Paramedical, Mashhad University of Medical Sciences, Mashhad, Iran. ¹⁰Hematology-Oncology and Stem Cell Transplantation Research Center, Tehran University of Medical Sciences, Tehran, Iran. ¹¹These authors contributed equally: Zahra Mousavi Kouzehkanan and Sepehr Saghari. ✉email: reshad.hosseini@ut.ac.ir

WBCs	% In blood	Nucleus	Cytoplasm	Size (μm)
Neutrophils	60%	It is divided into 2 to 5 segments and stains dark purple (multi-lobed)	It is pale pink to tan with pink-purple granules	12–16
Eosinophils	3%	It is blue and is divided into 2 segments	It is full of pale pink tan with large orange and red granules	14–16
Basophils	1%	It has 2 lobes that each stains purple, and is difficult to be seen	It is pale pink-tan but contains large purple/blue-black granules which obscure the cell nucleus	14–16
Monocytes	6%	It is singular and is kidney shaped (convoluted shape), bean shaped or horseshoe shaped with deep indentation	It stains a blue-gray color and is "ground glass" with tiny granules, Vacuoles are sometimes present in it	14–20
Lymphocytes	30%	It is large, round or oval, and is dark staining	It is not present or very small, and is pale blue in color, and occasionally has purple-reddish granules	8–15

Table 1. Characteristics of white blood cells⁴.

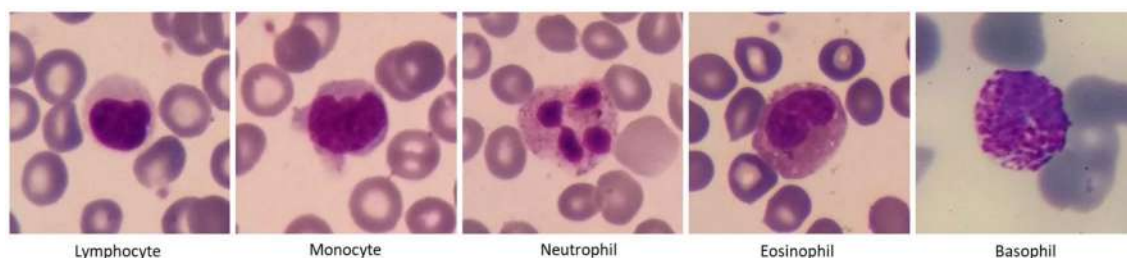


Figure 1. Five types of white blood cells in the normal peripheral blood.

White blood cell	Increase	Decrease
Lymphocyte	Acute and chronic leukemia, hypersensitivity reaction, viral infection	AIDS, influenza, sepsis, aplastic anemia
Monocyte	Autoimmune disease, fungal and protozoan infection	Aplastic anemia, hairy cell leukemia, acute infections
Neutrophil	Chronic inflammation, Infection	Chediak-Higashi syndrome, Kostman syndrome, Auto-immune neutropenia
Eosinophil	Allergic reaction, parasitic infection, malignancy	Cushing syndrome, shock or trauma driven stress
Basophil	Leukemias	Hyperthyroidism and acute infections

Table 2. White blood cells alterations and related different diseases⁵.

decrease in the number of white blood cells. For example, in allergic diseases, the number of basophils increases, or in blood malignancies, we can see an increase in the number of precursors of blood cells and changes in their shape and size. Therefore, determining the correct type and number of white blood cells is very important for diagnosing various diseases.

At present, manual (microscopic evaluation) and automated methods (using automatic hematology devices) are used to evaluate blood cells. Automated methods include devices which evaluate blood cells based on light scattering or electrical impedance such as Sysmex XP-300, Nihon Kohden Blood Cell Counter, and DH36 3-Part Auto Hematology Analyzer.

In electro-optical analyzers, a light-sensing detector measures the optical scattering. The size of the detected pulses corresponds to the size of the blood cells. Furthermore, in electrical impedance or Coulter principle cell counter, the passage of cells through an aperture in which an electric current is applied causes change in the electrical resistance. Pulses the height of which corresponds to the volume of the cell are counted, and this is considered as the basis of Coulter's principle working⁴. Besides these methods, microscopic hyperspectral imaging technology, as an emerging imaging modality, is currently being used. This method is a combination of spectroscopy and 2D imaging^{6–8}.

One of the serious drawbacks of these devices apart from their high cost is the simple act of counting cells without them being evaluated qualitatively from a structural and morphological point of view. As a result, after evaluating the blood sample by the mentioned cell counters, it is necessary to prepare a smear and evaluate it microscopically by the laboratory staff to achieve an accurate and correct diagnosis.

On the other hand, issues such as the lack of specialists and laboratory equipment, heavy workload, inexperience, and incorrect diagnosis affect the test results. Misdiagnosis affects the treatment regime, and consequently, can result in the malpractice and an increase of associated costs. However, the use of new technologies such as artificial intelligence and image processing allows quantitative and qualitative evaluations to improve the quality of diagnosis⁹.

Over the past 20 years, the techniques for automated imaging of the blood-stained slides have been introduced by computer-connected microscopes capable of assessing blood cell morphology. With the development

of technology, companies such as Cellavision, Westmedica, Siemens, etc. have made it possible to differentiate the count of normal from abnormal blood cells¹⁰. In fact, today, deep neural networks are one of the most widely used machine learning methods for the classification and segmentation of medical images. Shahin, A et al.¹¹ used DNNs to classify white blood cells. In addition, these networks are used for the classification of red blood cells to detect a sickle cell anemia¹². Deep neural networks are also used for the segmentation of the pancreas in the CT scan images^{13,14} and the segmentation of the MRI images¹⁰.

Data have the most important role in the development of machine learning models. In order to train deep neural networks and increase their generalizability, we need a lot of diverse precise data and confident labels. The process of labeling medical data should be carried out by professionals and is, therefore, a time-consuming and challenging procedure. As a result, medical databases are of high significance in smartening medical diagnoses. Unfortunately, researchers, today, have limited access to a variety of medical data for various reasons. Examples of available medical image databases are¹⁵ and¹⁶. The database¹⁵ contains 82 3D CT scans in which the Grand Truths of the pancreas for all slices were manually extracted by medical students and finalized by a specialist radiologist. Camelyon¹⁶ is another dataset with 1399 whole-slide images of the lymph node smear samples with and without metastases, for which the labels were checked twice.

The morphological diversity of white blood cells is very high and in some cases, it is very challenging, even for an expert, to distinguish some classes from each other. On the other hand, many artificial intelligence articles have adopted two approaches to evaluate their proposed method regarding segmentation and classification of white blood cells: They have either collected small databases to the best of their ability^{17–20} or used the small databases available^{21–23}. Therefore, a database with a large amount of diverse data and reliable labelling is truly necessary to evaluate and compare different methods with each other. Such a reference database will allow more artificial intelligence scientists to enter the field and will help the advancement of intelligence differentiation of white blood cells. The most important characteristics of the Raabin-WBC dataset that distinguishes it from similar datasets are as follows:

- **Large number of data:** We tried to collect as much data as possible for each class in order for them to be appropriate for all machine learning techniques, especially deep learning. (Approximately 40,000 white blood cell images)
- **Precise labels:** We considered more detailed labels than five types of white blood cells. In fact, labels contain the most important subgroup of each type. For example, we considered the meta and band which are subgroups of neutrophils and are valuable in diagnosis. In the next section, more information about the labels will be presented.
- **Double labeling:** For more insurance, most of the cells are labeled by two experts.
- **Free public access:** Since we aim at helping the development of artificial intelligence in hematology, the Raabin-WBC dataset is freely available for all.
- **Data cleaning:** In the process of data collection, the existence of duplicate cell images is not inevitable. The first problem is that the duplicate cell images are not exactly the same. For example, there is a possibility of the cell being somewhat moved. The second issue is that having more than two versions of one cell image is also possible. Hence, we developed a fast graph-based image processing method that can accurately remove as many duplicate cells as possible. Despite this, it is still probable for some duplicate images to exist, albeit being significantly different.
- **The ground truths of nuclei and cytoplasm:** The ground truths of the nuclei and cytoplasm are extracted for 1145 selected cells. In order to extract the ground truth of nuclei, we developed a toolkit that by using image processing techniques makes the ground truth extraction process much easier.
- **Diversity of the microscope and camera:** Although most of the data were collected by a fixed type of microscope and camera, we collected some data with another type of camera and microscope, as well. In the section of experiments, you will see how new test data help to evaluate the generalization power of our trained models. In other words, the diversity of the dataset assists us in selecting a model that has correctly learned the manifold of cell images.

The rest of the paper is as follows: In “[The characteristics of Raabin-WBC](#)” section, we will elaborate more on the details regarding the dataset. In “[Data collection](#)” section, the data collection process will be explained completely. In “[Experiments](#)” section, we will do some machine learning experiments and discuss the generalization power of the models.

The characteristics of Raabin-WBC

In this section, more information is provided about the Raabin-WBC dataset. About 73 peripheral blood films were used for collecting this dataset. After imaging stained blood films, we tried to mine the most possible useful information from the raw data. For instance, the bounding box of all white blood cells and artifacts were extracted, cropped and labeled, successively. It is worth noting that a significant number of WBCs and artifacts were labeled by two experts. Furthermore, we provided the ground truth of the nucleus and cytoplasm for some of the cropped cells. The full details of the data collection steps are explained in Sect. 3. In Table 3, some general and useful information of the Raabin-WBC dataset is provided. Note that these numbers have been computed after the cleaning phase.

Labels. In the Raabin-WBC dataset, more detailed labels are considered than just five general types of white blood cells. For example, besides the mature neutrophil, we have evaluated two other ancestors of this white blood cell: Metamyelocytes and Band. An increase in the number of band forms and metamyelocytes is one of

Number of all films (smear)	73
Number of CML films	1
Number of normal-anemia films	2
Number of normal-eosinophilia films	2
Number of normal films	68
Number of microscopic large images	20,936
Number of bounding boxes (including WBCs and artifacts)	40,763
Number of 0 labeled WBCs	10,385
Number of 1 labeled WBCs	4971
Number of 2 labeled WBCs	25,408
Number of ground truths for lymphocytes (including nucleus and cytoplasm)	242
Number of ground truths for monocytes (including nucleus and cytoplasm)	242
Number of ground truths for neutrophils (including nucleus and cytoplasm)	242
Number of ground truths for eosinophils (including nucleus and cytoplasm)	201
Number of ground truths for basophils (whole cell)	218

Table 3. Raabin-WBC information table.

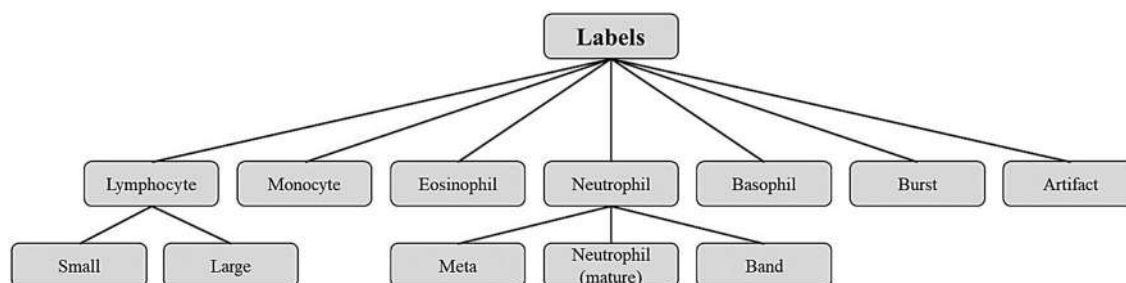


Figure 2. Diagram of labels in the Raabin-WBC dataset.

the features of reactive neutrophilia (an increase in the number of circulating neutrophils to levels greater than $7.5 \times 10^9/L$)⁵. In addition, lymphocytes are divided into small (the main agents of the acquired immune system including B and T cells) and/or activated lymphocytes (activated small lymphocytes referred to as large lymphocytes or lymphoblasts). Bursts refer to smudge cells that are leukocyte remnants formed during blood smear preparation. Beside the leukocytes we considered drying artifacts as new labels, because artifacts are commonly seen after staining the samples. In Fig. 2, the diagram of the labels is presented.

In Table 4, the number of labels associated with two experts is shown. The rows and columns of the table belong to the first and second experts, respectively, noting that 9015 cells have not been labeled yet. We asked our experts to label the cells as unrecognizable if they had any doubts. Indeed, we have 1099 cells labeled as not recognized by the two experts. In Table 4, you can see the amount of disagreement for each pair of different labels (Non-diagonal elements of the matrix). For example, large and small lymphocytes are confused a lot. Also, seem bands have often been mistaken with mature neutrophils. Other examples of confusing pairs are artifact and burst, large lymphocyte and monocyte, and small lymphocyte and burst. The high numbers in the rows and columns labeled as not recognized indicate that it is very challenging to identify the type of white blood cell.

Data structure. The Raabin-WBC dataset consists of images that were taken from blood films (similar to Fig. 5). Corresponding to each microscopic image, a dictionary (.json format) file containing the following information about that image was provided:

- Information about the blood elements in the image including their coordinates and labels. Most of the elements are labeled by two experts.
- Information about the blood smears including staining method and the type of the disease. Note that all blood smears have been prepared from normal samples. Only a Chronic Myeloid Leukemia (CML) sample has been used to extract basophils.
- Information about the microscope includes the type of microscope and its magnification size.
- The type of camera used.

There is also a subset of the database called double-labeled Raabin-WBC which includes cropped images of the five main types of WBCs and were labeled the same by both of the experts. We will explain more about this sub-dataset in the experiment section.

	Artifact	Band	Basophil	Burst	Eosinophil	Large lymph	Meta	Monocyte	Neutrophil	Small lymph	Not recognized	Not labeled
Artifact	3489	0	0	14	2	1	0	0	6	4	96	225
Band	0	311	0	2	0	0	2	0	32	0	16	71
Basophil	0	0	308	0	0	0	0	0	0	0	2	0
Burst	29	0	0	2673	1	11	0	4	1	32	96	525
Eosinophil	0	0	0	9	1466	0	0	0	1	1	13	607
Large lymph	0	0	0	1	0	2153	0	4	1	172	23	163
Meta	0	0	0	1	0	0	11	1	2	0	6	12
Monocyte	0	0	0	2	0	24	0	874	1	0	36	104
Neutrophil	0	134	0	29	3	1	0	2	11,726	1	109	1078
Small lymph	1	0	0	1	0	31	0	0	0	1833	20	370
Not recognized	65	5	0	9	5	744	10	81	127	332	1099	268
Not labeled	5	0	0	0	0	1	0	0	9	4	3	9015

Table 4. The number of labels associated with two experts. The rows and columns belong to the first and second experts, respectively.

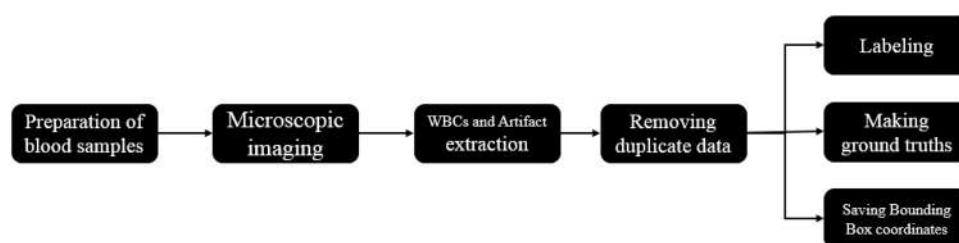


Figure 3. Main steps of the Raabin-WBC dataset collection.

Data collection

The data were collected from patients and the ethical approval was gotten from the ethical committee of Hematology, Oncology, and SCT Research center of Shariati hospital. We confirm that all methods were performed in accordance with the relevant guidelines and regulations. The steps of data collection (Fig. 3) include preparing blood smears and photographing them, extracting the bounding box of white blood cells, data cleaning, and finally labeling the data and extracting ground truths. More details are explained in the rest of this section.

Preparation of blood smears and imaging. 72 normal peripheral blood films (male and female samples from ages 12 to 70) have been used to collect neutrophils, eosinophil, monocyte, and lymphocyte images. On the other hand, due to the very low presence of basophils in normal specimens (<1–2%)⁹, basophils of one CML-positive sample have been imaged. Owing to the widespread use of Giemsa in medical labs⁹, all samples were stained by Giemsa. It should be noticed that all samples were taken from collaborator medical laboratories (Razi Hospital in Rasht, Gholhak Laboratory, Shahr-e-Qods Laboratory, and Takht-e Tavous Laboratory in Tehran, Iran) and we did not deal directly with patients. It is worth noting that in Iran, it is necessary to get the approval of patients only in clinical trials. But in retrospective studies, this is not necessary in Iran. The process of imaging the slides was performed by the help of two types of microscopes, namely Olympus CX18 and Zeiss at a magnification of 100×. Since determining the Diff area to evaluate and count different types of white blood cells is of utmost importance, an expert lab staff had supervised the cell imaging process.

With smart phones being widely used in society, a rapidly growing trend has emerged to adapt them to medical diagnostics^{24,25}. The availability, ease of use and low cost of high-pixel density cameras available in smart phones make them widely used in various science fields^{26–32}. Therefore, in compiling this database, the cameras available on smart phones have been used, the details of which are given in Table 5. Smartphones can be adapted for microscopic imaging using some accessory equipment^{33–35}. To facilitate the use of smart phones in microscopic imaging in this dataset, an adapter was designed and made by 3D printing to mount the smart phone on the microscope ocular lens (Fig. 4). The designed adapter has somewhat managed to minimize the drawbacks of the commercial models available in the market such as restrictions on the size of the phone and ocular lenses, as well as the difficulty of the adjustment.

Extraction of white blood cells from images. In total, about 23,000 images were taken from blood films. There exist many red blood cells in each blood smear image. It is also probable that one or more other blood elements such as white blood cells and sometimes color spots exist in the image. The bounding box of these blood elements should be somehow identified. For this purpose, two approaches have been considered. Due to the distinct color of the nucleus in white blood cells, in the first approach, several white blood cells were

Smartphone	Release date	Sensor model	Sensor type	No. of pixels	Aperture	Sensor size	Pixel size
Samsung Galaxy S5	2014	Samsung S5K2P2XX ISOCELL	CMOS	16 MP	f/2.2 31 mm	1/2.6"	1.12 μm
LG G3	2014	Sony IMX135 Exmor RS	CMOS	13 MP	f/2.4 29 mm	1/3"	1.12 μm

Table 5. Smartphone camera specifications used for data collecting.

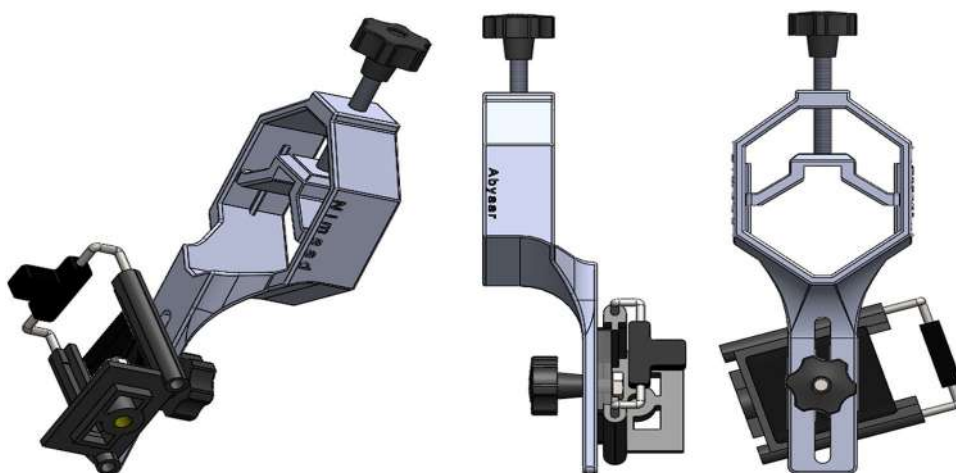


Figure 4. Designed adapter to mount smart phones on the ocular lens of a microscope to make the act of capturing the photos from the samples quicker and easier. Experts work with a microscope manually and see the images on the mounted smartphone and take photos.

extracted manually as Grand Truth data, and a color filter was trained to separate the white blood cells from the background. The aforementioned color filter was applied to the main images, and the approximate position of the white blood cells was marked. Finally, a 512 by 512 square with the center of the cell is considered as a bounding box. In the second approach, extracted bounding boxes with the help of the first approach were used, and a Faster RCNN network³⁶, which can determine the exact location of the white blood cells in the original image, was trained. Eventually, about 43,000 blood elements were obtained.

Data cleaning. In the process of imaging from the blood smears, a white blood cell may be placed in more than one image (Fig. 5). Therefore, duplicate cell images exist among cropped images. The major problem is that the two images of one cell are not necessarily very similar. Consequently, a simple mean square error on the value of the pixels is not enough to detect duplicate cell images. Indeed, a cell can be repeated more than twice. In Fig. 6, an example of three images of one cell is represented. As you can see, the qualities of the three images are different.

Manual comparison of these images in pairs is practically impossible. Hence, an artificial intelligence algorithm, fast and accurate, has been developed to remove duplicate cell images. We used the Python ImageHash library, in this regard. First, for all pairs of cropped images, the Average Hash (AHash) and Perceptual Hash (PHash) values are calculated very quickly. Paired images, the AHash and PHash distances of which are less than those of the specific thresholds, are the same, and one of them should be removed. The thresholds of the Average Hash and Perceptual Hash are set manually through trial and error (See appendix 1 for more details).

Since an image may exist more than twice, a two-by-two comparison is not sufficient. For this purpose, a solution to the problem is presented from the Graph's point of view. In fact, we have a graph with N nodes (N is the number of cropped cell images from blood film). There exist edges between the nodes that satisfy the same-ness condition. In this case, the connected components of the graph form equal images. Connected components of a graph can be calculated with the help of the breadth-first search algorithm very swiftly (See appendix 2 for more details). If a connected component has $n > 1$ images, $n - 1$ of them must be removed. To enhance the quality of the database, the image with the highest resolution remains out of n images, and the rest are deleted. The OpenCV³⁷ library is used to compare the resolution of images. In this regard, Sobel horizontal and vertical filters³⁸ are applied to the images and the gradient magnitude is calculated for each pixel. Finally, the image with the highest average gradient magnitude is selected, because it is the sharpest one.

As described in Sect. 2, to offer full information, we provide our data in the format of large and not cropped images (like Fig. 5). For each large image, the coordinates and the labels of the containing cells are provided. We tried to remove as many duplicates as possible from large images. Indeed, we remove a large image in which all containing cells are inside another image. For example, in Fig. 5, image b is removed.

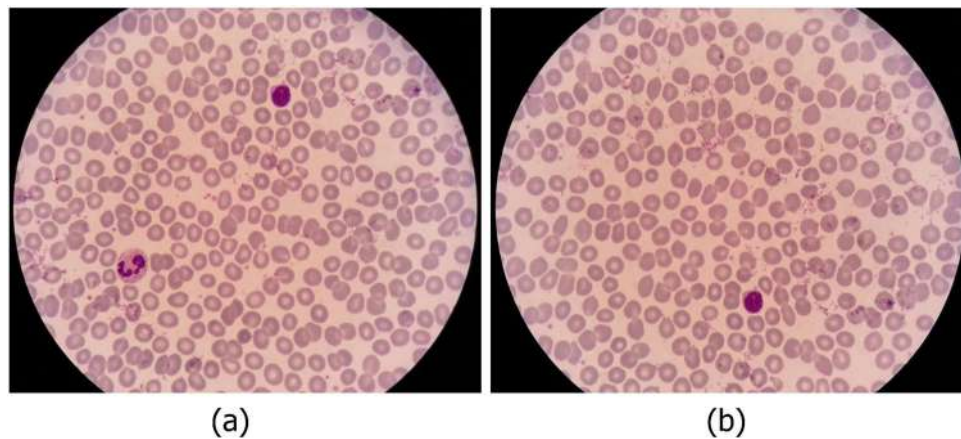


Figure 5. An example of two overlapped microscopic images.

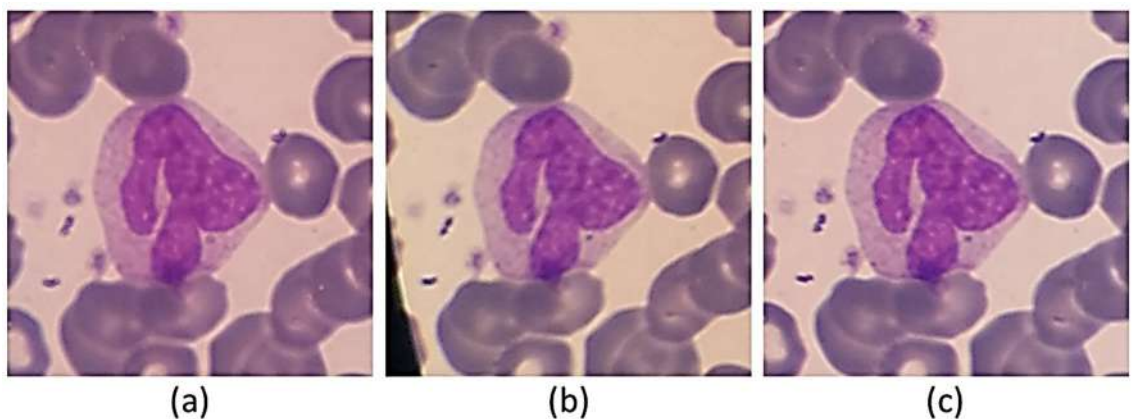


Figure 6. One sample that had been repeated three times.

Labeling process

This section describes the labeling process, which involves determining the cell types and the ground truth of the nucleus and cytoplasm. As you can see in Table 1, the characteristics of the nucleus and cytoplasm can significantly affect determining the type of the cell. Some papers^{19,39} extract different features from the nucleus and the cytoplasm to classify white blood cells. These features usually describe the shape and the color of the nucleus and the cytoplasm.

Cell type labeling. For labeling cells, two applications were developed for Android (Fig. 7). One application is for labeling cropped cells (Fig. 7-part b), and the other is for selecting the location and type of each cell (Fig. 7-part a). Furthermore, a desktop application with the help of the Python Tkinter library⁴⁰ was developed for manually selecting the location and type of the cells (Fig. 8). It is worth mentioning that most of the images were labeled by two experts.

Ground truths of the nucleus and the cytoplasm. In recent years, many researchers have developed segmentation algorithms for the cytoplasm and nucleus of the white blood cells^{3,18–21,23}. Hence, we tried to prepare the ground truths of the cytoplasm and the nucleus for a proper number of cropped white blood cells. For this purpose, 1145 cropped images including 242 lymphocytes, 242 monocytes, 242 neutrophils, 201 eosinophils, and 218 basophils were randomly selected, and their ground truths were extracted by an expert. It is worth mentioning that we only prepared the ground truth of the whole cell for basophils, and we were not able to produce the ground truths of the nucleus and cytoplasm for basophils. This is because the basophils are usually covered by very purple granules, and the border between cytoplasm and nucleus is not easily visible. Figure 9 shows some samples of the cells along with their ground truths.

To produce the ground truths of nuclei, a newly published software called Easy-GT⁴¹ was employed. This software has been developed to extract the ground truths of nuclei. In Easy-GT software, a nucleus is determined by a relatively accurate segmentation method, and if necessary, the user can adjust the ground truth of the nucleus by modifying the final threshold⁴¹ (Fig. 10). In the segmentation process, the RGB image is first color-balanced⁴¹ and converted to the CMYK color space. Secondly, the two-class Otsu's thresholding algorithm⁴² applied to the

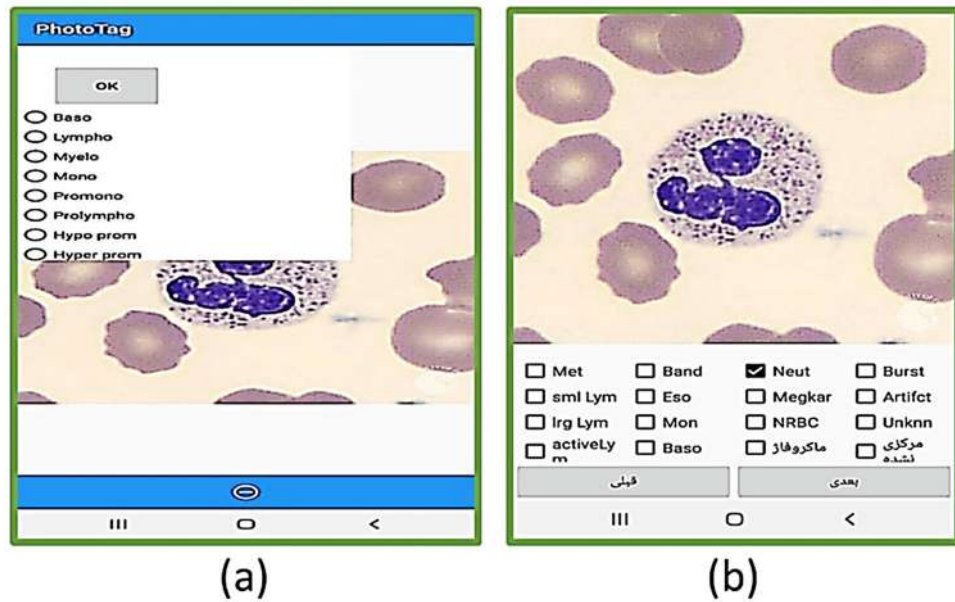


Figure 7. The user interface of the two android applications that were designed to selecting and labeling the white blood cells.

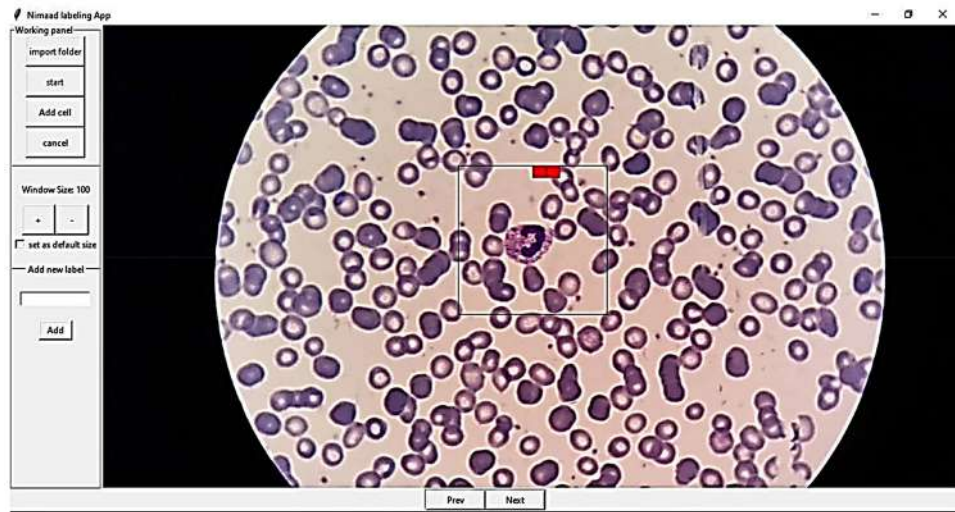


Figure 8. The user interface of the desktop application designed for labeling white blood cells.

M channel gives us a threshold (th^{2class}). Again, the three-class Otsu thresholding algorithm is applied to the M channel and the two lower and upper thresholds (th_{low}^{3class} , th_{up}^{3class}) are extracted. Finally, the ultimate threshold value is obtained by computing the convex combination of th^{2class} and th_{up}^{3class} .

To make the ground truth of the cytoplasm, a light pen was used, and the ground truth of the whole cell was specified by an expert. Finally, by removing the nucleus part obtained from Easy-GT, only the cytoplasm remains.

Experiments

In this section, we are going to do some machine learning experiments on the Raabin-WBC data. Due to the diversity of information in the database, many research lines can be developed. Yet, we consider the most common possible experiment. We classify five classes of white blood cells, and we leave the rest to those who are interested in this field. For this purpose, we used the double-labeled cropped cells and considered only five main classes including mature neutrophils, lymphocytes (small and large), eosinophils, monocytes, and basophils. We called this sub-dataset Double-labeled Raabin-WBC. In the following, we will compare this database with some

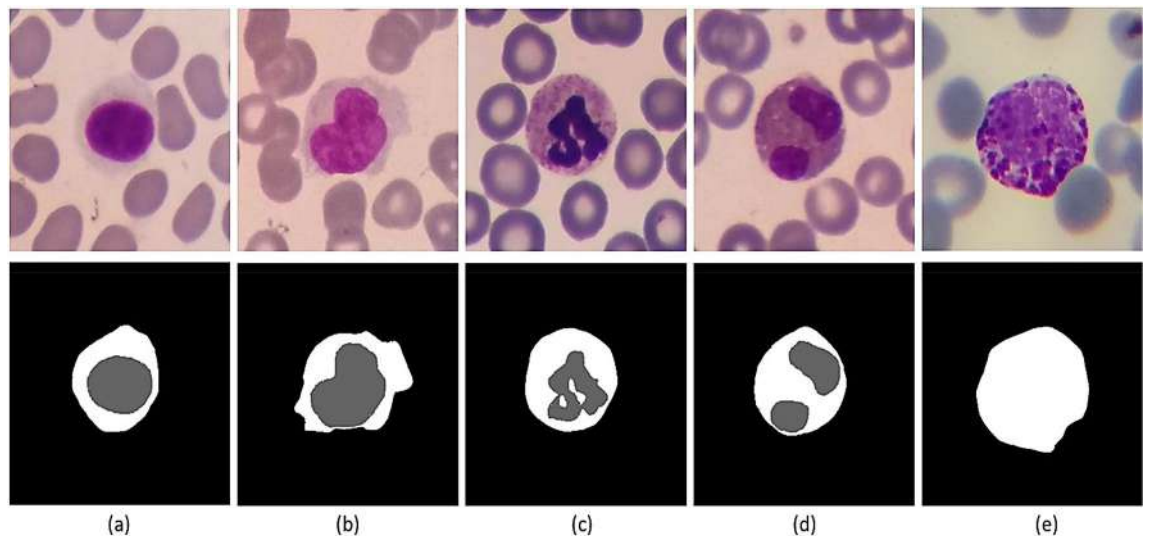


Figure 9. Some samples of ground truths provided in the Raabin-WBC dataset. First row contains the original cropped images of white blood cells. Second row contains the ground truths of some nuclei and cytoplasm. The columns (a), (b), (c), (d), and (e) show lymphocyte, monocyte, neutrophil, eosinophil, and basophil, respectively.

existing 5-class databases and train some deep popular neural networks. We will also discuss the generalization power of the models.

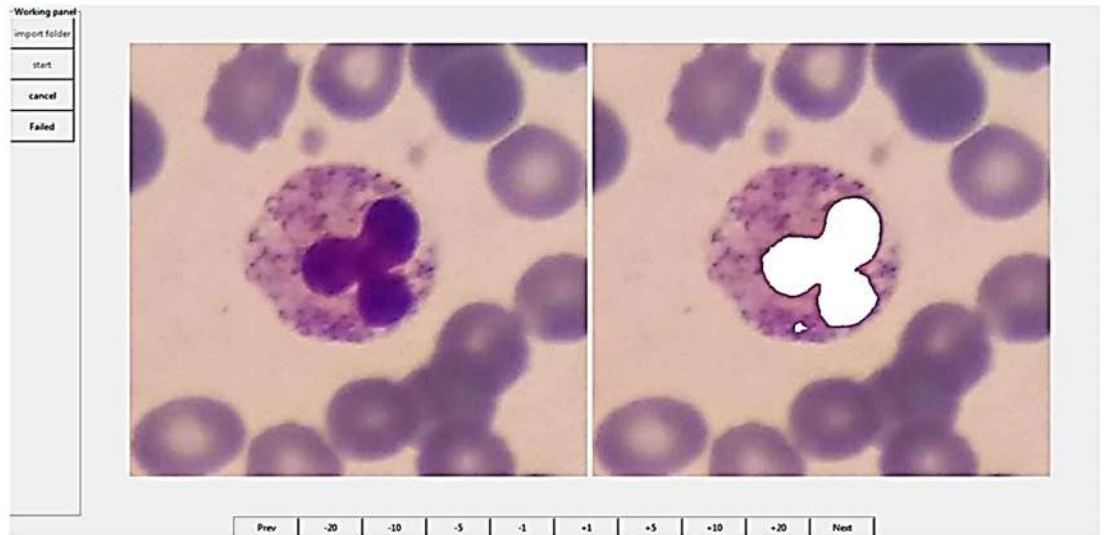
A comparison with similar datasets. Various datasets of normal peripheral blood with different properties exist, but in general, most of them have a small number of samples. This is due to the fact that in the medical field, data collection and labeling are complicated. On the other hand, in the field of Hematology, artificial intelligence models are usually sensitive to some specifications of the dataset such as the number of data, the staining technique, the microscope and camera used, and the magnification. So, by altering the aforementioned characteristics, the accuracy of the models may be reduced. In Table 6, the characteristics of some datasets are presented and compared with Double-labeled Raabin-WBC. As you can see, our database is far better in several ways including data number, label assurance, ground truth, camera, and microscope variety. Most importantly, this database is available to everyone for free.

Utilized models. Some popular pre-trained deep neural networks were trained on Double-labeled Raabin-WBC to classify five types of white blood cells. VGG16⁴³ is the oldest CNN model consists of alternating convolutional and pooling layers. From deep residual network families, Resnet18⁴⁴, Resnet34⁴⁴, Resnet50⁴⁴, and Resnext50⁴⁵ were tested. In Resnet architecture, identity shortcut connections that skip one or more layers are used⁴⁴. Resnext is an extension of Resnet in which the residual block is replaced by a new aggregation component⁴⁵. In mentioned aggregation component, the input feature map is projected to some lower-dimensional representations, and their outputs are aggregated⁴⁵. Another CNN used in the experiments is DenseNet121⁴⁶ which consists of dense blocks. At each dense block, each layer is fed from all previous layers, and its outputs are transferred to all next layers.

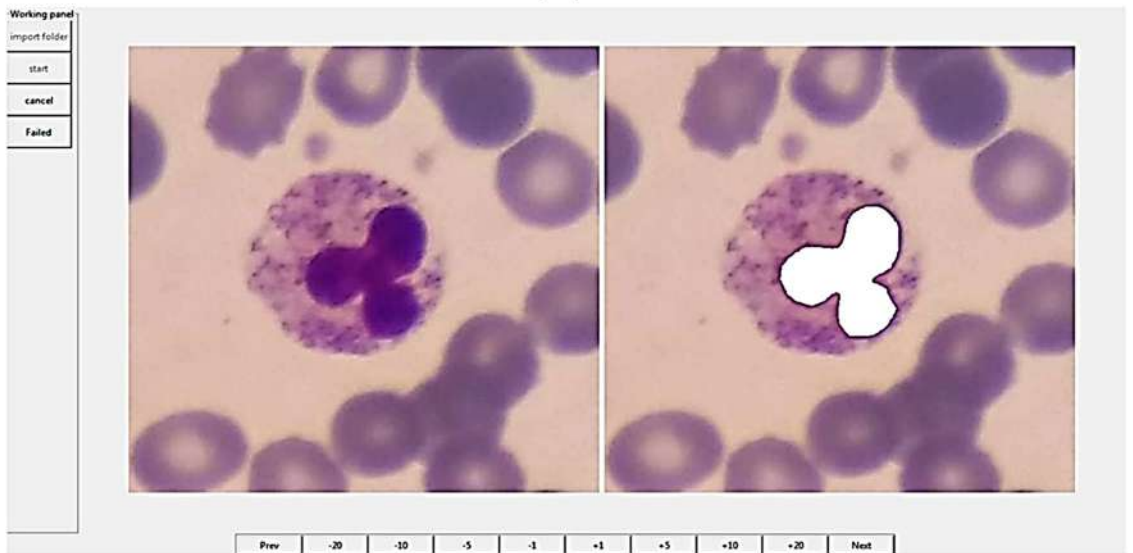
Another tested deep architecture is MobileNet-V2⁴⁷ which is suitable for mobile devices. The building block of MobileNet-V2 is an inverted residual block, and non-linearities are removed from narrow layers. MnasNet1⁴⁸ and ShuffleNet-V2⁴⁹ are other light-weight CNNs for mobile devices. In MnasNet, reinforcement learning is employed to find an efficient architecture⁴⁸. In ShuffleNet-V2⁴⁹ at the beginning of the basic blocks, a split unit divides the input channels into two branches, and at the end of the block, concatenation and channel shuffling occur. Besides the aforementioned neural networks, we also utilized a feature-based method⁵⁰ in which the nucleus was segmented at first, and its convex hull was then obtained. After that, shape and color features were extracted using the segmented nucleus and its convex hull. Finally, WBCs were categorized by an SVM model.

Classification results. The generalization power of the models described in the former section is to be examined at two levels. For this purpose, we split data into three groups of training data, test-A, and test-B, the properties of which can be observed in Table 7. The quality of the images in the test-A dataset is similar to that of the training dataset, but the images in the test-B dataset have different qualities in terms of camera type and microscope type. Unfortunately, the test-B data only contains double-labeled neutrophils and lymphocytes.

The training data are not balanced, in other words, the number of cells in each class is imbalanced. Hence, the training set was augmented and moderated using augmentation methods such as horizontal flip, vertical flip, rescaling, and a combination of them. In order to evaluate the models, four metrics are considered for each class:



(a)



(b)

Figure 10. The user interface of Easy-GT software⁴¹. This software was developed for extracting the ground truths of nuclei in white blood cells.

precision (P), sensitivity (S), F1-score, and accuracy (Acc). The aforementioned criteria are obtained through the Eqs. (1), (2), (3), and (4).

$$Precision = \frac{TP}{TP + FP} \tag{1}$$

$$Sensitivity = \frac{TP}{TP + FN} \tag{2}$$

$$F1 = 2 \times \frac{Prec \times Sens}{Prec + Sens} \tag{3}$$

$$Acc = \frac{TP + TN}{TP + FP + TN + FN} \tag{4}$$

In Eqs. 1–4, TP, FP, TN, and FN are true positive, false positive, true negative, and false negative, respectively. In Tables 8 and 9, the results on the test-A and test-B datasets are presented. Also, in the last row of Tables 8 and

Dataset	Number of WBCs						Access	Staining	Microscope and zoom	Camera	Label	Ground truths	
	Lymp	Mon	Neut	Eos	Bas	Total						Nucleus	Cytoplasm or whole cell
LISC ²⁰	59	55	56	42	54	266	Public	Gismo-right	Axioskope40 Zoom : 100X	Sony-SSCD-C50AP	One expert	266	266
BCCD ¹⁸	33	19	208	86	3	349	Public	Gismo-right	Regular light microscope Zoom : 100X	CCD color camera	One expert	×	×
Hegde et al. ³⁹	33	23	30	22	14	122	Private	Leishman	OLYMPUS CX31 Zoom : 100X	N/A	One expert	122	×
MISP ¹⁹	36	33	38	42	0	149	Public	N/A	Canon optical microscope Zoom : 100X	Canon V1	One expert	×	×
ALL-IDB ¹⁷	60	3	18	2	1	84	Public	N/A	N/A Zoom : 300X–500X	Canon Power-Shot G5	N/A	×	×
Zheng et al. ²¹ (CellaVision)	37	18	30	12	3	100	Public	N/A	N/A Zoom : N/A	N/A	One expert	100	100
Zheng et al. ²¹	53	48	176	22	1	300	Public	A newly developed method ²¹	N800-D motorized autofocus Zoom : N/A	Motic moticam pro 252A	One expert	300	300
Double-labeled Raabin-WBC	3609	795	10,862	1066	301	17,965	public	Giemsa	1. Olympus Cx18 2. Zeiss microscope Zoom : 100	1. Camera phone Samsung galaxy S5 2. Camera phone LG G3	Two experts	1145	1145

Table 6. Comparing some datasets with double-labeled Raabin-WBC. Double-labeled Raabin-WBC does not contain the repeated samples as well as includes only five general cells (lymphocyte, monocyte, neutrophil, eosinophil, and basophil).

Sets	Lymph	Mono	Neut	Eos	Bas
Training data	2427	561	6231	744	212
Test-A	1034	234	2660	322	89
Test-B	148	0	1971	0	0

Table 7. The number of samples in training data, test-A, and test-B.

Methods	Lymph			Mono			Neut			Eosi			Baso			Acc (%)
	P (%)	S (%)	F1 (%)	P (%)	S (%)	F1 (%)	P (%)	S (%)	F1 (%)	P (%)	S (%)	F1 (%)	P (%)	S (%)	F1 (%)	
ResNet18 ⁴⁴	98.84	99.23	99.08	96.96	95.30	96.12	99.66	99.36	99.5	95.77	98.45	97.09	100	100	100	99.06
ResNet34 ⁴⁴	98.66	99.52	99.09	96.93	94.44	95.67	99.85	99.14	99.49	94.67	99.38	96.97	100	100	100	99.01
ResNet50 ⁴⁴	98.47	99.52	98.99	97.36	94.44	95.88	99.74	99.40	99.57	96.94	98.45	97.69	100	100	100	99.10
ResNext50 ⁴⁵	98.01	100	98.99	99.53	91.45	95.32	99.74	99.55	99.64	97.85	98.76	98.30	100	100	100	99.17
MnasNet1 ⁴⁸	98.93	98.65	98.79	94.14	96.15	95.14	99.66	98.76	99.21	92.15	98.45	95.20	100	100	100	98.59
MobileNet-V2 ⁴⁷	98.85	99.32	99.08	96.93	94.44	95.67	99.66	99.40	99.53	96.97	99.38	98.16	100	100	100	99.12
DenseNet121 ⁴⁶	98.38	99.61	98.99	97.72	91.45	94.48	99.70	99.14	99.42	94.40	99.38	96.82	100	100	100	98.87
ShuffleNet-V2 ⁴⁹	98.09	99.32	98.70	96.49	94.02	95.24	99.74	99.36	99.55	97.85	98.76	98.30	100	100	100	99.03
VGG16 ⁴³	98.26	98.26	98.26	95.09	91.03	93.01	99.43	98.57	99	89.01	98.14	93.35	100	100	100	98.09
Tavakoli et al. ⁵⁰	97.23	95.07	96.14	84.87	86.32	85.59	98	95.60	96.78	72.24	91.30	80.66	96.59	95.51	96.05	94.65

Table 8. The results of different pre-trained models as well as Tavakoli et al.⁵⁰ on the test-A dataset.

9, the results of the feature-based classification presented in the paper⁵⁰ are showed. In Fig. 11, the plots of the accuracy and the loss of training data and validation data related to nine pre-trained models are shown.

The results are surprising, and all methods have an acceptable outcome on the test-A data. Yet, the performance of most of the models on the test-B data experience a dramatic decrease. The feature-based method⁵⁰ had the least performance reduction, despite having the lowest accuracy on the test-A data. Among deep neural networks, the VGG16⁴³ network has relatively more generalizability. It can be said that the feature-based method could extract more meaningful features from cell images than the deep neural networks. If we had not tested the models on the test-B data, we would have thought that we have trained a strong classification model; yet, this was

Method	Lymp			Neut			Acc (%)
	P (%)	S (%)	F1 (%)	P (%)	S (%)	F1 (%)	
ResNet18 ⁴⁴	21	94	34	100	2	3	8
ResNet34 ⁴⁴	24	94	38	100	27	43	32
ResNet50 ⁴⁴	21	95	35	100	2	4	8
ResNext50 ⁴⁵	24	90	38	100	3	5	9
MnasNet1 ⁴⁸	20	88	33	100	0	0	6
MobileNet-V2 ⁴⁷	72	50	59	100	0	0	4
DenseNet121 ⁴⁶	43	64	51	100	8	15	12
ShuffleNet-V2 ⁴⁹	42	75	54	100	0	0	5
VGG16 ⁴³	96	89	93	100	65	79	66
Tavakoli et al. ⁵⁰	94	54	69	100	92	96	90

Table 9. The results of different pre-trained models as well as Tavakoli et al.⁵⁰ on the test-B dataset.

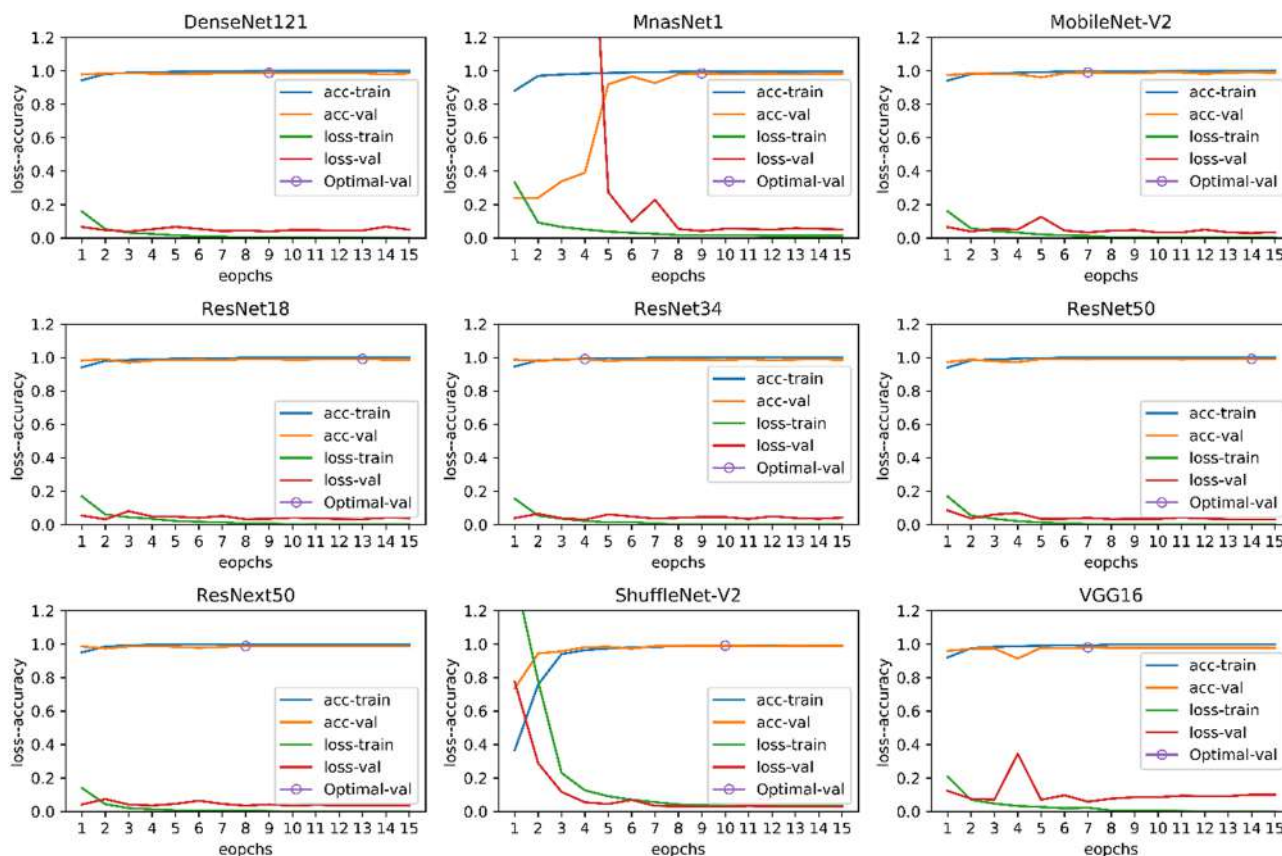


Figure 11. The plots of the accuracy and loss of training data and validation data related to nine pre-trained models.

not the case. In this experiment, we do not want to conclude that deep neural networks have less generalization power than feature-based methods. If we applied some appropriate pre-processing on the images before training or used some smarter image augmentation methods, the performance of deep neural networks would be better. In this experiment, you can easily understand the role of the dataset in the training of machine learning models.

All training processes were carried out using a single NVIDIA GeForce RTX 2080 Ti graphic card and were handled by Python 3.6.9 and Pytorch library version 1.5.1. We considered 15 epochs for the training process and the starting learning rate, and the batch size were 0.001 and 10, respectively. The learning rate was decayed by the ratio of 0.1 and step size 7. Stochastic gradient descent was utilized as the optimization method. We used the Torchvision library in order to load pre-trained networks on the ImageNet dataset⁵¹. The output size of the last linear layer was changed from 1000 to 5.

Conclusion

By evaluating the peripheral white blood cells, a wide range of benign diseases such as anemia and malignant ones such as leukemia can be detected. On the other hand, early detection of some of these abnormalities, such as acute lymphoid leukemia, despite its lethality, can help its treatment process. Therefore, it is important to adopt methods that can be effective in the early detection of different diseases. The role of machine learning methods in intelligent medical diagnostics is becoming more and more prominent these days. Indeed, deep neural networks are revolutionizing the medical diagnosis process and are considered as one of the state-of-the-arts.

Since deep neural networks usually have a huge number of training parameters, the overfitting problem is not highly unlikely. Therefore, the diversity of training data is necessary and cannot be ignored. In medical diagnostics, in particular, this diversity gets bolder, because the medical devices can be very diverse. For example, in the field of hematology, the type of microscope and camera is very influential. To this end, we collected a huge free available dataset of white blood cells from normal peripheral blood so as to relatively satisfy the mentioned diversity. This multipurpose dataset can serve as a reference dataset for the evaluation of different machine learning tasks such as classification, detection, segmentation, and localization.

Data availability

The Raabin-WBC dataset is publicly available through the following link: <https://www.raabindata.com/free-data/>.

Received: 29 May 2021; Accepted: 10 December 2021

Published online: 21 January 2022

References

1. Ten years in public health, 2007–2017: Report by Dr Margaret Chan, Director-General, *World Health Organization* (2017).
2. Serio, L. Importance of clinical lab testing highlighted during medical lab professionals week. <http://www.acla.com/importance-of-clinical-lab-testing-highlighted-during-medical-lab-professionals-week/> (2014).
3. Putzu, L., Caocci, G. & Di Ruberto, C. Leucocyte classification for leukaemia detection using image processing techniques. *Artif. Intell. Med.* **62**, 179–191 (2014).
4. McPherson, R. A. *Henry's Clinical Diagnosis and Management by Laboratory Methods: First South Asia Edition_e-Book* (Elsevier, 2017).
5. Hoffbrand, A. V. & Steensma, D. P. *Hoffbrand's Essential Haematology* (Wiley, 2019).
6. Wang, Q., Wang, J., Zhou, M., Li, Q. & Wang, Y. Spectral-spatial feature-based neural network method for acute lymphoblastic leukemia cell identification via microscopic hyperspectral imaging technology. *Biomed. Opt. Express* **8**, 3017–3028 (2017).
7. Lu, G. & Fei, B. Medical hyperspectral imaging: A review. *J. Biomed. Opt.* **19**, 1–24 (2014).
8. Mehdorn, M. *et al.* Hyperspectral imaging (HSI) in acute mesenteric ischemia to detect intestinal perfusion deficits. *J. Surg. Res.* **254**, 7–15 (2020).
9. Bain, B. J., Bates, I. & Laffan, M. A. *Dacie and Lewis Practical Haematology E-Book: Expert Consult: Online and Print* (Elsevier Health Sciences, 2016).
10. Avendi, M. R., Kheradvar, A. & Jafarkhani, H. A combined deep-learning and deformable-model approach to fully automatic segmentation of the left ventricle in cardiac MRI. *Med. Image Anal.* **30**, 108–119 (2016).
11. Shahin, A. I., Guo, Y., Amin, K. M. & Sharawi, A. A. White blood cells identification system based on convolutional deep neural learning networks. *Comput. Methods Programs Biomed.* **168**, 69–80 (2019).
12. Xu, M. *et al.* A deep convolutional neural network for classification of red blood cells in sickle cell anemia. *PLOS Comput. Biol.* **13**, 1–27 (2017).
13. Oktay, O. *et al.* Attention u-net: Learning where to look for the pancreas. arXiv 2018. arXiv Preprint [arXiv:1804.03999](https://arxiv.org/abs/1804.03999) (2018).
14. Roth, H. R. *et al.* DeepOrgan: Multi-level deep convolutional networks for automated pancreas segmentation. In *Medical Image Computing and Computer-Assisted Intervention* 556–564 (2015).
15. Roth, H. R. *et al.* Data from pancreas-CT. *Cancer Imaging Arch.* <https://doi.org/10.7937/K9/TCIA.2016.tNB1kqBU> (2016).
16. Litjens, G. *et al.* 1399 H&E-stained sentinel lymph node sections of breast cancer patients: The CAMELYON dataset. *Gigascience* <https://doi.org/10.1093/gigascience/giy065> (2018).
17. Labati, R. D., Piuri, V. & Scotti, F. All-IDB: The acute lymphoblastic leukemia image database for image processing. In *IEEE International Conference on Image Processing* 2045–2048 (2011).
18. Mohamed, M., Far, B. & Guaily, A. An efficient technique for white blood cells nuclei automatic segmentation. In *IEEE International Conference on Systems, Man, and Cybernetics* 220–225 (2012).
19. Sarrafzadeh, O., Rabbani, H., Talebi, A. & Banaem, H. U. Selection of the best features for leukocytes classification in blood smear microscopic images. In *Medical Imaging 2014: Digital Pathology*, Vol. 9041 159–166 (SPIE, 2014).
20. Rezatofghi, S. H. & Soltanian-Zadeh, H. Automatic recognition of five types of white blood cells in peripheral blood. *Comput. Med. Imaging Graph* **35**, 333–343 (2011).
21. Zheng, X., Wang, Y., Wang, G. & Liu, J. Fast and robust segmentation of white blood cell images by self-supervised learning. *Micron* **107**, 55–71 (2018).
22. Habibzadeh, M., Jannesari, M., Rezaei, Z., Baharvand, H. & Totonchi, M. Automatic white blood cell classification using pre-trained deep learning models: ResNet and inception. In *International Conference on Machine Vision*, Vol. 10696. 274–281 (SPIE, 2018).
23. Andrade, A. R. *et al.* Recent computational methods for white blood cell nuclei segmentation: A comparative study. *Comput. Methods Programs Biomed.* **173**, 1–14 (2019).
24. West, D. How mobile devices are transforming healthcare. *Issues Technol. Innov.* **18**, 1–11 (2012).
25. Ozcan, A. Mobile phones democratize and cultivate next-generation imaging, diagnostics and measurement tools. *Lab Chip* **14**, 3187–3194 (2014).
26. Kim, H. *et al.* Smartphone-based low light detection for bioluminescence application. *Sci. Rep.* **7**, 40203 (2017).
27. Majumder, S. & Deen, M. J. Smartphone sensors for health monitoring and diagnosis. *Sensors* **19**, 2164 (2019).
28. Mochida, K. *et al.* Computer vision-based phenotyping for improvement of plant productivity: A machine learning perspective. *Gigascience* <https://doi.org/10.1093/gigascience/giy153> (2018).
29. Wei, Q. *et al.* Detection and spatial mapping of mercury contamination in water samples using a smart-phone. *ACS Nano* **8**, 1121–1129 (2014).
30. Breslauer, D. N., Maamari, R. N., Switz, N. A., Lam, W. A. & Fletcher, D. A. Mobile phone based clinical microscopy for global health applications. *PLoS ONE* **4**, 1–7 (2009).
31. Martinez, A. W. *et al.* Simple telemedicine for developing regions: Camera phones and paper-based microfluidic devices for real-time, off-site diagnosis. *Anal. Chem.* **80**, 3699–3707 (2008).

32. Mudanyali, O. *et al.* Integrated rapid-diagnostic-test reader platform on a cellphone. *Lab Chip* **12**, 2678–2686 (2012).
33. Tseng, D. *et al.* Lensfree microscopy on a cellphone. *Lab Chip* **10**, 1787–1792 (2010).
34. Contreras-Naranjo, J. C., Wei, Q. & Ozcan, A. Mobile phone-based microscopy, sensing, and diagnostics. *IEEE J. Sel. Top. Quantum Electron.* **22**, 1–14 (2016).
35. Roy, S. *et al.* Smartphone adapters for digital photomicrography. *J. Pathol. Inform.* **5**, 24 (2014).
36. Ren, S., He, K., Girshick, R. & Sun, J. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**, 1137–1149 (2017).
37. Bradski, G. The opencv library. *Dr Dobbs J. Softw. Tools* **25**, 120–125 (2000).
38. Duda, R. O. & Hart, P. E. *Pattern Classification and Scene Analysis* Vol. 3 (Wiley, 1973).
39. Hegde, R. B., Prasad, K., Hebbar, H. & Singh, B. M. K. Development of a robust algorithm for detection of nuclei and classification of white blood cells in peripheral blood smear images. *J. Med. Syst.* **42**, 110 (2018).
40. Van Rossum, G. The Python Library Reference, release 3.8. 2. *Python Softw. Found.* **36** (2020).
41. Kouzehkanan, S.-Z. M., Tavakoli, I. & Alipanah, A. Easy-GT: Open-source software to facilitate making the ground truth for white blood cells nucleus. arXiv preprint [arXiv:2101.11654](https://arxiv.org/abs/2101.11654) (2021).
42. Otsu, N. A threshold selection method from gray level histograms. *IEEE Trans. Syst. Man. Cybern.* **9**, 62–66 (1979).
43. Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image recognition. arXiv Preprint. [arXiv:1409.1556](https://arxiv.org/abs/1409.1556) (2014).
44. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition* 770–778 (2016).
45. Xie, S., Girshick, R., Dollár, P., Tu, Z. & He, K. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition* 1492–1500 (2017).
46. Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K. Q. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition* 4700–4708 (2017).
47. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A. & Chen, L.-C. Mobilenetv2: Inverted residuals and linear bottlenecks. In *IEEE Conference on Computer Vision and Pattern Recognition* 4510–4520 (2018).
48. Tan, M. *et al.* Mnasnet: Platform-aware neural architecture search for mobile. In *IEEE Conference on Computer Vision and Pattern Recognition* 2820–2828 (2019).
49. Ma, N., Zhang, X., Zheng, H.-T. & Sun, J. ShuffleNet v2: Practical guidelines for efficient CNN architecture design. In *European Conference on Computer Vision* 116–131 (2018).
50. Tavakoli, S., Ghaffari, A., Mousavi Kouzehkanan, S.-Z. & Hosseini, R. New segmentation and feature extraction algorithm for the classification of white blood cells in peripheral smear images. *Sci. Rep.* **11**, 19428. <https://doi.org/10.1038/s41598-021-98599-0> (2021).
51. Deng, J. *et al.* ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition* 248–255 (2009).

Acknowledgements

The authors would like to thank the laboratory of Razi Hospital in Rasht, Gholhak laboratory, Shahr-e-Qods laboratory and, Takht Tavous laboratory in Tehran for their efforts in collecting samples of the dataset. Ms. Anahita Ahouraei and Mr. Arman Yekta are also to be appreciated for the efforts they put in editing the Raabin database contract as our attorneys at law. The authors would like to acknowledge the financial support of University of Tehran Science and Technology Park for this research under Grant Number 130067, as well.

Author contributions

Z.M.K. is the main author and the supervisor and executor of the artificial intelligence and image processing parts. S.S. is also the main author and the head of the medical and data collection department. S.T. and P.R. helped with manuscript revision and data collection. M.A., F.M., E.S.S., M.G., F.G., and S.M. have helped with the data collection and labeling process. R.H. is the advisor of the team and the corresponding author.

Competing interests

Raabin-WBC dataset has been compiled under the supervision of Nimaad Health Equipment Development Company. The authors and the company declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-04426-x>.

Correspondence and requests for materials should be addressed to R.H.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022