

# Comparison of Different Convolutional Neural Network Architectures for Satellite Image Segmentation

Vladimir Khryashchev, Leonid Ivanovsky, Vladimir Pavlov  
 P.G. Demidov Yaroslavl State University  
 Yaroslavl, Russian Federation  
 v.khryashchev@uniyar.ac.ru, leon.ivanovsky@yahoo.com,  
 vladimir@lpavlov.com

Anna Ostrovskaya  
 People's Friendship University of Russia  
 (RUDN University)  
 Moscow, Russian Federation  
 ostrovskaya-aa@rudn.ru

Anton Rubtsov  
 Russian Space Systems,  
 Moscow, Russian Federation  
 rubtsov493@gmail.com

**Abstracts**—Convolutional neural networks for detection ge-objects on the satellite images from DSTL, Landsat-8 and PlanetScope databases were analyzed. Three modification of convolutional neural network architecture for implementing the recognition algorithm was used. Images obtained from the Landsat-8 and PlanetScope satellites are used for estimation of automatic object detection quality. To analyze the accuracy of the object detection algorithm, the selected regions were compared with the areas by previously marked by experts. An important result of the study was the improvement of the detector for the class “Forest”. Segmentation of satellite images has found application at urban planning, forest management, climate modelling, etc.

## I. INTRODUCTION

The progress in the development of high-performance computers with graphics processing units (GPU) allowed researchers to work with convolutional neural networks (CNNs) that have millions of parameters. In solving modern problems of computer vision, modern approaches based on CNN exceed traditional methods and work of image analysis experts in different cases. CNNs demonstrated their advantage in tasks of image classification, object detection and scene recognition. Currently deep learning methods are used for solving almost all problems of computer vision [1]. Image segmentation is one of these tasks.

The problem of satellite images segmentation is challenging. In machine learning applications, aerial image interpretation is usually formulated as a pixel labeling task. Nowadays the object detection for aerial high-resolution photos is in the focus of research community. Meanwhile, the most approaches to solve this problem is the use of a CNNs. The features in these networks are formed automatically in the process of training.

Nowadays, per-pixel satellite image segmentation requires the use of deep learning algorithms. The use of such methods instead of traditional approaches is non-trivial for some reasons. Unique methods are needed to solve the problem of

the spatial extent of the detected objects and the invariance to the rotation or the scale of the images [2]. These algorithms must adjust for:

- Taking into account the small spatial extent of objects. Detection of small objects in large images is one of the main problems in the satellite images analysis. Unlike the large objects captured in the ImageNet database, objects on satellite images are often very small, but they are densely grouped. The reason is that the resolution of satellite image is determined based on the distance to the ground. It determines what is captured on one pixel of the image. The size of captured area on these images usually ranges from 30 cm<sup>2</sup> to 16 m<sup>2</sup>. This means that the size of an object such as a car will only take 15 pixels.
- Being invariant to rotation. Objects on satellite images have different orientation. For example, ships can be rotated to any angle, while trees in the forest are located vertically.
- Having sufficient amount of training images. For most available datasets, such as LandSat [3] and Inria [4], there is a shortage of annotated images. However, some efforts to create a large number of learning samples, such as SpaceNet [5], can solve this problem.
- Being able to work with high-resolution images. Satellite images, which are used by algorithms of machine learning, are high-resolution. For example, some images from the DigitalGlobe satellite cover more than 64 km<sup>2</sup>, which includes more than 250 million pixels.

Image segmentation is the separation of image into significant areas, which can be considered as a task of per-pixel classification. The simplest (and slow) approach to solve this problem is manual segmentation of images. Nevertheless, it is a laborious and long process, which usually leads to make mistakes. Currently, the great interest of researchers in the field of machine learning is associated with the development of

automatic image segmentation systems. This type of segmentation allows to process images immediately after the receiving. Such automatic systems must provide the necessary accuracy to be useful in practical applications.

This article consists of seventh parts. The second part is devoted to the related works concerning the problem of satellite images segmentation. The third part is devoted to available databases of satellite images. The fourth section describes developed CNN architectures for image segmentation. Also in this part there were described tools for building classifiers, as well as the peculiarities of training process. The fifth and sixth sections show the results of numerical experiments for developed models on different databases. Then, in conclusion there could be found summarizing and suggestions about the possible improvement of proposed classifiers.

## II. RELATED WORKS

In recent years, various methods for creating CNN have been proposed, which can produce segmentation for the entire input image. One of the most successful algorithms is based on fully convolutional networks (FCNs). The basic idea of this approach is to use CNN to extract the necessary feature values, while replacing the fully connected layer with a convolution layer with the output in the form of feature maps instead of classification results [6]. This method allows you to train CNN for the segmentation of images of different sizes.

Following this way, the authors of [7] present the architecture of CNN named Mask R-CNN, using pre-trained weights of COCO dataset. This algorithm is composed by two networks: a Region Proposal Network (RPN) and a FCN. The first model takes the whole input image and output the transformed image with bounding box proposals of detected objects. The second model uses the information from previous network and performs the segmentation for transformed image. This nonsimple method allowed to detect buildings on satellite photos of Inria Aerial Images Dataset [4] exactly. The best performance was 92.49%.

The method of using FCN was supplemented and now it is known as U-Net. In paper [8] there is presented U-Net architecture – a specific type of FCN, which had received a lot of interest for segmentation of biomedical 2D and 3D images [8, 9]. Later, this model has proven to be very efficient for the pixel-wise classification of satellite images [10]. The U-Net architecture uses skip-connections to combine low-level and higher-level feature maps, which provides accurate localization of objects. Using U-Net architecture, the authors of paper [10] get the value of Sorensen-Dice coefficient is equal to 0.75.

The authors of [11] hold the similar method to solve the problem of satellite images segmentation. They developed the U-Net like architecture, which is using ResNet-34 weights in the encoder. This algorithm shows excellent results of detecting roads on satellite images of DeepGlobe database [12]. The best public score is 0.64.

Classical neural network architectures, such as DeconvNet, which contains only coders and decoders without merging layers, are also very popular for solving similar problems [13].

Such neural networks are usually initialized by network weights preliminarily trained on a large dataset (for example, ImageNet or ResNet). This approach allows to significantly increase the efficiency of the learning process. In some practical applications, such as cancer detection or traffic safety, the accuracy of models plays the crucial role.

This article is a continuation of previous author's work [14] devoted for object detection and segmentation on aerial images.

## III. DATABASES OF SATELLITE IMAGES

A standard database of images is the important part for learning, efficiency evaluation and comparative analysis of different machine learning algorithms. Nowadays, there are some available databases of satellite images.

DeepSat database [15] contains two sets of annotated images from different satellites: 500,000 Sat-4 images, divided into 4 classes (“barren land”, “trees”, “grassland” and “other”), 405 000 Sat-6 images, divided into 6 classes (“barren land”, “trees”, “grassland”, “roads”, “buildings” and “water bodies”). All samples have a size of  $28 \times 28$  px at a spatial resolution of 1 m/px and contain 4 channels (red, green, blue and NIR - near infrared radiation). However, while this dataset is very useful for preliminary preparation of more complex models, it does not allow to take further steps for detailed analysis of developed algorithms. The examples of images from the DeepSat database are shown on Fig. 1.

Inria database [4] contains aerial orthorectified color images, which cover the area of 810 km<sup>2</sup> of 10 cities (180 images and 405 km<sup>2</sup> for training and testing set) with a spatial resolution of 0.3 m. Each photo has a size of 1000x1000 px. All images are divided into 2 classes: “buildings” and “not buildings”. The samples of the database cover dissimilar urban settlements, ranging from densely populated areas to alpine towns. As DeepSat database[15] this dataset is usable to assess the generalization power of techniques of image segmentation. The examples of images from the Inria database are shown on Fig. 2.

DeepGlobe database [12] contains images in RGB format, collected by DigitalGlobe's satellite. Each image has a size of 1024x1024 px. In the training dataset of this dataset, each image contains a mask for road labels. The mask is given in a grayscale format, with white standing for road pixel, and black standing for the background. The labels are not perfect due to the cost of annotating segmentation mask, especially in rural areas. In addition, in most cases small roads within farmlands are not annotated. The training set contains 6246 photos and the validation set contains 1243 photos. The examples of images from the DeepGlobe database are shown on Fig. 3.

DSTL dataset contains 1km x 1km satellite images in GEOTIFF formats. For the first time, this database was provided in Kaggle competition “DSTL Satellite Imagery Feature Detection” [15]. Images of DSTL dataset are labeled on 10 different classes: “buildings”, “manmade structures”, “roads”, “tracks”, “trees”, “crops”, “waterway”, “standing water”, “large vehicles” (e.g. lorries, trucks or buses) and

“small vehicles” (cars, vans or bikes). All 50 samples have a size more than  $3300 \times 3300$  px. In spite of little amount of images, extracting methods allow to crop smaller images. The

examples of images from the DSTL database are shown on Fig. 4.

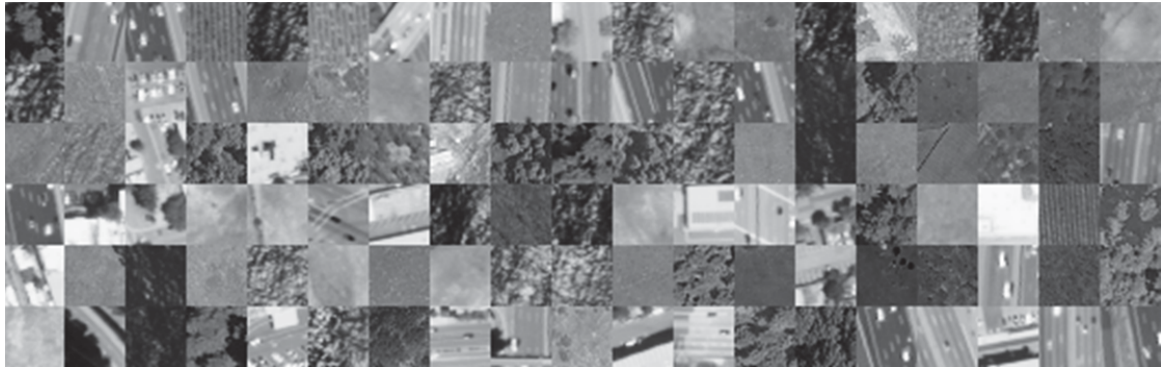


Fig. 1. Examples of images from the DeepSat database

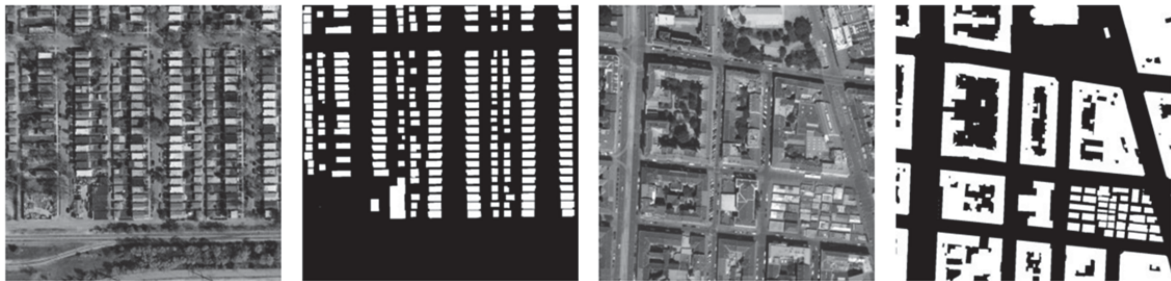


Fig. 2. Examples of images from the Inria database



Fig. 3. Examples of images from the DeepGlobe database



Fig. 4. Examples of images from the DSTL database



Images obtained from the Landsat-8 satellites [20] are used for estimation of automatic object detection quality. Landsat-8 images have a resolution of 30 meters per pixel. This is the highest resolution from open sources aerial images. Also in this research we use PlanetScope satellite imagery provided by Russian Space Systems Agency. The PlanetScope satellite group has 10 times better resolution than Landsat-8 – 3 meters per pixel.

IV. CONVOLUTIONAL NEURAL NETWORK ARCHITECTURES

The paper presents developed models, based on CNNs – special architectures aimed at the rapid and qualitative detection of various objects [1]. CNNs are related to algorithms of deep machine learning, which are popular now to solve most modern problems of computer vision.

To implement comparative analysis of different algorithms for a segmentation of satellite images there were created three models, based on architectures of U-Net [18], SegNet [19] и LinkNet [20] respectively. The research of working of developed models continues the research, which was provided in paper [14]. All created networks were carried out using Keras library with Tensorflow framework as a backend. Keras is an open source library written in Python. It is built on Tensorflow framework and contains numerous implementations of commonly used neural network building blocks such as layers, activation functions, optimizers, and ready tools to preprocess images and text data. Keras offers a higher-level, more intuitive set of abstractions to develop deep learning models regardless of the used computational backend [21]. Moreover, this library allows to train networks on GPU.

As shown on Fig.5 U-Net consists of two parts: an encoder (on the left) and a decoder (on the right). The encoder represents the typical architecture of CNN and contains four blocks of layers. Every such block consists of two convolutional layers with a  $3 \times 3$  filter, following one by one, ReLU activation functions, followed after each convolution,

and a maxpooling operation with a filter size of  $2 \times 2$  in steps of 2 for downsampling. At each step of dimension reducing, the number of channels is doubled. The decoder contains the same amount of blocks as an encoder. Every such block consists of upsampling operation, which reduce the number of channels, using a  $2 \times 2$  filter (deconvolution), merging operation with the corresponding features map from an encoder, two convolutional layers with a  $3 \times 3$  filter and ReLU activation functions, followed after each convolution. The last layer uses a  $1 \times 1$  convolution to match each component vector to classes. In general, the network has 19 convolutional layers, 18 ReLU activation functions, 4 maxpooling operations, 4 upsampling operations and 4 merging operations.

As in the case of U-Net architecture, SegNet has an encoder, a decoder and a final pixelwise classification layer. The architecture of this model is shown on Fig. 6. The encoder consists of 13 convolutional layers, 13 batch normalization, 13 ReLU activation functions, 5 maxpooling functions for downsampling and 5 upsampling functions. All convolutional layers of encoder corresponds to the first 13 convolutional layers in the VGG16 network for object classification. SegNet is initialized by the weights of this network. Each layer of encoder has a corresponding layer in the decoder. So the decoder consists of the same layers as an encoder, except maxpooling functions, which were exchanged to the same number of upsampling operations. The final output layer is a multiclass softmax classifier, which help to predict class probabilities for each pixel independently.

Also there was developed LinkNet-like architecture (TLinkNet) based on the model from paper [18]. The difference between TLinkNet and the network from [19] consists in the absence of one encoder and one decoder block before Encoder Block 3. This fact is explained by the difference in the size of input images corresponding to the problem. As other developed algorithms, TLinkNet has two parts: an encoder and a decoder. Both parts, encoder and decoder, consist of 3 blocks.

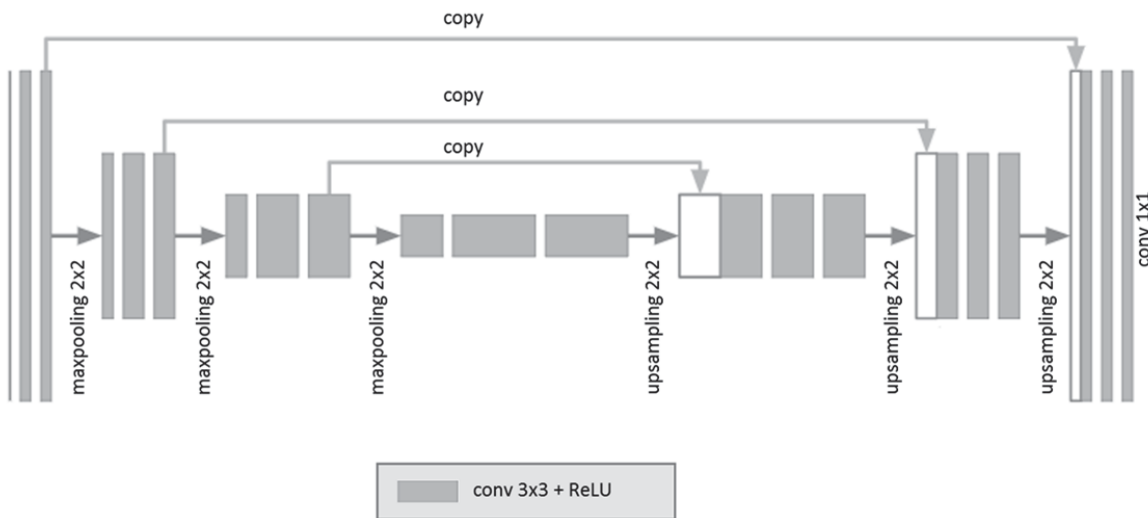


Fig. 5. U-Net neural network architecture

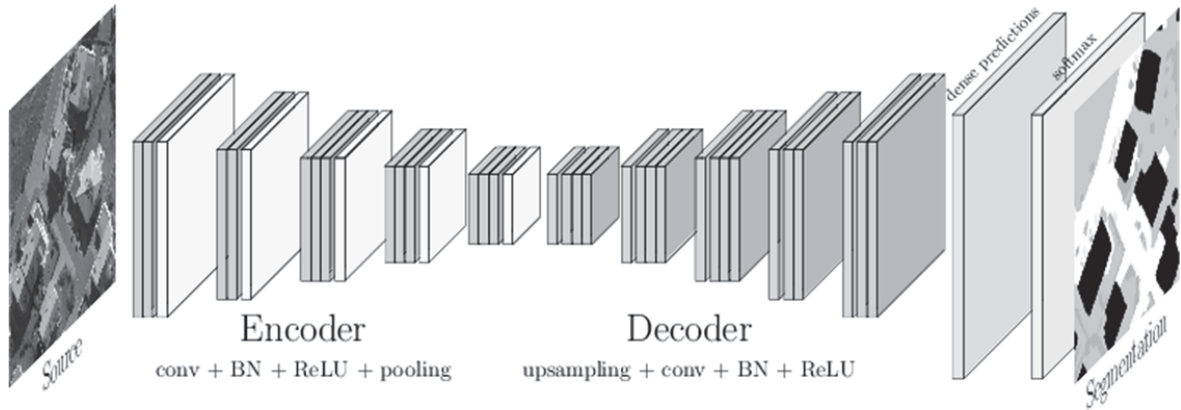


Fig. 6. SegNet architecture

Every encoder block contains 4 convolutional layers, 2 merging and 1 maxpooling operation. According to encoder block, each decoder block has the same architecture, except a maxpooling function, which was exchanged to an upsampling operation. The architecture of modified TLinkNet is shown on Fig. 7.

The approach based on CNNs has high resource consumption. To accelerate neural network operations, the training and testing processes were performed on a large number of independent streams on GPU using parallel computing technology NVIDIA CUDA. This technology is cross-platform and is supported by all modern NVIDIA graphics cards [21]. The developed CNNs were launched on the graphic processor of the video card. The learning rate was set equal to  $10^{-3}$ . As a numerical optimization algorithm, Adaptive Moment Estimation optimizer (Adam) was chosen.

This effective optimizer uses averages and the second moments of the gradients to maintain a learning rate that improves performance on problems with sparse gradients [22]. As a loss function, binary cross entropy function was chosen. This function is a generic approach to combinatorial optimization of weights for machine learning algorithms [23]. On every training iteration the model updated its weights using the batch of 64 samples. The classifier ended its training after completing 256 epochs.

#### V. SIMULATION RESULTS ON DSTL IMAGE DATABASE

Numerical experiments for developed algorithms were performed on images of DSTL database. For an experiment from the initial dataset there were extracted smaller images. In spite of little amount of training images, pulling methods allow to crop smaller images with size of  $160 \times 160$  pixels and corresponding masks, which were received from csv file. As a result the training set contains 3955 photos and the test set contains 600 photos. Train and test samples did not have same pictures. For prepared images there were taken into account only 3 classes: “trees”, “crops” and “waterway”.

The launch of the CNNs was carried out on the supercomputer NVIDIA DGX-1 and lasted around 1 hour. As a result of numerical experiments, accuracy (A) of model was calculated according to the following formula:

$$A = \frac{P}{N}, \quad (1)$$

where  $P$  is a quantity of right classified images and  $N$  is the size of test sample [24]. The results of numerical experiments on validation set cite in Table I.

TABLE I. TESTING RESULTS OF CONVOLUTIONAL NEURAL NETWORKS

Model	Accuracy (A)
SegNet	93.59%
TLinkNet	94.53%
U-Net	94.66%

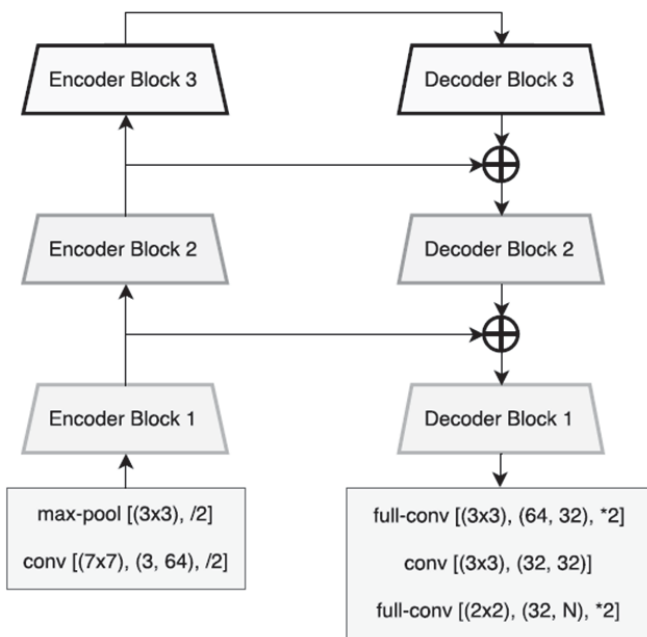


Fig. 7. Transformed LinkNet architecture

According to presented results, all algorithms show high value of accuracy, so this metric is not sufficient to measure the efficiency of networks. This fact can be explained by the little size of objects of different classes for segmentation. For each model the value of accuracy keeps the stability after some epochs (E) of training. Moreover for every network the value of loss function is insignificant, which decreases with the increase of completed training iterations.

As a rule, the quality of algorithms for image segmentation is evaluated by special coefficients for comparing the similarity of predicted and true masks. To estimate developed models there was used Dice similarity coefficient (DSC). This index is binary measure of similarity, possesses the value from [0, 1] and can be calculated by the following formula:

$$DSC = \frac{2I}{S}, \quad (3)$$

where  $I = |X \cap Y|$  is a power of intersection and  $S = |X| + |Y|$  is an sum of powers for real mask  $X$  and predicted mask  $Y$ . In other words,  $DSC$  equals twice the number of common elements to both sets divided by the sum of the number of elements for each set. In our task, numerator  $I$  and denominator  $S$  can be calculated by following formulas

$$I = \sum_{x \in X} \sum_{y \in Y} xy, \quad S = \sum_{x \in X} (x + y), \quad (2)$$

where  $x, y$  are values of pixels from [0, 1] for real mask  $X$  and predicted mask  $Y$  respectively. Graphs of dependency DSC value from the number of epochs for 3 types of neural networks were shown on Fig. 8.

According to testing results presented at Table II, the worst algorithm of image segmentation was SegNet, whereas the best result was shown by U-Net. This fact can be explained by the difficulty of architectures of developed networks. TLinkNet and U-Net architectures throw features from encoder to decoder as opposed to SegNet. This peculiarity of these models

allows to take advantage of using more useful information from input data.

TABLE II. TESTING RESULTS OF CONVOLUTIONAL NEURAL NETWORKS

Model	Dice similarity coefficient (DSC)
SegNet	0.45
TLinkNet	0.68
U-Net	0.75

## VI. SIMULATION RESULTS ON LANDSAT AND PLANETSCOPE DATABASES

For continue the investigations we have manually marked a new set of pictures for retraining a previously created detector based on U-NET convolutional neural network. The marking of satellite images was carried out by 3 independent experts in a web application "Supervise" [26]. A new training set of images contain 30 satellite scenes. We used an average contour for each coordinate in the training of the neural network. Each image in the training set is a tile of a satellite image with a width of 26 km (8600px) and a height of 17.5 km (5800px).

Initially, we were faced with the problem of low detection accuracy of the class "Forest" on satellite imagery of PlanetScope. This was due to the increased detail of tree crowns and the exact boundary of the forest belt. The average accuracy of the detector, which was trained on Landsat-8 satellite imagery, was 73.84% in forest class on satellite imagery PlanetScope. To improve the situation, we retrained the detector on a new training set, which contains images of Landsat-8 and PlanetScope mixed and split into patches for 300 images measuring 224x224 pixels. For the purity of the experiment, the network structure did not change and the training was conducted in 60 epochs.

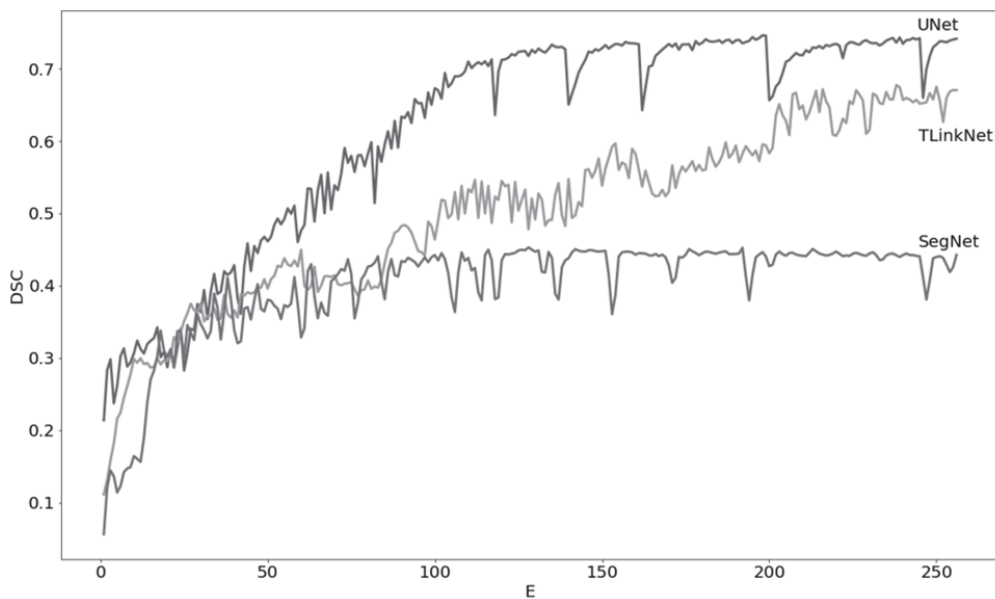


Fig.8. Dependencies of Dice coefficient to training epochs on validation set for convolutional neural networks

The approach for mixing images from 2 allowed to create a universal detector that works on Landsat-8 and PlanetScope. The results of testing a new detector on 3 classes of objects are presented in Table III.

As a result, the class "Forest" has an accuracy of 92.14% on satellite imagery of PlanetScope, which is 18.3% more accurate than the previous detector [14]. The remaining classes of this experiment also increased the accuracy of the detection due to the increase in the accuracy of the boundaries and the reduction in the error price, which decreased by a factor of 10 due to a reduction in the area per pixel of the image from 30m to 3m. Examples of detecting images of "Forest" class are shown in Fig. 9. These figures highlight areas not related to the forest class. The selection of precisely such areas usually causes the main difficulty in such experiments.

TABLE III. TESTING RESULTS OF DETECTOR BASED ON U-NET CONVOLUTIONAL NEURAL NETWORK

Class	Landsat-8		PlanetScope	
	Detection accuracy	Average percentage of intersection	Detection accuracy	Average percentage of intersection
Forest	89.03%	92.54%	92.14%	93.21%
Water	90.87%	81.64%	90.98%	82.14%
Agriculture	94.35%	96.32%	96.52%	96.88%

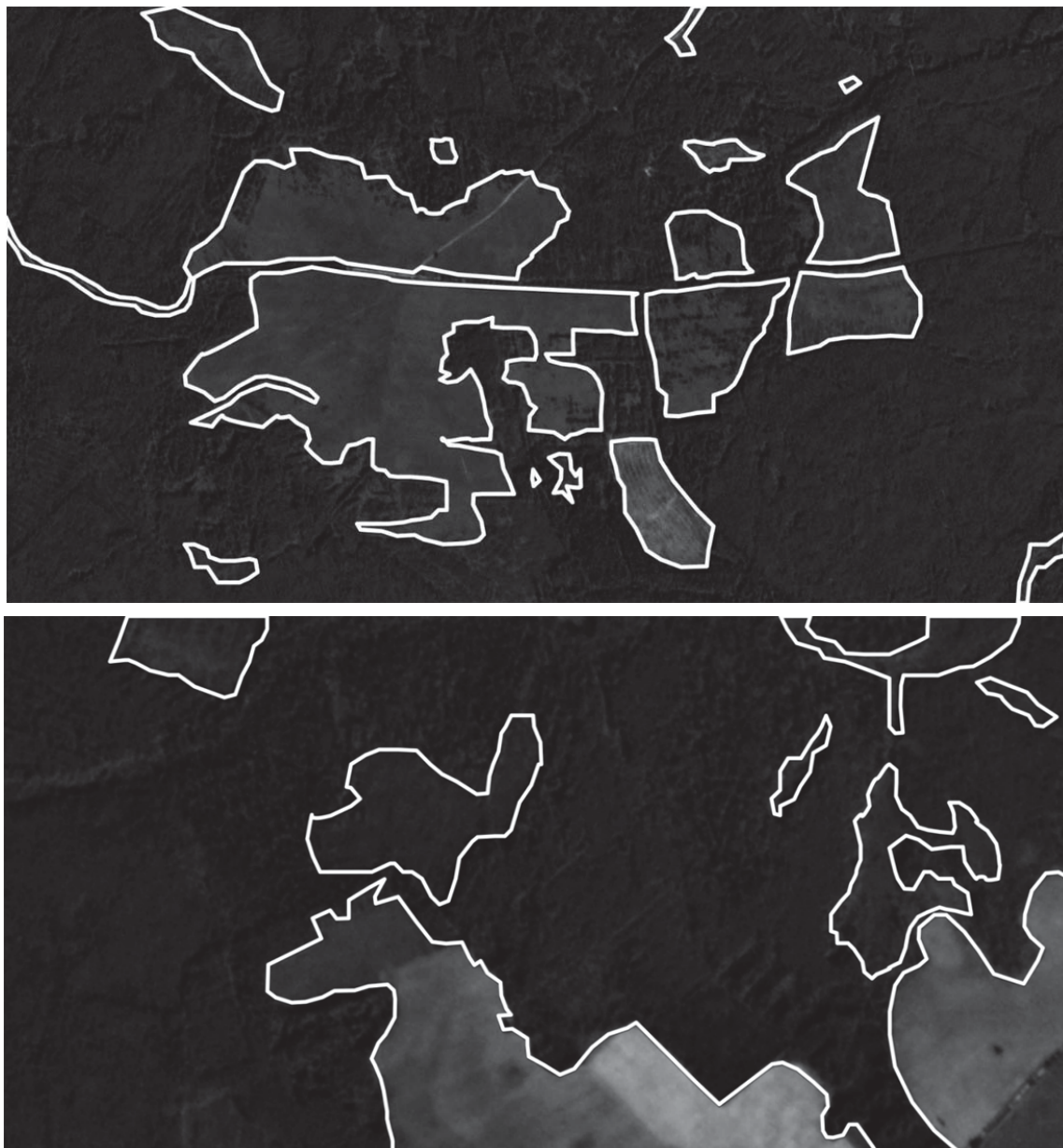


Fig. 9. Examples of detector operation on PlanetScope satellite images of "Forest" class



## VII. CONCLUSION

Experiments of efficiency evaluation for developed algorithms were performed for the aerial photos of DSTL database. Derived results show that the use of complicated CNN allows to increase the quality of segmentation of satellite images. Despite the high values of accuracy (A) for each model, Dice similarity coefficient (DSC), shows the difference in application of various developed algorithms. The greatest value of DSC is equal to 0.75 and was given by using U-Net.

For simulations on another databases we have manually marked a new set of pictures for retraining a previously created detector based on U-NET convolutional neural network. The marking of satellite images was carried out by 3 independent experts in a web application "Supervise". The average accuracy of the detector, which was trained on Landsat-8 satellite imagery, was 73.84% in forest class on satellite imagery PlanetScope. To improve the situation, we retrained the detector on a new training set, which contains images of Landsat-8 and PlanetScope mixed and split into patches for 300 images measuring 224x224 pixels.

As a result, the class "Forest" has an accuracy of 92.14% on satellite imagery of PlanetScope, which is 18.3% more accurate than the previous detector. The remaining classes of this experiment also increased the accuracy of the detection due to the increase in the accuracy of the boundaries and the reduction in the error price, which decreased by a factor of 10 due to a reduction in the area per pixel of the image from 30m to 3m.

## ACKNOWLEDGMENT

The paper was prepared with the financial support of the Ministry of Education of the Russian Federation in the framework of the scientific project No. 14.575.21.0167 connected with the implementation of applied scientific research on the following topic: «Development of applied solutions for processing and integration of large volumes of diverse operational, retrospective and the thematic data of Earth's remote sensing in the unified geospace using smart digital technologies and artificial intelligence» (identifier RFMEFI57517X0167).

The authors are grateful to AI-center of P.G. Demidov Yaroslavl State University for providing access to the supercomputer NVIDIA DGX-1.

## REFERENCES

- [1] Y. Goodfellow, Y. Bengio, A. Courville, "Deep Learning", *The MIT Press*, 2016, 800 p.
- [2] A. Van Etten, "You Only Look Twice: Rapid Multi-Scale Object Detection In Satellite Imagery", Web: <https://arxiv.org/abs/1805.09512>.
- [3] "LandSat Database", Web: <https://landsat.visibleearth.nasa.gov>.
- [4] E. Maggiori, Y. Tarabalka, G. Charpiat, P. Alliez, "Can Semantic Labeling Methods Generalize to Any City? The Inria Aerial Image Labeling Benchmark", *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, 2017.
- [5] "SpaceNet Database", Web: <http://explore.digitalglobe.com/spacenet>.
- [6] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation", *CoRR, abs/1411.4038*, 2014.
- [7] S. Ohleyer, "Building segmentation on satellite images", Web: [https://project.inria.fr/aerialimagelabeling/files/2018/01/fp\\_ohleyer\\_compressed.pdf](https://project.inria.fr/aerialimagelabeling/files/2018/01/fp_ohleyer_compressed.pdf).
- [8] O. Ronneberger, P. Fischer, T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation", *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Springer, LNCS, vol. 9351, 2015, pp. 234-341.
- [9] J. Patravali, S. Jain, S. Chilamkurthy, "2D-3D Fully Convolutional Neural Networks for Cardiac MR Segmentation", *Qure.AI*, 2017.
- [10] G. Chhor, C.B. Aramburu, "Satellite Image Segmentation for Building Detection using U-net", Web: <http://cs229.stanford.edu/proj2017/final-reports/5243715.pdf>.
- [11] A. Buslaev, S. Seferbekov, V. Iglovikov, A. Shvets, "Fully Convolutional Network for Automatic Road Extraction from Satellite Imagery", *CoRR, abs/1806.05182*, 2018.
- [12] "DeepGlobe. CVPR 2018 - Satellite Challenge", Web: <http://deepglobe.org>.
- [13] H. Noh, S. Hong, B. Han, "Learning Deconvolution Network for Semantic Segmentation", *ICCV*, 2015, pp. 1520 - 1528.
- [14] V. Khryashchev, V. Pavlov, A. Priorov, E. Kazina, "Convolutional Neural Network for Satellite Imagery", *Proceedings of the 22th Conference of Open Innovations Association FRUCT22*, Jyväskylä, Finland, 2018, pp. 344-347.
- [15] S. Basu, S. Ganguly, S. Mukhopadhyay, R. DiBiano, M. Karki, R. Nemani, "DeepSat - A Learning framework for Satellite Imagery", Web: <https://arxiv.org/abs/1805.09512>.
- [16] "DSTL Satellite Imagery Feature Detection", Web: <https://www.kaggle.com/c/dstl-satellite-imagery-feature-detection>.
- [17] "Landsat8", Web: <https://landsat.usgs.gov/landsat-8>.
- [18] O. Ronneberger, P. Fischer, T. Brox "U-Net: convolutional networks for biomedical image segmentation", *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Springer, LNCS, Vol. 9351: pp. 234-241, 2015.
- [19] V. Badrinarayanan, A. Kendall, R. Cipolla, "SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39(12), 2017, pp. 2481 - 2495.
- [20] A. Chaurasia, E. Culurciello, "LinkNet: Exploiting Encoder Representations for Efficient Semantic Segmentation", *IEEE Visual Communications and Image Processing (VCIP)*, 2017.
- [21] A. Gulli, S. Pal, "Deep Learning with Keras", *Packt Publishing*, 2017, 320 p.
- [22] J. Sanders, E. Kandrot, "CUDA by Example: An Introduction to General-Purpose GPU Programming", *Addison-Wesley Professional*, 2010, 320 p.
- [23] D. P. Kingma, J. Ba, "Adam: A Method for Stochastic Optimization", Web: <https://arxiv.org/abs/1412.6980>.
- [24] P. T. de Boer., D. Kroese, S. Mannor, R. Rubinstein, "A Tutorial on the Cross-Entropy Method", *Annals of operations research*, vol. 134(1), 2005, pp. 19-67.
- [25] J. VanderPlas, "Python Data Science Handbook: Essential Tools for Working with Data. First Edition", *O'Reilly Media*, 2016, 541 p.
- [26] "Cloud platform for computer vision", <http://supervise.ly>.