

Recognition of handwritten Latin characters with diacritics using CNN

Edyta LUKASIK^{2*}, Malgorzata CHARYTANOWICZ^{1,2}, Marek MILOSZ²,
Michail TOKOVAROV², Monika KACZOROWSKA², Dariusz CZERWINSKI²,
and Tomasz ZIENTARSKI²

¹Systems Research Institute Polish Academy of Sciences, ul. Newelska 6, 01-447 Warsaw, Poland

²Lublin University of Technology, ul. Nadbystrzycka 38D, 20-618 Lublin, Poland

Abstract. Convolutional Neural Networks (CNN) have achieved huge popularity in solving problems in image analysis and in text recognition. In this work, we assess the effectiveness of CNN-based architectures where a network is trained in recognizing handwritten characters based on Latin script. European languages such as Dutch, French, German, etc., use different variants of the Latin script, so in the conducted research, the Latin alphabet was extended by certain characters with diacritics used in Polish language. To evaluate the recognition results under the same conditions, a handwritten Latin dataset was also developed. The proposed CNN architecture produced an accuracy of 96% for the extended character set. This is comparable to state-of-the-art results found in the domain of identifying handwritten characters. The presented approach extends the usage of CNN-based recognition to different variants of the Latin characters and shows it can be successfully used for a set of languages based on that script. It seems to be an effective technique for a set of languages written using the Latin script.

Key words: handwritten documents; diacritics; neural networks; character recognition; deep learning.

1. Introduction

Character recognition using image analysis techniques is one of the most important concepts related to natural language processing. In today's world, the digital form of paper documents facilitates their automatic search for the specified criteria or classification according to the subject matter of their content. In the field of computer-aided document analysis, automatic translation techniques and speech synthesis systems must also be mentioned. Nowadays, the recognition of printed Latin characters is rather easily accomplished utilizing programs widely available on personal computers and can achieve success rates above 98% [1]. In such a practice, the application of dictionary methods has proven to be most effective for error detection [2]. For an incorrectly recognized character the dictionary is searched for the word most similar to the word misread, and this allows for further significant reduction of the error rate. Distinguishing real signatures from fake ones is one of the most important user verification elements. The modeling of flexible neural networks for handwritten signatures preprocessing was illustrated in [3].

However, despite the extremely dynamic development of information technology, papers written in natural languages do not always exist in a form suitable for direct processing by a computer. Indeed, most documents still exist only in the traditional paper version. Furthermore, the development of mobile devices brings online handwritten scripts as the most important

input to smartphones and tablets. Hence, the implementation of computer handwriting recognition is still a challenging field of current research. Unfortunately, in this case the results obtained are no longer of a very high quality and depend very much on the degree of neatness and readability of the handwriting. Character identification in handwritten documents is more difficult as compared to printed fonts due to their wide diversity, worse contrast, and ruling lines. We have discovered that for many documents written in different variants of the Latin script, the implementation of current optical character recognition systems is much more complicated. The problems result from the fact that living European languages are richer than the Latin alphabet – they contain diacritics. These are common in many European languages, for example, Dutch, French, German, Polish, etc. Software modules integrated with devices such as scanners, copiers with the option of scanning documents or multifunction devices do not always handle documents containing special characters from different languages. The simplest and naive solution to the problem of diacritical marks is to omit them. Consequently, this makes some expressions incomprehensible, and the lack of diacritics very often leads to considerable lexical and morphological ambiguity. However, the use of advanced recognition techniques based on artificial intelligence methods, such as convolutional neural networks, successfully overcomes these problems.

The main aim of this work is to assess the effectiveness of a CNN where a network is trained in recognizing handwritten characters based on the Latin script. Special attention is paid to the influence of diacritical marks. In the conducted research, the Latin alphabet was extended by certain characters with diacritics used in the Polish language. To evaluate the recognition results under the same conditions, a handwritten Latin data-

*e-mail: e.lukasik@pollub.pl

Manuscript submitted 2020-06-30, revised 2020-11-09, initially accepted for publication 2020-11-17, published in February 2021

set was also developed. In addition to show that CNNs perform well in the handwritten characters classification task, we demonstrate performance with 97.6% accuracy in the classification of uppercase Latin letters and uppercase Latin letters with diacritics and the accuracy of 95.6% for the extended character set including digits, special characters, uppercase Latin letters and uppercase Latin letters with diacritics.

2. Related works

Recognizing handwriting is a very broad issue. The following research problems can be distinguished in this area: handwritten text recognition, the recognition of handwritten characters within of a text, and handwritten isolated character recognition.

Deep learning techniques based on Artificial Neural Networks (ANNs) are already successfully used in character recognition and font classification. An analysis of the effectiveness of handwritten digits recognition by two ANNs with different architecture and parameters was presented in [4]. MNIST database (Modified National Institute of Standards and Technology database) was used in the experiment [5]. Pal et al. [6], for example, presented the effectiveness of applying ANNs in the task of recognizing handwritten English text. The system proposed by the authors incorporated image processing methods, classification, and character recognition. The network was trained with an error back-propagation algorithm [7, 8] and demonstrated a maximum efficiency of 94%. Related topics were presented by Perwej et al. in [9] and Pradeep et al. in [10], and the accuracy of the proposed algorithms was around 82.5% and 90.19%, respectively. Further experiments, including distinguishing between small and capital letters, were conducted in [11] and achieved an efficiency of 95.62%.

Nowadays, deep CNNs have become the standard for handwritten document recognition. The modern CNNs architectures are more efficient than the ones available a few years ago, but still a lot of improvement is required. Designing efficient architectures of CNN for speeding up their efficiency was presented in [12]. In paper [13], the design of a novel deep convolutional network for the classification of the Latin characters was described. The proposed network achieved an overall accuracy of 96%, which is now one of highest results. A novel offline handwriting recognition algorithm was introduced in [14] by Such et al. The fully convolutional handwriting model took in a handwriting sample of unknown length and outputted an arbitrary stream of symbols. Therefore, a three-staged approach was proposed for the recognition of handwritten characters. First a CNN was trained to quickly predict the word label for common words such as “the”, “her”, “this”, etc. In the second stage, a CNN was trained to predict the number of symbols in a word block. Finally, an arbitrary length sequence of characters from a variable length word block was predicted. The model achieved 92.4% accuracy on a subset of 12,000 word blocks consisting of English and French word groups and special characters generated from the National Institute of Standards and Technology (NIST) dataset [15].

Several approaches to character recognition have been proposed focusing on various applications, but these have not considered the specialty of some languages which contain diacritics. When using Latin script, people for various reasons sometimes write without diacritics, replacing them by the underlying basic character. For practical purposes they use their ASCII counterparts. However, diacritics are not additional or optional, they can reveal very useful information about the font type and are a crucial part of the language. In paper [16], the problem of Arabic font recognition based on diacritics features was considered. Two algorithms for diacritics segmentation were developed, namely flood-fill based and clustering-based algorithms. The experiments conducted proved that the chosen approach can achieve an average recognition rate of 98.73% on a typical database that contains 10 of the most popular Arabic fonts. The Amazigh language transcribed in Latin, which is distinguished by its diacritical characters, was studied by Gajoui et al. [17]. They proposed a system based on neural networks and compared its behavior to a diacritical language and a diacritic-free language with a different quality paper. The work opened an interesting prospect towards developing a particular module for diacritical languages. An analysis of the percentage of words with diacritics in languages using Latin script was presented by Naplava et al. in [18], and this substantiated that problem. They measured that approximately half of the words contain diacritics in languages using Latin script. The actual figures are Czech – 52.5%, Hungarian – 50.7% and Latvian – 47.7%. Around every third word contains diacritics in Polish – 36.9%, Romanian – 31.0%, Turkish – 30.0% and Irish – 29.5%. They then proposed a novel combination of a recurrent neural network and a language model for performing diacritics restoration, especially in Croatian, Slovenian, Serbian, and Czech. The proposed system is language agnostic as it is trained solely from parallel corpora of texts without diacritics.

Research on recognizing single characters was conducted by Grzelak et al. in [19]. They analyzed the effectiveness of Polish handwritten character identification based on neural machine learning technologies. They showed that the original EMNIST (an Extension of MNIST) dataset [20] is not sufficient for training a neural network to properly identify Polish diacritical characters. The conducted experiments indicated that when the EMNIST dataset was extended by two representative classes of Polish diacritics: A_2 and C' , the convolutional neural network gives acceptable responses. However, solutions to the problem of recognizing handwriting in diacritical languages transcribed in Latin script still remain an open issue. Thus, in this paper we investigate the usage of a CNN against Polish handwritten characters that is distinguished by certain characters with diacritics.

3. Materials

Many languages throughout the world use the Latin script, also known as the Roman script, albeit with various modifications depending on the language. The problem with the analysis of recognition methods for Latin handwriting with diacritics lies in available databases. Though there are several classes of specific

handwritten datasets, a large handwritten dataset with adequate diacritical characters does not exist. The most widely known MNIST dataset is derived from a small subset of the numerical digits contained within the NIST Special Database 19 [20]. The MNIST is remarkably easy to access on any platform and it has a very high classification accuracy that was achieved using deep learning. Thus, it has become a standard for extensive learning, and classification and computer vision systems [5], especially to validate and test the effectiveness of neural networks. The NIST Special Database 19 contains digits, as well as uppercase and lowercase handwritten letters collected from over 500 writers. There are 814,255 handwritten characters within the database, including [0–9], [a–z] and [A–Z]. It should be noted, however, that almost half of the total samples are handwritten digits. The complete collection was published in 1995. The re-released 2016 version has a modern file format with a more complex structure. A new classification benchmark has been introduced by the EMNIST dataset, the extended variant of the MNIST dataset [20]. It contains more image samples, more output classes, and a more varied classification task than MNIST, at the same time maintaining its structure and nature. Unfortunately, it cannot be used for Latin script detection and recognition tasks of languages which make use of characters with diacritics.

Diacritical marks are rare in English, but they are common in French, German, Italian, Spanish, Czech, Polish and other languages. Similar to punctuation marks, diacritics are objects of semantic importance. Omission or negligence of diacritics creates ambiguities that cannot be resolved based on a language model alone and is likely to cause problems in further processing activities, including full-text search or machine translation, or can result in misunderstanding because many word-pairs seem to have the same spelling if diacritics are omitted. Therefore, in our paper, in the conducted research, certain diacritical characters used in the Polish language were used to extend the Latin alphabet. To evaluate the recognition system under the same conditions, a handwritten Latin dataset was also developed.

The Polish Handwritten Characters Database (PHCD) is made up of 530,000 image characters written by over 2,000 people, both university students, graduates, and faculty [21]. The writers were of different ages, genders, and backgrounds in terms of field and level of education – which makes the dataset a good fit for research. The main purpose of this work was to create a dataset to be used for written character recognition based on the Latin script with diacritics, by providing training and testing sets so that the influence of the diacritics model could be studied.

The samples were collected through a form designed on the basis of the NIST form and scanned at a resolution of 600 dpi. Following this, abnormal samples were removed, and further pre-processing was carried out. This involved noise removal, conversion of RGB to greyscale and binarization according to Otsu's method [22].

The samples then existed either as separate letters or within a sentence. The separate characters were subsequently extracted, centered, and scaled into a box with a width of 20 px and a height of 32 px [21].

Each character had at least 6,000 manifestations. Sample images of lowercase handwritten characters are shown in Fig. 1, including letters with the diacritics: acute [ć, ń, ó, ś, ź], overdot [ż], little tail [ą, ę] and stroke [ł]. It should be noted, however, that some writers use [z] with a stroke [ẓ] instead of a [z] character with the dot accent. The PHCD in its current form contains over 530,000 images of 32×32 px centered on a 32 px by 32 px square, including: 10 digits [0–9], 26 lowercase and uppercase Latin letters [a–z] and [A–Z], 9 lowercase and uppercase characters with diacritic marks used in the Polish language [ą, ć, ę, ł, ń, ó, ś, ź, ż] and [Ą, Ć, Ę, Ł, Ń, Ó, Ś, Ź, Ż], and 9 special characters [+ , - , : , ; , \$, ! , ? , @ , .]. This makes up an open access database that can be used for academic purposes and for further research free of charge [23].

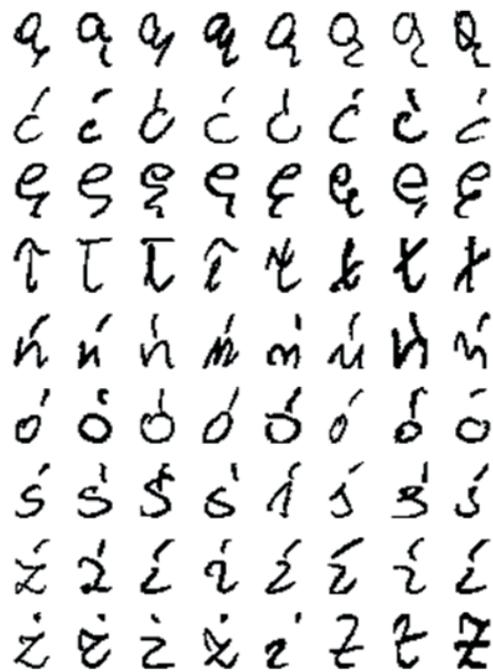


Fig. 1. Sample images of lowercase handwritten characters with diacritics

4. Methodology of CNN

The deep learning used in neural networks is widely applied in the design of pattern recognition systems, particularly for handwriting cognizance. The focus was on the use of CNN to elaborate an Optical Character Recognition (OCR) system that would be suitable for recognizing script built of Latin-based characters and containing diacritical marks. The choice of a classifier is the most important element of the OCR system, and neural networks are an excellent solution for the OCR classification due to their remarkable ability to derive meaning from complicated or imprecise data [17]. A CNN is a variation of a neural network that can extract topological properties from an image [5]. The architecture of this network consists of several main layers: convolutional layers using convolutional operations to

detect features, pooling, or sub-sampling layers to reduce data dimensions and passing this reduction onto the subsequent layer, fully connected layer(s) which connect all neurons of one layer with all neurons of the next layer and followed by a final RBF (Radial Basis Function) classification layer that receives selected input from the previous layer. The features are extracted from the raw image in the first layers and classified in the last layers. The first issue is to determine the order and number of layers, and the second is to choose the appropriate functions and parameters. In each layer, appropriate hyper-parameters should be selected, which are very important in the final effectiveness of the network.

Various network models were tested for both non-diacritical files and diacritical. The research was conducted for the following groups of characters:

- Digits
- Uppercase Latin letters
- Uppercase Latin and uppercase Latin letters with diacritics
- Uppercase all characters: digits, uppercase Latin, uppercase Latin letters with diacritics and special characters
- All characters: lower- and uppercase.

The studied sets were divided into a teaching and testing set in the ratio 80%:20%. Each character occurred about 6,000 times in the set.

The CNN in which the following parameters have been changed were tested:

- Activation function: ReLU (Rectified Linear Unit) or sigmoidal function
- Kernel size of the 2D convolution window: 3×3 or 5×5
- Pooling operation for 2D spatial data: max or average.

All combinations of the indicated parameters were assessed. The best network model was chosen. This achieved the highest accuracy for each of the mentioned data sets. Tests were also carried out for the MNIST collection. The result obtained for MNIST is above 0.990 so the quality of the classification is of the same order as the research results reported in the literature [5].

5. Performance evaluation

In order to produce an optimal configuration, in the conducted studies, the number of layers and coefficient values were changed. However, adding more layers of convolutions did not increase efficiency, while the change of the dropped coefficient did not induce improvement. Its value of 0.25 is optimal. The chosen model with best achievable results, both for a set without diacritical marks and with these marks had the following configuration.

The architecture of the concrete CNN is shown in Fig. 2. The input is a 32×32 binarized matrix. The input is then propagated through 12 adaptable layers. First come two convolutional layers having 32 filters with the size of 3×3 and stride 1. Secondly, the output of the convolutional layer is fed to the ReLU function. The output is down-sampled using a max-pooling operation with a 2×2 stride. Next, the dropout technique is used with the coefficient 0.25. The four operations (two convolutions, nonlinearity, max-pooling, and dropout) are repeated, using 64 filters for the convolutional layers. The output of the last layer is then flattened and fed through a fully connected layer with 256 neurons and ReLU nonlinearities, dropped out with the 0.25 coefficient, and a final output layer is fully connected with a Softmax activation function. The Adam optimizer and the cross-entropy loss function were used in the network [24]. The output is a probability distribution over 89 classes.

The proposed CNN has been implemented in Python using TensorFlow. TensorFlow is an open-source platform for machine learning. This is a library of artificial intelligence and neural networks from Google, characterized by good performance and scalability [25]. TensorFlow has a number of different models and algorithms. In addition, it provides good-quality documentation and an effective TensorBoard tool for displaying and analyzing data flow diagrams describing the calculations performed. The Keras, a high-level API to build and train models which includes first-class support for TensorFlow-specific functionality, was used.

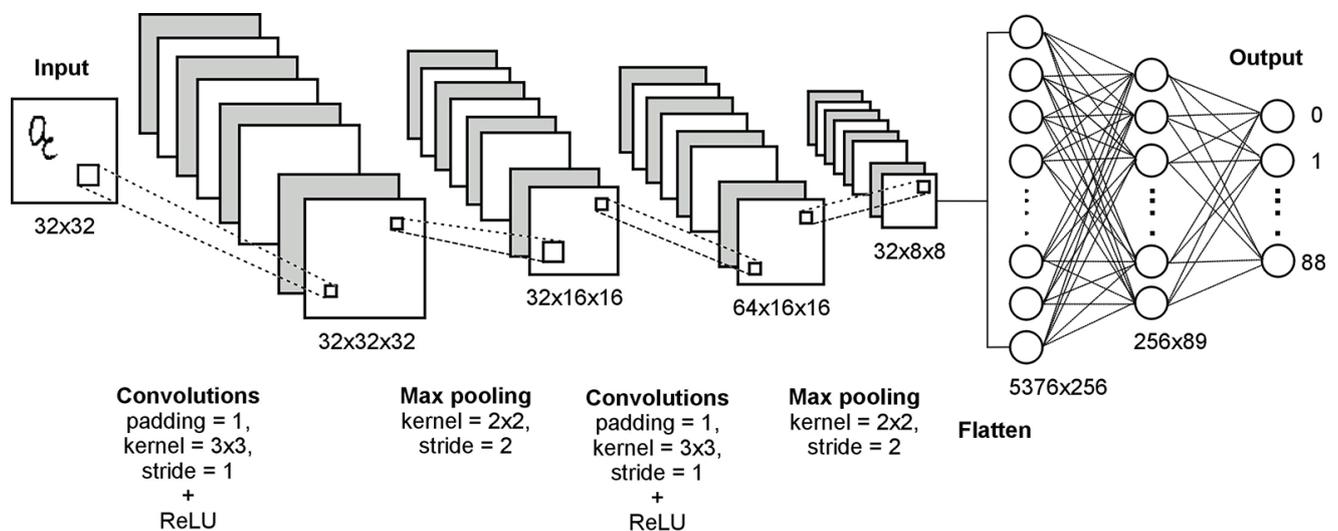


Fig. 2. Scheme of the convolutional neural network

Recognition of handwritten Latin characters with diacritics using CNN

The results of each experiment obtained for the proposed CNN on the PHCD database structure are presented in Figs. 3–6.

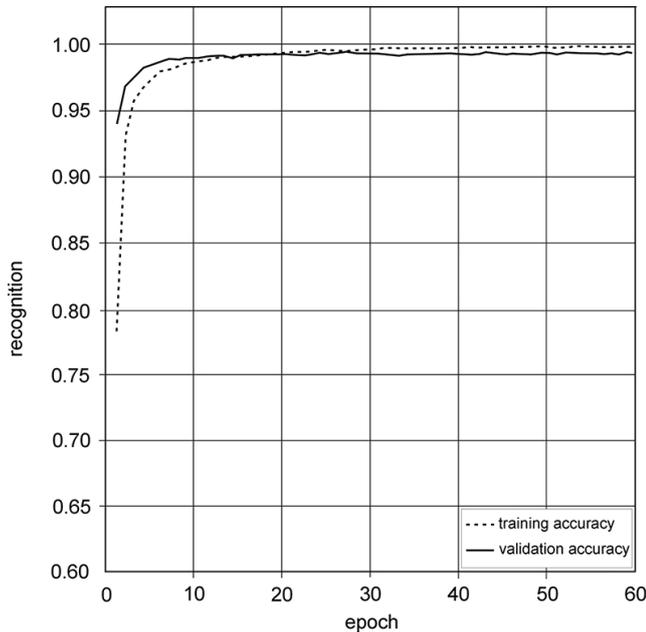


Fig. 3. CNN model accuracy on training and validation sets of digits

The training accuracy of set of digits was 0.998 and validation accuracy was 0.994. The obtained accuracy was very high. The obtained results are at the level of the best achieved on the basis of MNIST [5].

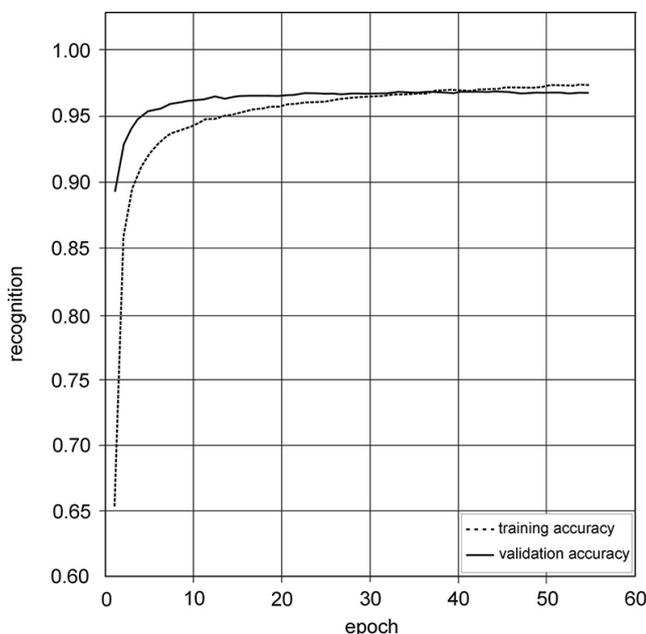


Fig. 4. CNN model accuracy on training and validation sets of uppercase Latin letters

The training accuracy of set of uppercase Latin letters was 0.973 and validation accuracy was 0.967. The accuracy was satisfactory.

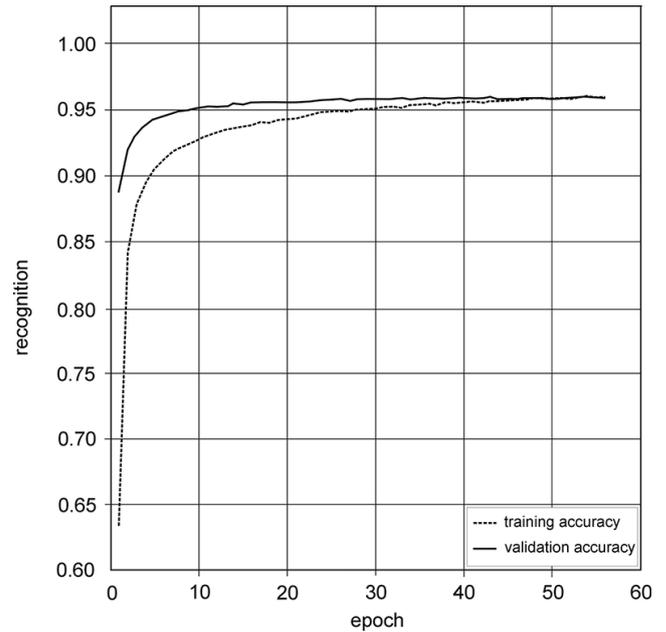


Fig. 5. CNN model accuracy on training and validation sets of uppercase Latin and uppercase Latin letters with diacritics

After adding Latin letters with diacritics to the set of letters, the obtained training accuracy was 0.961 and validation accuracy was 0.960. Accuracy after adding characters with diacritics has dropped very little by thousandths.

A detailed analysis of the recognition rate for each group of characters will be presented in the next section.

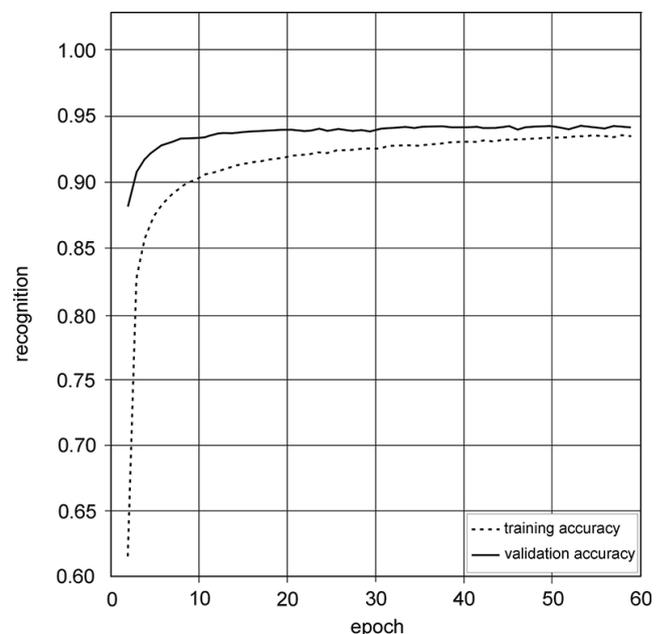


Fig. 6. CNN model accuracy on training and validation sets of uppercase all characters

The training accuracy of the set of all uppercase characters was 0.935 and validation accuracy was 0.942. The obtained accuracy is decent.

The addition of special characters to the set slightly reduced the efficiency of the proposed model to 0.941. For uppercase Latin letters, the result is 0.967. Adding diacritics did not reduce the network efficiency, with the outcome amounting to 0.960. The CNN that we designed achieves an overall accuracy of 96%, which is one of highest results obtained in the literature so far [13].

6. Experimental results

The proposed CNN-based system for Latin handwriting with diacritics was evaluated using the PHCD database, including digits, special characters, lowercase, and uppercase characters, as well as lowercase and uppercase characters with diacritical marks that may appear in the Polish language.

The research was conducted through five experiments, in which random selection occurred with 80% of the entries as the train set and the remaining 20% as the test set. Each of the 89 characters had almost 5,890 manifestations. The first two experiments concerned the initial phase of the study and referred to two commonly used sets of characters: 10 digits [0–9] and 26 uppercase Latin letters [A–Z]. These experiments facilitated a comparison to the results obtained for the well-known handwriting datasets that are used to validate and test the effectiveness of neural networks.

The second phase of the research incorporated characters with diacritical marks. The third experiment was conducted on uppercase Latin letters and uppercase Latin letters with diacritics. This was done to measure the influence of diacritical class in the recognition results.

Finally, digits, Latin letters, and Latin letters with diacritical marks, as well as special characters were considered. Two experiments were performed: the fourth included digits, special characters, and all uppercase letters, and the fifth incorporated lowercase letters.

To describe the results, two measures were considered: average prediction and recognition rate. For all examined characters, the prediction was calculated as the largest value of the vector of real values that sum to 1, turned by Softmax function. Next, the confusion matrix was created, and the recognition rate was calculated as the total number of correctly classified characters divided by the total number of these characters.

In the first experiment, the recognition rates of the PHCD digit test set ranged from 99.6 to 99.9% (see Fig. 7). The highest value was obtained for the digit '6' and the lowest for the digit '4'. The average prediction for correctly classifying digits ranged from 0.994 to 0.999. The results were comparable to the results achieved by most implementations of neural networks for well-known datasets, amounting to 98%-99% for correctly classifying handwritten digits.

In the second experiment, for the uppercase Latin letters, the recognition rates varied from 88.9% to 99.5%, with the lowest value obtained for the letter 'V' (see Fig. 8). The average prediction of this set for correctly classifying characters was in the range from 0.917 to 0.997, with the highest value obtained for the letter 'B' and the lowest value for the letters 'U' and 'V'.

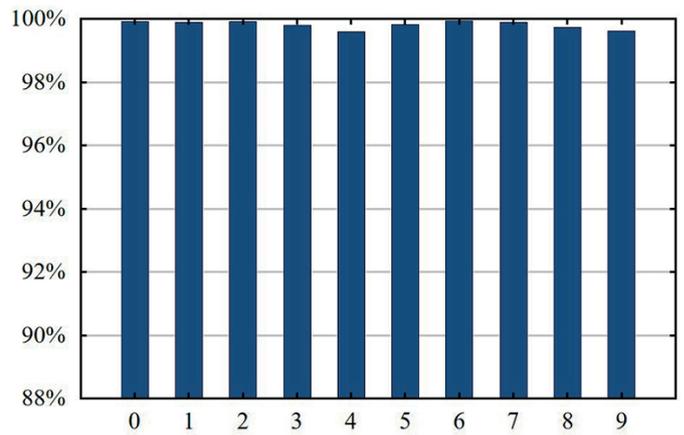


Fig. 7. The recognition rates of the PHCD test set: digits

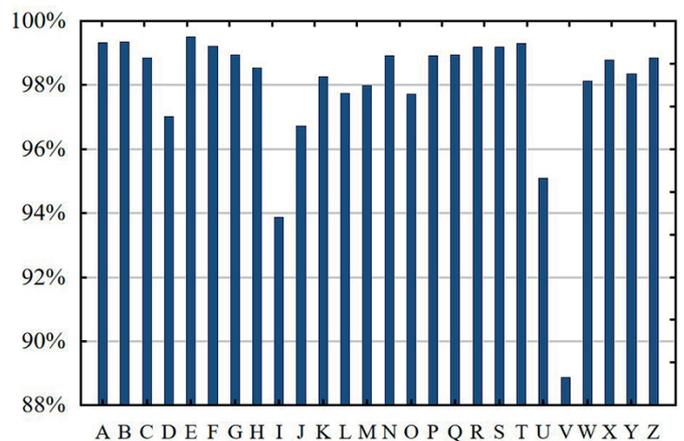


Fig. 8. The recognition rates of the PHCD test set: Latin letters

The third experiment included uppercase Latin letters and uppercase Latin letters with diacritical marks. This did not significantly change the results; the percentage of correctly recognized letters was of the same order and varied from 91.0% to 99.5%. The lowest result was also obtained for the letter 'V' (91%). The average prediction ranged from 0.899 to 0.995, with the highest value obtained for the letter 'F' and the lowest value for the letter 'Ž'.

The recognition rates for Latin letters were comparable to the outcome of Experiment 2. The recognition rates for Latin letters with diacritics are shown in Fig. 9.

In the fourth experiment, a set containing digits, uppercase Latin letters and uppercase Latin letters with diacritical marks, as well as special characters was considered. Therefore, special characters achieved high recognition rates ranging from 93.5% to 100% with the lowest values obtained for colon and semicolon.

The recognition rates for digits and letters excluding two letters: 'O', 'Z' and two digits: 0, 2 ranged from 89.6% to 98.9%. The letters 'O' and 'Z' achieved the worst results because of the high similarity with the digits '0' (recognition rate 84.6%) and '2' (recognition rate 85.8%). The average prediction was in the range from 0.554 to 0.992, with the lowest value obtained for the letter 'O' (0.553) and the digit '0' (0.593). Worse results were

Recognition of handwritten Latin characters with diacritics using CNN

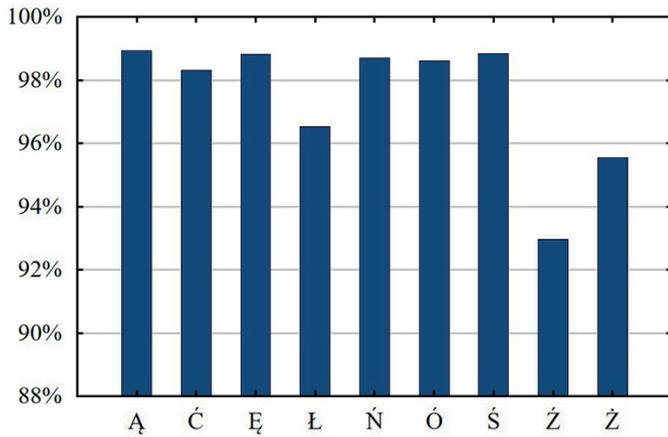


Fig. 9. The recognition rates of the PHCD test set: Latin letters with diacritics

also obtained for the letter 'Z' (0.871) and the digit '2' (0.888). This fact is obvious due to the very high similarity between these characters.

Finally, the fifth experiment was conducted for the extended set that included lowercase letters. Special characters achieved high recognition rates ranged from 95% to 100%. Lower results were caused by a very high similarity between certain pairs of characters. The recognition rates for digits and letters ranged from 46.0% to 99.5%, with the lowest value for the digit '0'. The average prediction for this set was in the range from 0.540 to 0.992, with the lowest value obtained for the digit '0' (0.540) and the letter 'O' (0.576). The recognition rate and the percentage of these characters recognized as comparable characters are presented in Table 1. This similarity is due to the coincident shape of each pair of such characters, but their meaning is completely different.

Table 1

Recognition rate of selected characters and percentage of these characters recognized as of similar character

Recognized character	Recognition rate	Similar character	Percentage of examined characters recognized as of similar character
0	46.14%	O	50.77%
2	84.84%	Z	13.52%
g	60.56%	9	24.98%
l	67.87%	I	15.80%
O	67.62%	0	28.17%
V	84.69%	U	13.17%
Z	84.69%	2	11.29%

The addition of lowercase letters worsened the recognition results due to their similarity to the corresponding uppercase letters and vice versa. Table 2 presents the results for such characters. In this situation, during the recognition process, the problem of incorrect capitalization is detected.

Table 2

Recognition rate of selected characters and percentage of these characters recognized as the lowercase characters and vice versa

Recognized character	Recognition rate	Similar character	Percentage of examined characters recognized as the lowercase characters and vice versa
Ć	69.89%	ć	28.78%
K	82.40%	k	11.73%
Ó	67.16%	ó	30.16%
P	64.01%	p	33.92%
Ś	73.04%	ś	22.83%
Ź	57.77%	ź	32.08%
Ż	72.10%	ż	21.37%
ć	71.72%	Ć	25.88%
k	80.09%	K	12.70%
ó	71.98%	Ó	23.45%
p	66.90%	P	30.27%
ś	63.39%	Ś	30.81%
y	75.78%	Y	13.27%
ź	73.26%	Ź	17.71%
ż	65.55%	Ż	28.00%

The recognition rates for Latin letters are shown in Figs. 10 and 11. These characters achieved mainly high recognition rates, lower results were caused by a very high similarity between certain pairs of characters (see Tables 1 and 2).

The detailed results, including the recognition rate and average prediction for particular characters, are presented in Appendix 1. The confusion matrix for all characters has been added to Appendix 2.

To summarize, the proposed CNN recognizing achieved 99.7% accuracy for the set of 10 digits [0–9] and 98.1% accuracy for the set of 26 uppercase Latin letters [A–Z]. A very

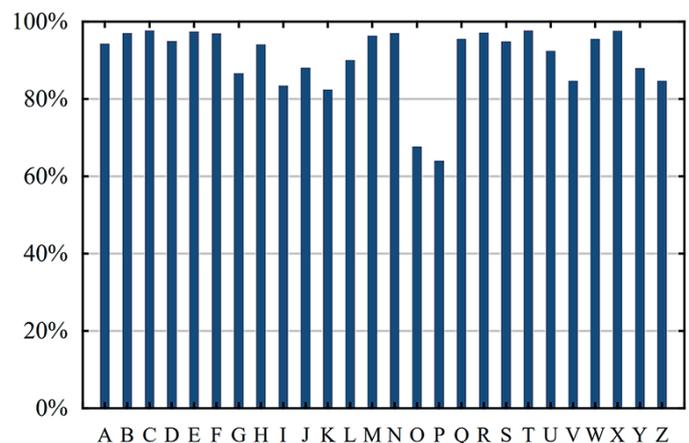


Fig. 10. The recognition rates of the PHCD test set: uppercase Latin letters

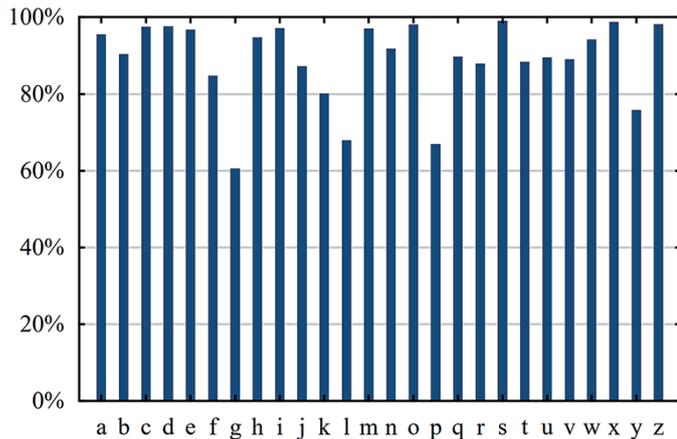


Fig. 11. The recognition rates of the PHCD test set: lowercase Latin letters

high recognition accuracy of 97.6% was achieved for the set of uppercase Latin letters and uppercase Latin letters with diacritics. All Latin letters with diacritics were very well recognized and did not lowered the accuracy. Finally, the recognition accuracy for the sets of all characters including digits, special characters, uppercase Latin letters and uppercase Latin letters with diacritics was evaluated in Experiment 4 and a result of 95.6% was achieved. Special characters were also well recognized; however, the result was lowered by pairs of the similar looking characters: colon and semicolon, 'O' and '0', and 'Z' and '2'. The last experiment was conducted on all characters used in Experiment 4 and incorporated lowercase Latin letters and lowercase Latin letters with diacritics. The recognition accuracy of 90.5% was a bit worse because of a very high similarity between fourteen pairs of characters, especially certain lowercase, and uppercase letters. For such a problem one should consider uppercase and lowercase representations of corresponding characters as belonging to the same classes and then the comparable result can be achieved.

Overall, it can be said that the results obtained for the proposed CNN structure are better than those obtained in study [13]. In the presented article, for Latin letters, the accuracy of 96.7% was achieved for 26 classes, while for them, it was 96% for 23 classes. In [6], an accuracy of 94% was achieved for English and 26 classes, whereas in [9] it was 82.5%. In the ANN for 52 characters of the handwritten English alphabet (lower- and uppercase) in [11], the accuracy was 95.62%. In our work, the accuracy of the designed and tested network on a set consisting of the uppercase Latin letters and the uppercase Latin letters with diacritical marks reached 96.02%.

7. Summary and final comments

Nowadays, deep learning techniques based on Convolutional Neural Networks have taken on an important role in handwritten text recognition. In this paper, a closer look was laid on the recognizing process of handwritten Latin script and especially of the diacritical marks that are a defining feature of many

European languages. The experimental research was conducted for two purposes: to ascertain the effectiveness of CNNs in enabling Latin script recognition and to make known the influence of diacritics on recognition accuracy.

The whole procedure has demonstrated a number of positive features in its approach to Latin handwriting recognition, including:

Lack of changes in the architecture of the CNN for the extended character set, including diacritics (as demonstrated, the network configurations known to work for standard data set can be used directly)

Complete experimental analysis of the impact of different classes of characters on the recognition results

The possibility of adding new classes of characters with the same level of recognition accuracy

The possibility of generating a proper choice of character classes despite the very high similarity between certain pairs of characters, especially certain lowercase, and uppercase letters

Bringing to light the large experimental significance for further research of the Polish Handwritten Characters Database that contains over 530,000 of separate .png files of character images

It should be emphasized that the proffered model is flexible and can be easily extended to other European languages that are based on Latin script. With this aim in mind, we investigated designing a European multi-languages recognition test, with the significance of diacritics being an obstacle in hand-writing recognition. With regard to the future directions of the study, it is worth developing the methods of anomalies detection and character discernment of handwritten script so as to access world knowledge that is archived in this form.

REFERENCES

- [1] E. Lukasik and T. Zientarski, "Comparative analysis of selected programs for optical text recognition", *J. Comput. Sci. Inst.* 7, 191–194 (2018).
- [2] P. Kusaj, M. Kosyra, and M. Charytanowicz, "Web-Page Classification Based on Wikipedia Structure. Recent Developments" in *Mathematics and Informatics, Contemporary Mathematics and Computer Science 2*, Part II, A. Zapała (red.), pp. 89–102, Wydawnictwo KUL, 2016.
- [3] D. Połap and M. Woźniak, "Flexible neural network architecture for handwritten signatures recognition", *Int. J. Electron. Telecommun.* 62, 197–202 (2016).
- [4] M. Milosz and J. Gazda, "Effectiveness of artificial neural networks in recognising handwriting characters", *J. Comput. Sci. Inst.* 7, 210–214 (2018).
- [5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition". *Proc. IEEE* 86(11), 2278–2324 (1998).
- [6] A. Pal and D. Singh, "Handwritten English character recognition using neural network", *Int. J. Comput. Sci. Commun.* 1(2), 141–144 (2010).
- [7] B.K. Verma, "Handwritten Hindi character recognition using multilayer perceptron and radial basis function neural network", *IEEE International Conference on Neural Network* 4, 2111–2115 (1995).

Recognition of handwritten Latin characters with diacritics using CNN

- [8] D. Singh, S.K. Singh, and M. Dutta, “Hand written character recognition using twelve directional feature input and neural network”, *Int. J. Comput. Appl.* 1(3), 94–98 (2010).
- [9] Y. Perwej and A. Chatirvedi, “Neural networks for handwritten English alphabet recognition”, *Int. J. Comput. Appl.* 20(7), 1–5 (2011).
- [10] J. Pradeep, E. Srinivasan, and S. Himavathi, “Neural network based handwritten character recognition system without feature extraction”, *2011 International Conference on Computer, Communication and Electrical Technology (ICCCET)*, Tamilnadu, 2011, pp. 40–44.
- [11] A.M. Obaid, H.M. El Bakry, M.A. Elodusuky, and A.I. Shehab, “Handwritten text recognition system based on neural network”, *Int. J. Adv. Res. Comput. Sci. Technol.* 4(1), 72–77 (2016).
- [12] V. Lebedev and V. Lempitsky. “Speeding-up convolutional neural networks: A survey”, *Bull. Pol. Ac.: Tech.* 66(6), 799–810 (2018).
- [13] D. Firmani, P. Merialdo, E. Nieddu, and S. Scardapane, “In codice ratio: OCR of handwritten Latin documents using deep convolutional networks”, in *AI* CH@ AI* IA, 2017*, pp. 9–16.
- [14] F.P. Such, D. Peri, F. Brockler, P. Hutkowski, and R. Ptucha. “Fully convolutional networks for handwriting recognition”. In: *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, IEEE, 2018, pp. 86–91.
- [15] P. Grother, “NIST special database 19 handprinted forms and characters database”, National Institute of Standards and Technology, Tech. Rep., 1995.
- [16] M. Lutf, X. You, Y. Cheung, and C.L.P. Chen, “Arabic font recognition based on diacritics features”, *Pattern Recognit.* 47, 672–684 (2014).
- [17] K.E. Gajoui, F.A. Allah, and M. Oumsis, “Diacritical Language OCR based on neural network: Case of Amazigh language”. *Procedia Comput. Sci.* 73, 298–305 (2015).
- [18] J. Náplava, M. Straka, P. Straňák, and J. Hajič, “Diacritics Restoration Using Neural Networks”, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC)*, 2018.
- [19] D. Grzelak, K. Podlaski, and G. Wiatrowski, “Analyze the effectiveness of an algorithm for identifying Polish characters in handwriting based on neural machine learning technologies”, *Journal of King Saud University – Computer and Information Sciences*, 2019, doi: 10.1016/j.jksuci.2019.08.001.
- [20] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, “EMNIST: an extension of MNIST to handwritten letters”. Retrieved from: <http://arxiv.org/abs/1702.05373>, 2017.
- [21] M. Tokovarov, M. Kaczorowska, and M. Milosz, “Development of Extensive Polish Handwritten Characters Database for Text Recognition Research”, *Adv. Sci. Technol. Res. J.* 14(3), 30–38 (2020), doi:10.12913/22998624/122567.
- [22] M. Charytanowicz and P. Kulczycki, “An Image Analysis Algorithm for Soil Structure Identification”, in: *Intelligent Systems’2014*, pp. 681–692, D. Filev, J. Jablkowski, J. Kacprzyk, I. Popchev, L. Rutkowski, V. Sgurev, E. Sotirova, P. Szykarczyk, S. Zadrozny (eds.), Springer, Berlin, 2014.
- [23] The Polish Handwritten Characters Database, [Online]. <https://cs.pollub.pl/phcd/?lang=en>.
- [24] D.P. Kingma and J.L. Ba, “Adam: A method for stochastic optimization”. arXiv:1412.6980v9, 2014.
- [25] M. Abadi *et al.*, “Tensorflow: A system for large-scale machine learning.” in *12th Symposium on Operating Systems Design and Implementation*, 2016, pp. 265–283.

Appendix 1

 Table 1
 Recognition rate and average prediction for digits, special characters, lowercase and uppercase letters, lowercase and uppercase letters with diacritics

Char	Recognition rate [%]	Average prediction	Char	Recognition rate [%]	Average prediction	Char	Recognition rate [%]	Average prediction	Char	Recognition rate [%]	Average prediction
0	46.14	0.540	N	97.03	0.972	+	100.00	1.000	n	91.79	0.943
1	96.46	0.984	O	67.62	0.576	-	100.00	1.000	o	98.05	0.992
2	84.84	0.910	P	64.01	0.674	:	95.61	0.962	p	66.90	0.651
3	97.21	0.989	Q	95.55	0.977	;	94.93	0.953	q	89.62	0.913
4	91.01	0.943	R	97.11	0.985	\$	98.11	0.984	r	87.90	0.955
5	93.57	0.961	S	94.82	0.936	!	99.92	1.000	s	98.93	0.992
6	93.89	0.936	T	97.72	0.932	?	99.24	0.997	t	88.37	0.875
7	96.94	0.981	U	92.36	0.926	@	99.18	0.996	u	89.42	0.907
8	95.59	0.972	V	84.67	0.930	.	100.00	0.999	v	89.04	0.862
9	90.36	0.764	W	95.45	0.982	A	95.43	0.980	w	94.15	0.972
A	94.30	0.968	X	97.58	0.986	B	90.31	0.943	x	98.70	0.991
B	97.04	0.980	Y	87.90	0.855	C	97.45	0.985	y	75.78	0.878
C	97.70	0.978	Z	84.69	0.865	D	97.53	0.977	z	98.15	0.992
D	94.91	0.963	Ą	96.30	0.976	E	96.71	0.971	ą	96.87	0.985
E	97.40	0.976	Ć	69.89	0.727	F	84.73	0.933	ć	71.72	0.705

