# Different Approaches for Automatic Nucleus Image Segmentation in Fluorescent in situ Hybridization (FISH) Analysis for HER2 Status Assesment

Denis Makhov, Andrey Samorodov
Bauman Moscow State Technical University
Moscow, Russian Federation
dennismak@yandex.ru, avs@bmstu.ru

Elena Slavnova
P. Hertsen Moscow Oncology Research Institute
Moscow, Russian Federation
mnioict@mail.ru

*Abstract*— according to American Cancer Society breast cancer is the most common cancer type in women. For most effective treatment choice and patients' state of health prediction it is necessary to make a differential diagnosis to determine breast cancer subtype. The tumor subtype is determined by immunohistochemical or immunocytochemical studies, which evaluate the expression levels of steroid hormone receptors, proliferative protein Ki-67, and oncoprotein CerbB-2 (HER2/neu). HER2-positive subtypes are most adverse (about 25-30% of all cases). In case of indefinite CerbB-2 expression fluorescence in situ hybridization (FISH) investigation is utilized. In most cases, this study is held by visual estimation of fluorescent image parameters by pathologist and thus is subjective. We need to employ automatization techniques to decrease human factor impact and increase reproducibility of the analysis result. FISH analysis automatization for HER2 amplification can be divided into three tasks: nucleus segmentation, signal detection and presentation of the results according to ASCO/CAP recommendations. In this article results for nucleus segmentation task using different machine learning algorithms are presented. The image database for investigations consisted of RGB fluorescent images, as well as gray scale images for each individual fluorophore. The best result was achieved using the random forest algorithm on gray-scale images of individual fluorophores.

## I. INTRODUCTION

Oncological diseases take second place in the list of leading mortality causes according to WHO [1]. Latest open cancer disease statistics states more than 18 million new cases of cancer in 2018 and according to the forecast for 2040 this number will increase to about 40 million [2]. Most common cancer type in women is breast cancer (BC) [3]. BC mortality is about 30% from the number of new cases [3].

Decreasing of BC mortality requires the development of early disease detection, differential diagnosis methods and effective treatments. The choice of treatment tactics depends on the breast cancer subtype which is assigned via differential diagnosis. Breast cancer subtype is determined by the expressions of steroid hormone receptors (estrogen and progesterone receptors, ER and PR, respectively), proliferative protein Ki-67, and oncoprotein CerbB-2 (HER2/neu). Reaction for ER, PR and Ki-67 expression evaluation gives nuclear staining and reaction with antibodies to CerbB-2 gives membrane staining. Automatization problem for nuclear staining was considered in detail by Dobrolyubova Daria [4].

HER2/neu positive (25–30% of all cases) and basal cell tumors (8–20% of all cases) are the most adverse subtypes of BC. The ASCO/CAP 2018 recommendations for HER2-status evaluation in IHC study are listed in table I [5].

TABLE I. HER2-SATUS ACCORDING TO ASCO/CAP RECOMMENDATIONS

| Score | Description |
|---|---|
| 0 | No staining is observed or membrane staining that is incomplete and is faint/barely perceptible and in ≤10% of tumor cells |
| 1+ | Incomplete membrane staining that is faint/barely perceptible and in >10% of tumor cells |
| 2+ | Weak to moderate complete membrane staining observed in >10% of tumor cells |
| 3+ | Circumferential membrane staining that is complete, intense and in >10% of tumor cells |

In the case of an indefinite reaction (2+), either a repeated study or determination of the amplification of the HER2 gene using the fluorescence in situ hybridization method (FISH) is required to clarify the HER2 status.

Special fluorescent probes that are complementary to specific DNA sites (for example, DAKO HER2FISH) are used for FISH studies. Amplification of the HER2 gene is assessed by counting the number of signals referred to HER2 gene probe (red signal) and centromeric region of the 17th chromosome probe (green signal) (Fig. 1). An average copy number of the HER2 gene per core and the ratio of red and green labels is determined in at least 20 cores and HER2 amplification is assessed (table II) [5].

TABLE II. HER2 GENE AMPLIFICATION IN FISH STUDY ACCORDING TO ASCO/CAP RECOMMENDATIONS

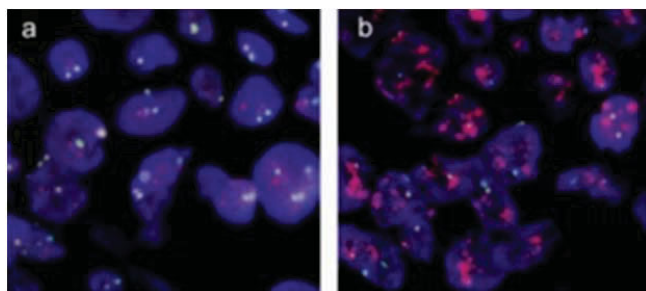| Amplification | Description |
|---|---|
| HER2 gene amplification | Average HER2 copy number ≥ 6.0 signals/cell |
| Undetermined | Average HER2 copy number ≥ 4.0 and < 6.0 signals/cell |
| No HER2 gene amplification | Average HER2 copy number < 4.0 signals/cell |

Fig. 1. Examples of fluorescent images with (right) and without (left) HER2 gene amplification

Image registration during a FISH study is carried out using a fluorescence microscope and a highly sensitive digital camera. In most laboratories, an oncologist implements signal counting visually. To increase reproducibility of the analysis, it is necessary to apply automated image analysis methods.

Automatization for HER2 amplification assessment in FISH analysis can be divide into following steps:

- nucleus segmentation and separation individual nuclei in agglomerations;
- signal detection inside individual nuclei;
- medical decision support by representing calculated parameters.

To assess the state of FISH analysis automation for determining the amplification of the HER2 gene, we considered the work of the world's leading teams in this field [6–10].

The automation problem can be solved using, for example, simple threshold segmentation method proposed by Xingwei Wang and colleagues in [6]. The authors of the article proposed a method for automating FISH cervical cancer research. However, this approach cannot be used for FISH images of breast biopsy specimens, since in some cases autofluorescence of pepsin underprocessed cytoplasm is observed (Fig. 2). It is related to the fact that pepsin treatment is a standardized procedure that does not take into account the characteristics difference in biological tissue.
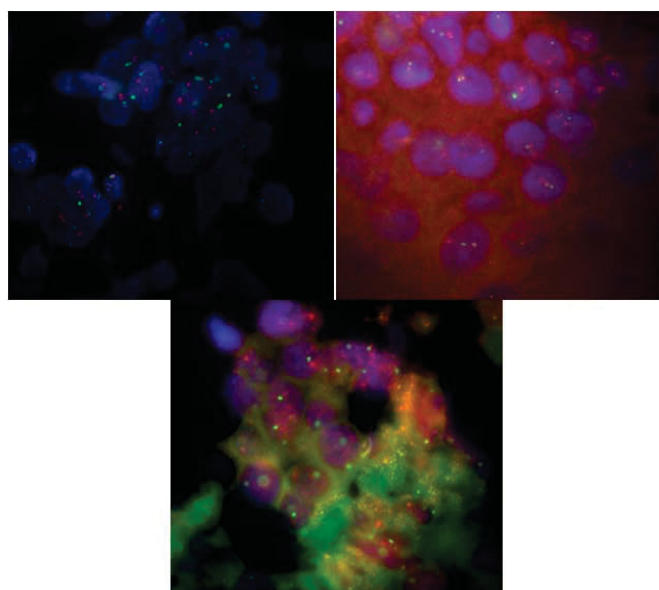


Fig. 2. Different quality of real-life fluorescent images

Tomasz Les and colleagues [7] used a watershed method to segment nuclei. Disadvantage of this method is big number of false boundaries, and skipping some nuclei with low contrast relative to the background (Fig. 3).
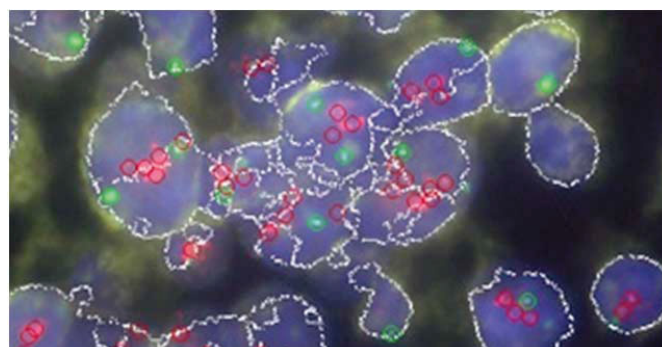


Fig. 3. Nuclei segmentation using watershed algorithm

The proposed approach can give an inaccurate estimation of signals distribution statistics in nuclei, since some nuclei are falsely divided into two or more parts and some are skipped. The article indicates the need for an expert to select informative nuclei after automatic image labeling. In this case a complete agreement between the results of the algorithm and the expert comment on a sample of 10 patients was obtained.

Falk Zakrzewski and colleagues used a large selection of images of high-quality slides, consisting of 299 images annotated with bounding box and class for each nucleus (5 classes in total: low, normal and high grade, uncertain, and artifact), and 301 images annotated with nuclei and signals bounding boxes [8]. In the study two convolutional neural networks (CNN) with the RetinaNet architecture implemented using the Keras library were used. As a result, on test database consisted of 57 images the accuracy of the nucleus detection algorithm was 48–97% and the accuracy of the signal detection algorithm was 44–97%, the accuracy of determining the amplification of the HER2 gene was 96%. One of the key advantages according to the authors is a two-stage nucleus classification, which leads to an increased algorithm robustness.

Henning Höfener and colleagues created an automated density-based algorithm for counting FISH amplification signals for HER2 status assessment [9]. Signals on FISH images regularly can't be easily detected as they make up clusters and thus the quality of HER2 status assessment depends on observer's experience. In the article authors proposed density-based CNN for counting signals that showed the best results among difference of Gaussian, CNN-detect and CNN-accumulate methods. This approach showed more robust on signal clusters and consequently mean normalized absolute errors for CEN17 signals, ERBB2 signals, and all signals combined were the least.

Gedmante Radziuviene with colleagues made a comparison study between automated and manual CEP17 and HER2 counting [10]. Automated signal counting was made after manual nucleus segmentation and StrataQuest v.205 software application. The results showed that automated methods are still insufficient for clinical use as they underestimate both CEP17 and HER2 signals.

An analysis of these studies allows us to formulate the following statements:

- information on the results of listed algorithms is not enough to select the most accurate solution to the problem of nucleus segmentation and separation of agglomerates;
- major efforts were applied to deploy robust algorithm for automated signal counting on FISH images;
- due to the peculiarities of fluorescent images, there is a trade-off between the quality of the images used for automation and the accuracy of the segmentation and classification algorithms.

In this article various algorithms for solving the nucleus segmentation problem in FISH images of various quality (Fig. 2) will be discussed.

## II.    MATERIALS AND METHODS

### A. Equipment and database

Image database was collected in Hertsen Moscow Oncology Research Institute. Slides for FISH were prepared using DAKO fluorescent probes, their spectral characteristics are shown on the figures 4, 6, 8 for DAPI, FITC and TexasRed, respectively [11]. Images were captured by Zeiss Axio Imager A1 fluorescence microscope with a MetaSystems CoolCube1 digital camera (resolution 1360 x 1024 pixels). There are Chroma SP100V2 (DAPI), Chroma MF101 (FITC) and Chroma SP103V1 (Texas Red) filters installed in the microscope, their transmission characteristics are presented on Fig. 5, 7, 9, respectively [12]. MetaSystems Isis 5.0 software was used for image registration. Due to the software features, the exported images had resolution of 990 x 878 pixels.
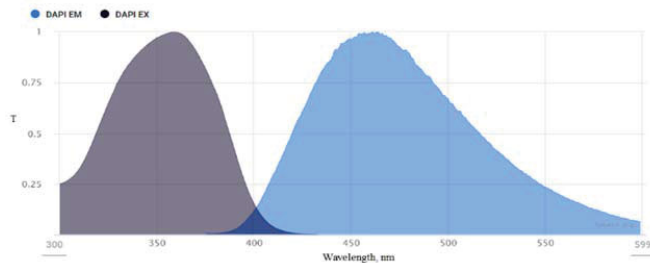


Fig. 4. Excitation (left) and emission (right) spectral characteristics for DAPI
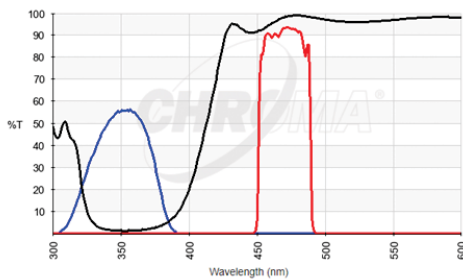


Fig. 5. Transmission characteristics of Chroma SP100V2 (DAPI), blue – excitation, red – emission, black – dichroic mirror
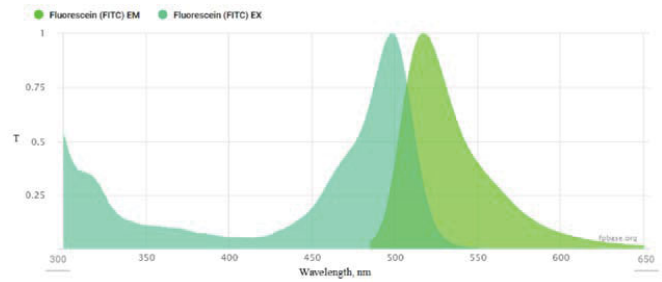


Fig. 6. Excitation (left) and emission (right) spectral characteristics for FITC
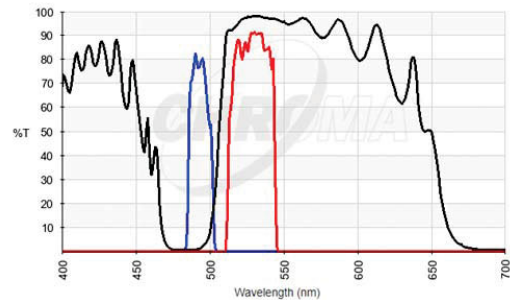


Fig. 7. Transmission characteristics of Chroma MF101 (FITC), blue – excitation, red – emission, black – dichroic mirror
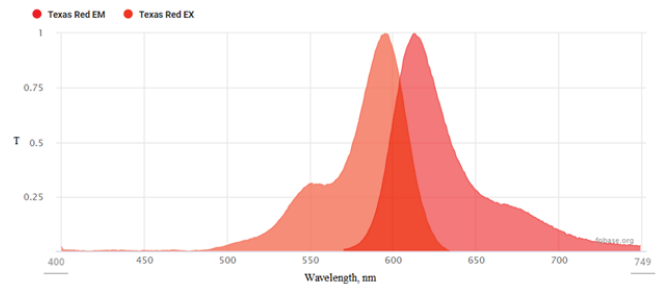


Fig. 8. Excitation (left) and emission (right) spectral characteristics for TexasRed
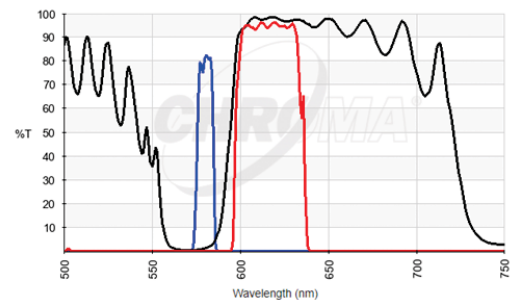


Fig. 9. Transmission characteristics of Chroma SP103V1 (Texas Red), blue – excitation, red – emission, black – dichroic mirror

Image database consists of RGB images and gray-scale images for each of the DAPI, FITC, and TexasRed (TR) fluorophore channel. All images were saved without any color-correction operations. Median (size 3x3) and Gauss filters (size 3x3, sigma 0.5) were applied. Further, these images were manually annotated with edges of nuclei by a pathologist. This annotation was carried out on the DAPI channel images, since nuclei contours were most contrast and sharp. The resulting database contained 30 series of images: RGB, DAPI, FITC, TR.

## B. RGB image nucleus segmentation

Due to the fact that exporting one RGB image instead of 3 images for different channels is less laborious, we first consider what results can be obtained using RGB images.

All images were collected in a single array, in which the columns represent the color coordinates (features) and the class label – target (1 – for nuclei and 0 - background). We suggest that visualization could be helpful to understand types of distribution in target classes and assess the degree of segmentation task complexity (Fig. 10).
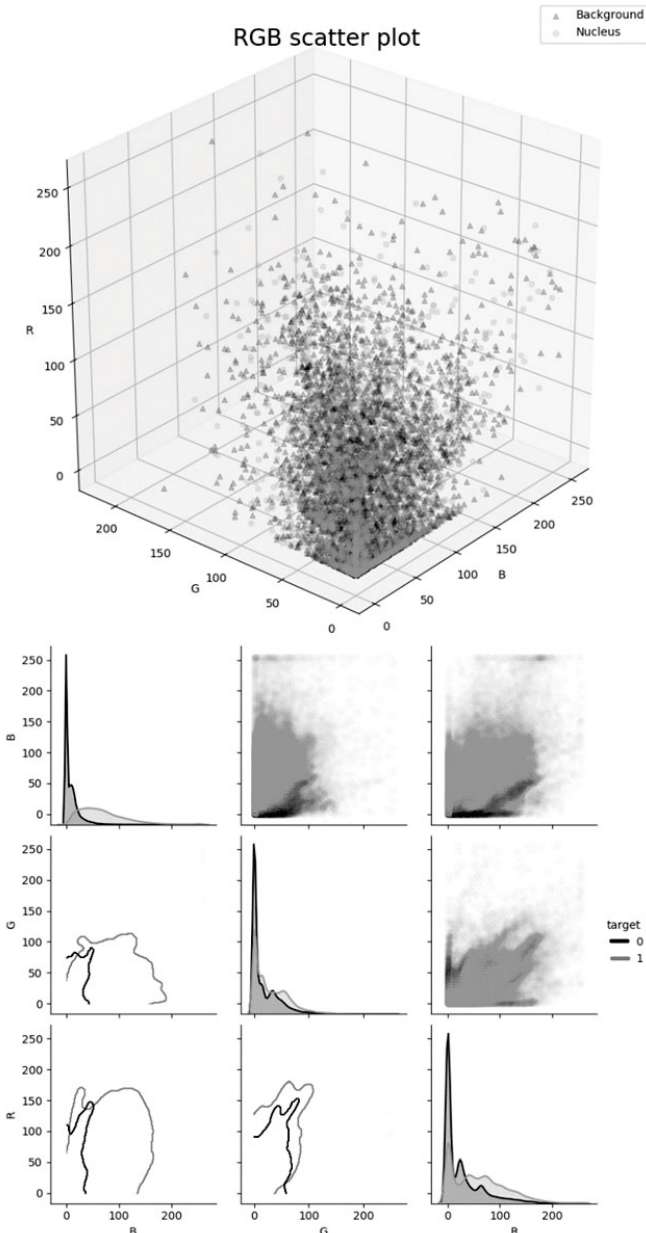


Fig. 10. Scatter plot and pair plot for RGB color space

The scatter plot of RGB components shows that the clouds of pixels belonging to different classes strongly overlap. From the pair plots we can conclude that separation of the distributions in the channels R and G is difficult. Hereinafter plots in diagonal illustrates histograms of each color component (feature)

according to its class label (target). Plots under diagonal illustrates edges for 95% of histogram weight for concrete component pair and target (1 – nuclei, 0 – background). Plots above diagonal illustrates scatter plot projection on 2D feature axis.

## C. Nucleus segmentation using additional color spaces

The original RGB color space was complemented by the HSV, Lab, LUV, XYZ, YCrCB, YUV spaces. For machine learning algorithms, it is preferable to have independent features.

The correlation coefficients for some pairs of color components modulo was close to 1, since the transformations of these color components pairs are linearly dependent. Color coordinates with an absolute correlation coefficient more than 0.9 were excluded from further analysis. Thus, for further research, we used only the R, G, B channels of the RGB space, the H, S, V channels of the HSV space, the a, b channels of the Lab space, and the Cr channel of the YCrCb space (Fig. 11).
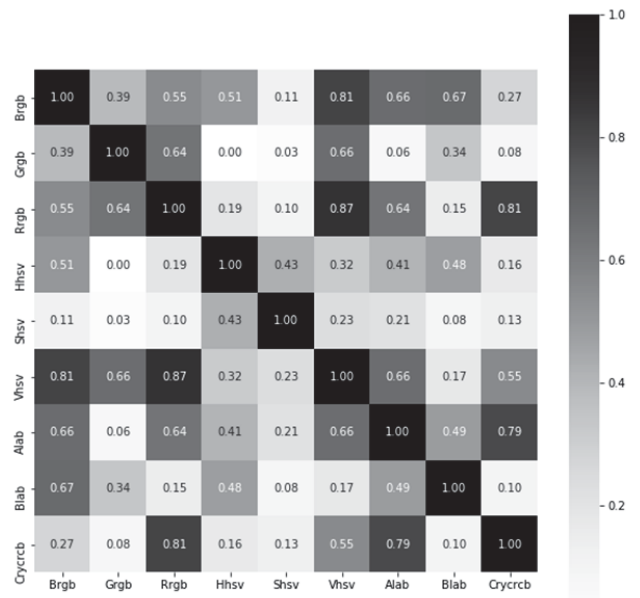


Fig. 11. Correlation coefficient heatmap for color components with absolute value of correlation coefficient <0.9

To increase the learning speed of the algorithms, dimensionality reduction using the principal component analysis (PCA) was applied. Explained variance ratio without removing components is presented in Table III.

TABLE III. EXPLAINED VARIANCE RATIO FOR ALL PCA COMPONENTS

| Component № | Explained variance ratio |
|---|---|
| 1 | 0.4728 |
| 2 | 0.2351 |
| 3 | 0.1357 |
| 4 | 0.0992 |
| 5 | 0.0498 |
| 6 | 0.0052 |
| 7 | 0.0018 |
| 8 | 0.0003 |
| 9 | 0.0000 (8e-6) |

It was decided to leave only 5 components, since the number of features was halved, and the amount of information was reduced only by 0.7%. The figure below (Fig. 12) shows the scatter plot of the background (triangles) and the nuclei (circles) pixels for the first three components of the PCA.
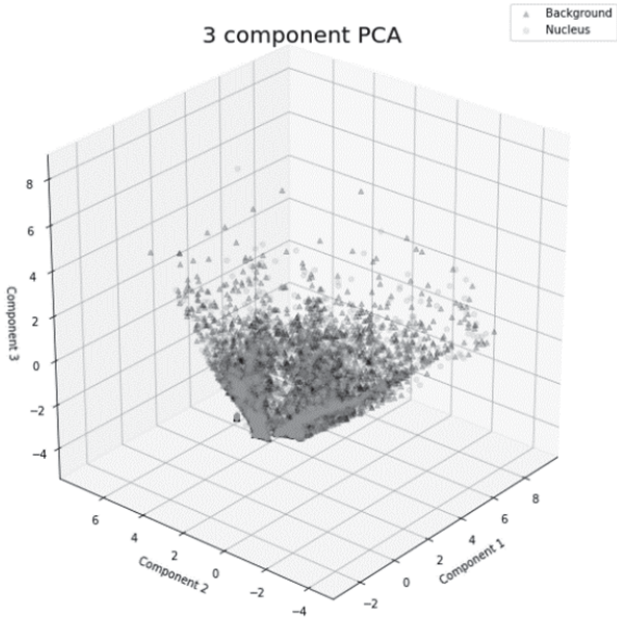


Fig. 12. Scatter plot for 3 components of PCA

Visualizing only three components represents a strong overlap of classes.

A visualization of the distribution of samples across five components is shown in Fig. 13.
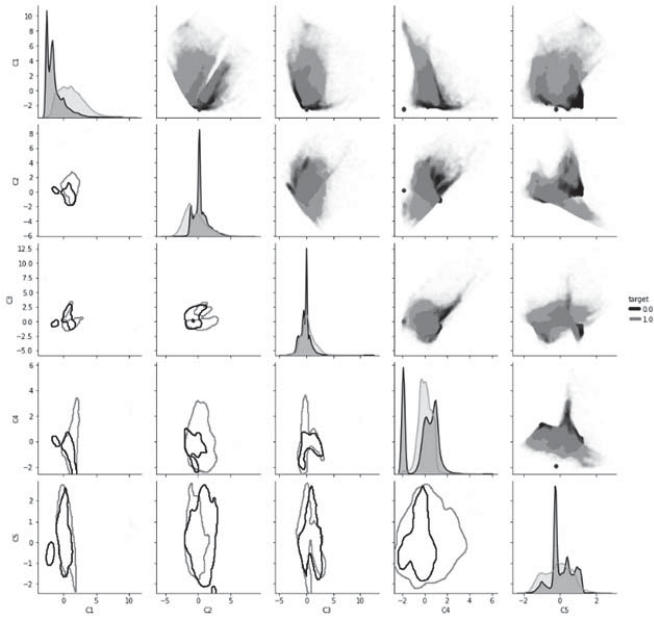


Fig. 13. Pair plot for 5 components of PCA

Fig. 13 shows that visual separation of classes after artificially expanding the number of features and then reducing the dimension using PCA is not improved compared to a model using only RGB coordinates.

*D. Nucleus segmentation using phluorophores channels*

We did primary visual analysis with scatter and pair plot just like we did it for RGB color space (Fig. 14).
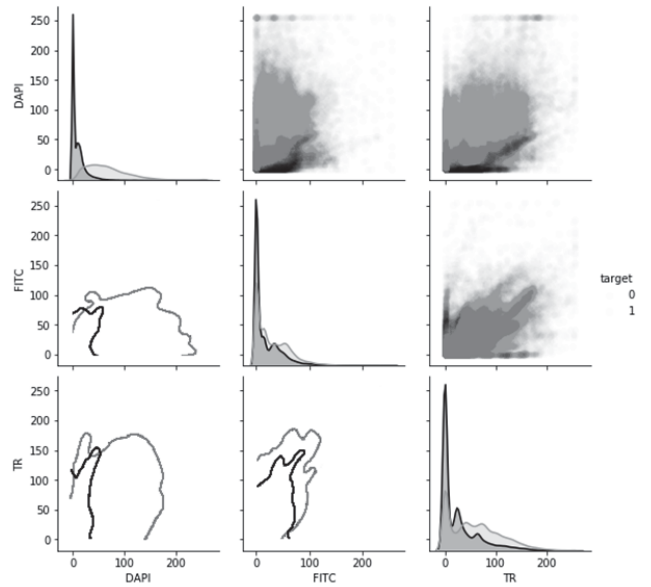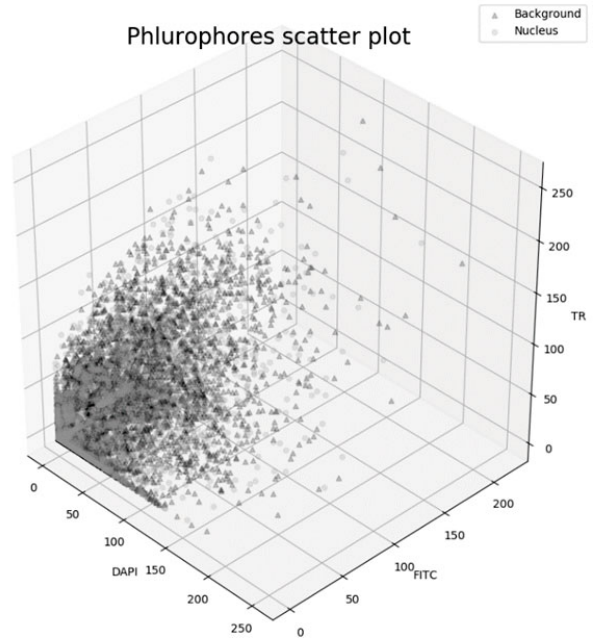


Fig. 14. Scatter plot and pair plot for phluorophores channels

Plots are similar to those for RGB color space: classes are still hardly separable (especially in FITC and TR channel).

*E. Proposed algorithms*

To find the best algorithm for solving the segmentation problem, linear classification algorithms (Linear, Ridge, Lasso, Logistic Regression), nonlinear algorithms (Linear, Ridge,

Lasso, Logistic Regression with polynomial feature transformation), ensemble algorithms (Balanced random forest, RUSBoost) and a convolutional neural network (CNN) based on the Unet architecture were used.

1) Linear regression classifier

One of the most simple and fastest classification algorithms [13]. The main idea is to solve regression task to minimize the following functional:

$$\min_{w} \| Xw - y \|_2^2$$

where X - feature matrix, w – weight vector, y – target vector.

Finally, the classification task is to determine which side of the inflating hyperplane lies concrete sample with its concrete features.

This approach in some cases can lead to overfitting so we need to enter some penalties, i.e. L1 and L2.

2) Lasso regression classifier

This classifier is linear regression classifier improved with L1 penalty. The lasso coefficients minimize a penalized residual sum of squares [13]:

$$\min_{w} \frac{1}{2n_{samples}} \| Xw - y \|_2^2 + a \| w \|_1$$

where a is regularization parameter.

3) Ridge regression classifier

This classifier is built on ridge regression. The ridge coefficients minimize a penalized residual sum of squares [13]:

$$\min_{w} \| Xw - y \|_2^2 + a \| w \|_2^2$$

where X - feature matrix, w – weight vector, y – target vector.

The classifier first converts binary targets to {-1, 1} and then treats the problem as a regression task, optimizing the same objective as above. The predicted class corresponds to the sign of the regressor's prediction.

In practice this method is very fast and its results are similar to more complex logistic regression.

4) Logistic regression classifier

Logistic regression classifier is a linear model, that minimizes next functional in case of L1 regularization [13]:

$$\min_{w,c} \| w \|_1 + C \sum_{i=1}^{n} \log(\exp(-y_i(X_i^T w + c)) + 1)$$

where C – inverse of regularization strength, c – bias.

And in case of L2 regularization [13]:

$$\min_{w,c} \frac{1}{2} w^T w + C \sum_{i=1}^{n} \log(\exp(-y_i(X_i^T w + c)) + 1)$$

5) Ensemble algorithms

Combining a large number of weakly correlated classifiers into a single classifier allows to avoid overfitting and get an unbiased error on the test data.

The two most popular ensemble algorithms are random forest and ADABoost [14].

In this study, we use the imbalanced learn library with its balanced random forest (BRF) implementation (the number of estimators is equal to 100, unlimited maximum depth), and instead of ADABoost, the RUSBoost algorithm is used, which is optimized for working with unbalanced samples (the number of estimators is equal to 50 and the maximum estimator depth equals to 1).

6) Convolutional neural network

In this research we utilize CNN with Unet architecture proposed in [15] with adding batch normalization layers before each activation layer.

For deep learning the image database was divided into training, validation and testing parts in the ratio of 70:20:10, respectively. The training process was carried out with mini-batches of 3 images. Data augmentation was not used, since the overfitting was insignificant in this sampling, and any data augmentation led to a validation and test accuracy decrease (Fig. 15).
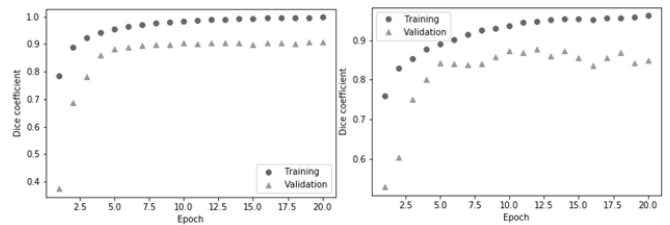


Fig. 15. CNN learning process without (left) and with (right) augmentation

*E. Metric and training process*

Sorensen-Dice coefficient as training and validation metric was used, it is calculated by the following formula [16]:

$$Dice = 2 * |True \cap Pred| / True \cup Pred$$

where Dice – Dice coefficient, True – true segmentation mask and Pred – predicted segmentation mask.

The training process of linear, nonlinear and ensemble algorithms was carried out for a stratified 5-fold split with balancing class weights. Previously, data was scaled by subtracting the mean and dividing by the variance.

For the comparison we also deployed the transition to polynomial features to improve the quality of linear algorithms. Due to the ascending computational costs with increasing degree of polynomial transition, the 3-d degree was chosen.

## III. RESULTS

### A. RGB color space

The results obtained using the listed algorithms are presented in table IV.

TABLE IV. AVERAGE DICE COEFFICIENT FOR DIFFERENT ALGORITHMS IN RGB COLOR SPACE

| Algorithm | Average Dice coefficient |
|---|---|
| Linear regression classifier | 0.6664 |
| Linear regression classifier (polynomial features, n=3) | 0.6664 |
| Lasso regression classifier | 0.6664 |
| Lasso regression classifier (polynomial features, n=3) | 0.6664 |
| Ridge regression classifier | 0.8116 |
| Ridge regression classifier (polynomial features, n=3) | 0.8896 |
| Logistic regression classifier | 0.8808 |
| Logistic regression classifier (polynomial features, n=3) | 0.8965 |
| BRF | 0.9270 |
| RUSBoost | 0.8903 |
| Unet | 0.9089 |

The table shows that in some cases polynomial features did not increase the Dice coefficient. For linear regression classifier and Lasso regression classifier, this can be explained with the fact that weights of additional polynomial features in the trained model are orders of magnitude less than initial features.

### B. RGB and additional color spaces

Average Dice coefficients for the same set of algorithms for sampling obtained using the sequential expansion of the number of features with components from other color spaces and reduction of dimension by the PCA are presented in table V.

TABLE V. AVERAGE DICE COEFFICIENT FOR DIFFERENT ALGORITHMS USING COMPONENTS FROM DIFFERENT COLOR SPACES

| Algorithm | Average Dice coefficient |
|---|---|
| Linear regression classifier | 0.6664 |
| Linear regression classifier (polynomial features, n=3) | 0.6664 |
| Lasso regression classifier | 0.6664 |
| Lasso regression classifier (polynomial features, n=3) | 0.6664 |
| Ridge regression classifier | 0.8385 |
| Ridge regression classifier (polynomial features, n=3) | 0.9061 |
| Logistic regression classifier | 0.8822 |
| Logistic regression classifier (polynomial features, n=3) | 0.9081 |
| BRF | 0.9265 |
| RUSBoost | 0.8923 |
| Unet | 0.9110 |

Lack of change in Dice coefficient for linear and Lasso regression classifiers in this case is explained analogously to the results interpretation for RGB color space.

### C. Phluorophores channels

It was determined that the RGB images' channels exported by Isis software do not coincide with the gray-scale images of individual fluorophore channels. Therefore, the training algorithms were applied to an array consisting of gray-scale images for three fluorophores. The results are shown in table VI.

In this case, adding polynomial features led to the Dice coefficient increase for linear and Lasso regression classifiers,

since some of the new generated features had similar weights in comparison with the initial features. This means that part of the polynomial features made it possible to obtain some new information about the data. Perhaps an increase in the number of polynomial features can lead to a further Dice coefficient increase for these two algorithms. A Dice coefficient decrease coefficient for logistic regression may be explained with overfitting on the training data, which led to a decrease in the Dice coefficient on the test sampling in cross-validation.

TABLE VI. AVERAGE DICE COEFFICIENT FOR DIFFERENT ALGORITHMS USING PHLUOROPHORES CHANNELS

| Algorithm | Average Dice coefficient |
|---|---|
| Linear regression classifier | 0.6664 |
| Linear regression classifier (polynomial features, n=3) | 0.6743 |
| Lasso regression classifier | 0.6664 |
| Lasso regression classifier (polynomial features, n=3) | 0.7428 |
| Ridge regression classifier | 0.8119 |
| Ridge regression classifier (polynomial features, n=3) | 0.8555 |
| Logistic regression classifier | 0.8819 |
| Logistic regression classifier (polynomial features, n=3) | 0.8543 |
| BRF | 0.9300 |
| RUSBoost | 0.8934 |
| Unet | 0.9113 |

## IV. CONCLUSIONS

Despite the presence of a sufficiently large number of works on the automation of FISH studies to determine the amplification of the HER2 gene, it can be noted that currently there is not enough information for choosing optimal solution to the problems of nucleus segmentation, agglomerate separation and signal detection.

In this paper various classes of algorithms and different approaches in image preprocessing for solving the nucleus segmentation problem were compared. As a result of the research on the available annotated image database, the following conclusions can be drawn:

- using separate grayscale phluorophores images leads to the segmentation quality improvement in comparison with the results obtained in RGB and extended RGB feature spaces. The Dice coefficient in this case is the maximum for almost all applied algorithms;
- artificially expanding the color channels of RGB images with subsequent compression by the PCA gives a slight segmentation quality improvement in most cases;
- Balanced random forest is the algorithm with the highest Dice coefficient. The Dice coefficient for this algorithm was the largest for all three data compilation approaches; the maximum value was 0.93.

## REFERENCES

[1] World Health Organization official website, Cancer, Web: https://www.who.int/news-room/fact-sheets/detail/cancer.

[2] International Agency for Research on Cancer, Cancer Tomorrow, Web: https://gco.iarc.fr/tomorrow/home.

[3] F. Bray, J. Ferlay, I. Soerjomataram, R.L. Siegel, L.A. Torre; A. Jemal, "Global Cancer Statistics 2018: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries", CA *CANCER J CLIN*, vol. 68(6), 2018, pp. 394–424.

[4] D.A. Dobrolyubova, et al., "Automatic image analysis algorithm for quantitative assessment of breast cancer estrogen receptor status in immunocytochemistry.", *Pattern Recognition and Image Analysis* Vol. 26.3, 2016, pp. 552-557.

[5] College of American Pathologists, HER2 Testing in Breast Cancer - 2018 Focused Update, Web: https://www.cap.org/protocols-and-guidelines/cap-guidelines/current-cap-guidelines/recommendations-for-human-epidermal-growth-factor-2-testing-in-breast-cancer.

[6] X. Wang et. al. "Automated analysis of fluorescent in situ hybridization (FISH) labeled genetic biomarkers in assisting cervical cancer diagnosis", *Technology in cancer research & treatment*. Vol. 9(3), 2010, pp. 231-242.

[7] T. Les, T. Markiewicz, S. Osowski, M. Cichowicz, W. Kozlowski, "Automatic evaluation system of FISH images in breast cancer", *International Conference on Image and Signal Processing, Springer,* 2014. pp. 332-339.

[8] F. Zakrzewski, et al., "Automated detection of the HER2 gene amplification status in Fluorescence in situ hybridization images for the diagnostics of cancer tissues", *Scientific reports*, vol. 9.1, 2019, pp. 1-12.

[9] H. Höfener, et al., "Automated density-based counting of FISH amplification signals for HER2 status assessment.", *Computer methods and programs in biomedicine*, Vol. 173, 2019, pp. 77-85.

[10] G. Radziuviene, et al., "Automated image analysis of HER2 fluorescence in situ hybridization to refine definitions of genetic heterogeneity in breast cancer tissue.", *BioMed research international 2017*, 2017.

[11] FPbase Fluorescence Spectra Viewer Web: https://www.fpbase.org/spectra/.

[12] Chroma Technology | Optical Filters, Custom & OEM Filter Design, Web: https://www.chroma.com/.

[13] Scickit-learn 0.22.1 documentation, Linear models, Web: https://scikit-learn.org/stable/modules/linear_model.html.

[14] Imbalanced-learn 0.5.0 documentation, Ensemble of samplers, Web: https://imbalanced-learn.readthedocs.io/en/stable/ensemble.html.

[15] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," *Medical Image Computing and Computer-Assisted Intervention, Lecture Notes in Computer Science*, vol. 9351, November 2015, pp. 234–241.

[16] D. Parpulov, et al., "Convolutional neural network application for cells segmentation in immunocytochemical study," *Proceedings of the 2018 Ural Symposium on Biomedical Engineering, Radioelectronics and Information Technology (USBEREIT 2018)*, 2018, pp.87-90.