

# Scene Text Script Identification with Convolutional Recurrent Neural Networks

Jieru Mei<sup>†</sup>, Luo Dai<sup>\*</sup>, Baoguang Shi<sup>\*</sup>, Xiang Bai<sup>\*</sup>

<sup>†</sup>School of Automation, <sup>\*</sup>School of Electronic Information and Communications  
Huazhong University of Science and Technology, Wuhan, China 430074  
{meijieru, dailuo, shibaoguang}@gmail.com, xbai@hust.edu.cn

**Abstract**—Script identification for scene text images is a challenging task. This paper describes a novel deep neural network structure that efficiently identifies scripts of images. In our design, we exploit two important factors, namely the image representation, and the spatial dependencies within text lines. To this end, we bring together a Convolutional Neural Network (CNN) and a Recurrent Neural Network (RNN) into one end-to-end trainable network. The former generates rich image representations, while the latter effectively analyzes long-term spatial dependencies. Besides, on top of the structure, we adopt an average pooling structure in order to deal with input images of arbitrary sizes. Experiments on several datasets, including SIW-13 and CVSI2015, demonstrate that our approach achieves superior performance, compared with previous approaches.

## I. INTRODUCTION

Script identification is one of the most important components of Optical Character Recognition, which depends on the language of scripts. In this area, script identification on documents [1], [2], [3], [4] or videos [5], [6] is the focus of the community. In document script identification, many approaches have been proposed. Spitz et al [1] use spatial relationships of features related to the upward concavities in character structures and distribution of optical density for page-wise script identification. Chaudhuri et al. [2] propose the use of projection profile and spatial features for Indian language text level script identification. Hochberg et al. [3] propose a new method using cluster-based templates to identify unique characteristic shapes. Tan et al. [4] propose a new method using rotation invariant texture features for identification of Chinese, English, Greek, Russian, Malayalam and Persian text. In the area of scene text recognition, Gomez et al. [7] propose a novel method based on the use of ensembles of conjoined networks to jointly learn discriminative stroke-parts representations and their relative importance in a patch-based classification scheme. Singh et al. [8] take use of mid-level features which are pooled from densely computed local features and off-the-shelf classifier to identify the script of the text image. They all achieve great performance.

All the methods mentioned above specifically analyze the text script identification on printed document images. However, script identification for scene text images remains a challenging task, because of the complicated backgrounds, different fonts and colors of texts etc.

In this paper, to deal with the challenges in scene text script identification, we exploit two important factors, namely the

image representation, and the spatial dependencies within text lines.

Image representation is a powerful description of text images, because of its complexity and high-dimensionality. We make use of image representation to handle various fonts and colors of texts and complicated backgrounds. We adopt Convolutional Neural Network to generate image representation.

Many languages, e.g. English, Russian and Greek, share same subset of letters. So directly using holistic representation produced by convolution layers to identify script in these three languages can be difficult. However, we can exploit the spatial dependencies within text lines to solve this problem. Although having many same letters, languages have unique usages of letters. For instance, "K" and "H" are seldom seen together in words in English. But in Russian, this is a common usage. In this way, "K" and "H" have distinct spatial dependencies in English and Russian. Under this scenario, analysis of spatial dependencies within text lines can help us identify text containing "KH" as Russian. In addition, in some languages, there are letters very similar with letters in other languages. Due to the minor differences, holistic representation often has poor performance when handling text images with this kind of letters. However, spatial dependencies of this kind of letters are unique compared with the similar letters in other languages. This uniqueness provides vital information to identify the text images. Figure 1 shows examples of these two situations. Recurrent Neural Network is used to analyze long-term spatial dependencies.

In this paper, we propose a novel approach for scene text script identification, which combines Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN). The CNN generates rich image representations, and RNN effectively analyzes long-term spatial dependencies.

Previously, hand-crafted features such as SIFT [9], HOG [10], are widely used in image classification, followed by an encoding method such as Autoencoder [11] and locality-constrained linear coding (LLC) to reduce the dimension of the feature vector. Recently, CNN has been used widely.

CNNs are one type of neural networks strengthened by the inclusion of operation of convolution. CNNs have a substantially more sophisticated structure than standard representations, comprising several layers of non-linear feature extractors. With a stacked structure, CNN's capacity can be controlled by varying depth and breadth. Features produced



Fig. 1. Description of the effective spatial dependencies. As we can see from the image (a), "K" is not often adjacent to "H" in English. But in Russian, many words have these letters adjacent. This can be discriminative for separating these two languages with each other. In image (b), the letter within the red bounding box can be easily confused with "H". However, with the help of the context including previous letter "T" and later letter "H", it can be correctly recognized. Because of the uniqueness of the letter in Russian, it provides vital information to classify the image to be Russian.

by CNN are learned directly from data not handcrafted, with lower dimension and more discriminative details.

CNNs have achieved great success in multiple fields including image classification [12], [13], [14], [15], object detection [16], [17], and image semantic segmentation [18], [19], etc. With the great capability of CNN, using it on script identification is natural.

In our model, we also use Recurrent Neural Networks. Recurrent neural networks are feedforward neural networks with edges that connect adjacent time steps, introducing a notion of time to the network. Recurrent neural network has a wide applications in NLP [20], [21], speech recognition [22] and image caption [23] due to its capability of capturing the dependencies of the sequence.

In comparison with previous sequence analyzing models like Hidden markov model (HMM) which is a generative model, RNN is discriminative. RNN in general performs better than conventional models for discriminating task. What's more, previous models can only produce non-normalized maximum likelihoods while discriminative model's outputs are normalized. At last, discriminative functions learned from data are more powerful than mixture diagonal Gaussian used in standard HMM which model the relative independent variables but limiting its own performance.

We evaluate our new model on widely used SIW-13 and CVSI2015 datasets. The result demonstrates that our approach achieves superior performance compared with previous conventional methods.

This paper is organized as follow: In section II, we describe our novel end-to-end architecture for script identification. The implementation details and improvements analysis are described in section III. Finally, we draw our conclusion in section IV.

## II. METHODOLOGY

The architecture of the end-to-end network mainly consists of three components.

In order to acquire accurate translation invariant image representations, we use a stacked convolutional layers structure as the first part of our model. After the convolutional layers, RNN layers take arbitrary length feature maps produced by CNN layers as inputs to exploit the spatial dependencies in script images. To bring together the output of RNN layers,

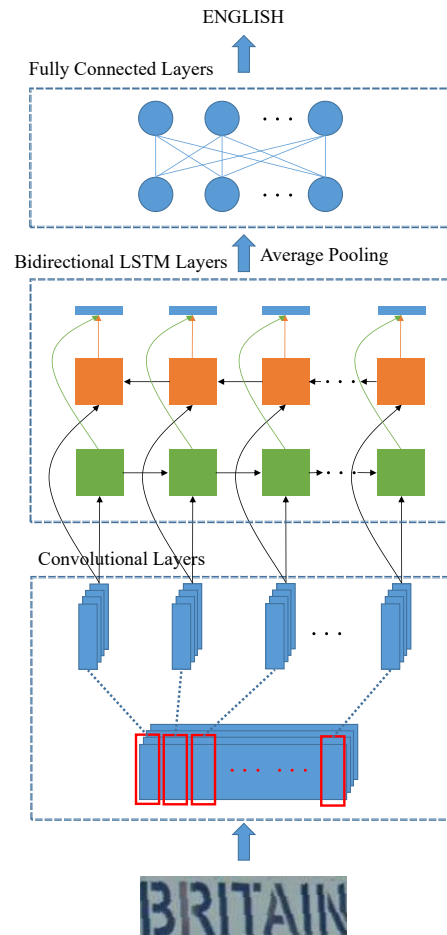


Fig. 2. Overall architecture of our method

we add an average pooling layer after RNN layers. The third part is a fully connected layer widely used in classification problems. The overall architecture of our method is described in Figure 2. As the input of the whole model, the text images are rescaled to have fixed height while having varied width to keep same aspect ratios.

Containing three components, our model can still be trained end-to-end by back-propagation due to the differential property of all these components.

### A. Image Description

As we have mentioned, image description is of great significance for our architecture. Previous widely used SIFT or HOG are hand-crafted and cannot capture the task specialized feature which may be discriminative for special task. Features directly learned from data are supposed to outperform the hand-crafted features. This advantage leads to our usage of a trainable feature generator.

Though trainable, multilayer perceptrons (MLP) suffer from the huge amount of parameters which make the model hard to train. Another disadvantage of MLP is that it treats the images as vectors, ignoring the 2D property of natural images.

Convolutional neural network (CNN), a variant of MLP, is inspired by biological processes, to simulate the role of visual vertices. With local connectivity and shared weights, it reduces parameters dramatically compared with MLP, while catching 2D property of natural images.

With this great power, adopting it in our architecture is natural.

Nonetheless, direct use of state-of-art CNN model in script identification is not possible due to the extreme length variance of the aspect ratio. All images need to be rescaled to the same size because inputs of fully connected layer require fixed dimensions. However, rescaling the image in both sides and feeding them into a holistic CNN classifier can extremely worsen the image description, leading to poor performance.

The structure of the convolutional layers adopted in our method is similar to the most commonly used structure proposed by Lecun et al. [12]. We stack the convolutional layer and max pooling layer but drop the fully connected layer due to the reason described above.

We use this CNN structure to extract the image representations of the text images. Then the image representations are fed into the following LSTM layers. Images are required to be rescaled to the same height. We keep the ratio of height and its width for consistent.

Each column vector in the feature maps actually corresponds to a region in the original input text image. This region, known as receptive field, is always associated with a vector, which can be regarded as local descriptor of the region. As the operations of convolution and max pooling are translation invariant, these image descriptions are robust for small distortion.

### B. Spatial dependencies analysis

Previous works deal with the problem that feature comes with different size with different approaches. Shi et al. [24] adopts the feature extractor described above, followed with a discriminative clustering method for finding the discriminative patches proposed by Singh et al. [25]. SPPnet [17] proposes a novel pooling method named Spatial Pyramid Pooling to fight against this problem.

Though these works fix the problem elegantly, they fail to take the spatial dependencies within text lines into consideration, which may be discriminative for script identification as explained in Figure 1.

Instead of the methods mentioned above, we propose a novel solution for this problem. Recurrence Neural Network (RNN) models are designed to deal with sequences, and it doesn't require fixed length of input text images, which naturally solve the problem while capturing the spatial dependencies within text lines.

Given a set of input vector  $X = (x_1, x_2, \dots, x_T)$ , The most common method to use RNN is :

$$h_t = f(x_t, h_{t-1}), \quad (1)$$

The hidden state  $h_t$  depends not only on the current input  $x_t$ , but also previous time hidden state  $h_{t-1}$  stored in the

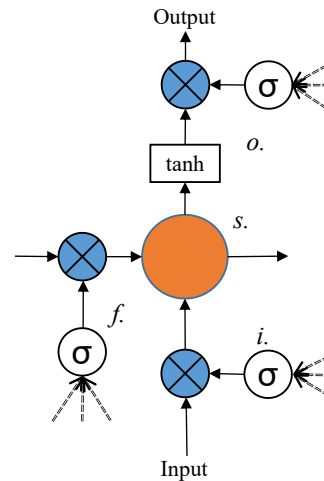


Fig. 3. The architecture of the LSTM block, i, o and f stand for input gate, output gate and forget gate respectively. By introduction these gates to RNN, gradients can be kept for a long time in the same order of magnitude. In this way, the problem of gradient vanishing or explosion is fixed.

RNN block, which gives RNN ability to catch time/spatial dependence in sequences. Output is linear transform of the hidden states.

Although having great power to resolve sequence-based problems, RNN suffers from gradient vanishing or gradient explosion during back-propagation [26], limiting the capability of dealing with relatively long time/spatial dependence. Gradient vanishing and explosion are bottle-necks in this task, because text images are tend to be long in script identification. Long Short Term Memory (LSTM) [27] is proposed to overcome it. LSTM solves this problem by introducing three gates: input gate, forget gate, output gate. They are integrated together to capture spatial dependencies for a long time by keeping the gradient in the same order of magnitude during back-propagation. The input gate controls how much information in input should be saved in hidden state. The output gate determines how much information in hidden state should be the output for current time stop. The forget gate is used to control what information to drop from hidden state. All these gates are controlled according to the current input and previous hidden state. See Figure 3 for vivid illustration.

Back Propagation Trough Time (BPTT) is adopted for learning of the parameters. In practice, gradients are always truncated for simplicity without hurting the performance evidently.

To be clear, RNN and LSTM rely on the assumption that input sequences have a direction. If we only use forward RNN and LSTM, the model can only capture the dependencies from left to right. However, right to left dependencies can be helpful. We refer Figure 1 for a concrete sample. In order to fix this problem, we adopt Bidirectional LSTM.

Single layer of LSTM may not have enough abstraction ability. This is proved in the following experiments. Following method in Shi et al. [28], LSTM layers are stacked for better performance.

### C. Classification

The outputs of LSTM have varied length and fixed dimensions which are determined by previous RNN layers. For simplicity, we combine the features generated by the bidirectional LSTM by average pooling. Average pooling layer takes all of the spatial dependencies within text lines into consideration. To be clear, the bidirectional LSTM weights the features implicitly, so there is no need for additional weights.

Finally, a softmax layer, which outputs the normalized probabilities for each class, is built on top of the LSTM layers. The cross-entropy classification loss function is adopted because of its requirement for confident prediction which has gradient even when every sample has been correctly classified.

## III. EXPERIMENTS

We evaluate our approach on common script identification dataset, including SIW-13 [24] and CVSI2015 [29] which have standardized splits. We also compare our method with previous state-of-art methods on these datasets. The same architecture and training process are used for both datasets except for the last fully connected layer due to the different classes in two datasets. The experiments validate the effectiveness of our method.

### A. Implementation details

Six-hierarchy convolutional layers are stacked for sequence feature generation. Following Simonyan et al. [14], we use 3x3 receptive fields with fewer parameters and fine-grained features. Similar with Shi et al. [28], 1x2 pooling is adopted to make better use of the power of RNN, instead of common square one pooling which leads to a 4x on the width. Batch normalization [30] is adopted after the 3th and the 6th convolutional layer, speeding up the training significantly. The overall configurations are listed in the Table I. With the deeper and wider convolutional layers, more features and more complicated abstracts can be extracted.

Images are simply normalized by subtracting 128 and dividing by 256. Data augmentation is adopted to enlarge the dataset to avoid over-fitting. We crop the images horizontally and flip the image horizontally with probability 0.5. Then we wrap the image into a background with the size of the original image filled with 128 which will be subtracted in the normalization.

Loss function is biased to balance the training examples without inclining to certain kinds of languages.

Experiments are conducted on a server with Intel(R) Xeon(R) CPU E5-2609 CPU and Nvidia Geforce GTX 780 GPU. The network is optimised with ADADELTA [31].

As the operating (both forward and backward) time of an epoch grows linearly with the length of the longest image in the batch, direct training with varied length images is not effective and hard to converge as we observed in our experiments. To deal with this problem, we first rescale all gray-scale text images to the fix size 32x128 and then feed them into our network. The network converges after 20000 iterations with batch size 32.

TABLE I  
OVERALL CONFIGURATION OF OUR MODEL

Layer Name	Configurations
Input	32×Width gray-scale image
Convolution MaxPooling	out-maps:64, kernel:3 × 3, padding:1, stride:1 × 1 Window: 2 × 2, stride:2 × 2
Convolution MaxPooling	out-maps:128, kernel:3 × 3, padding:1, stride:1 × 1 Window: 2 × 2, stride:2 × 2
Convolution BatchNorm	out-maps:128, kernel:3 × 3, padding:1, stride:1 × 1 -
Convolution MaxPooling	out-maps:256, kernel:3 × 3, padding:1, stride:1 × 1 Window: 1 × 2, stride:1 × 2
Convolution MaxPooling	out-maps:256, kernel:3 × 3, padding:1, stride:1 × 1 Window: 2 × 2, stride:2 × 2
Convolution BatchNorm MaxPooling Split	out-maps:512, kernel:3 × 3, padding:1, stride:1 × 1 - Window: 1 × 2, stride:1 × 2 -
Bi-LSTM	hidden units: 256
Bi-LSTM	hidden units: 256
Ave-Pooling SoftMax	- -

Then we feed the network with images scaled to the same height 32 while keeping the aspect ratio. To be clear, we do not feed varied length one-by-one into the model. This will consume a 16x of time, making the total training time too long. In contrast, images with similar aspect ratios are integrated in one batch for efficiency. Some images are pretty long, we have to reduce the batch size to 16 to fit the GPU's memory, and another 40000 iterations are performed. By testing one by one, each image takes approximately 92ms on average on GTX 780 which is relatively longer than previous methods. But with the great improvement, this sacrifice is worthy.

### B. Experiments on SIW-13

In this section, we evaluate our method on SIW-13 dataset.

SIW-13 is a scene text dataset which consists of 13 types scripts. It contains totally 16291 images extracted from natural images. 6500 of them are for testing, 500 for each language, and the other are for training.

As a comparison, we also show the results of the previous state-of-art Shi et al. [28] and naive CNN (see below for details) on it. We attain an improvement about 3.3% on average compared with Shi et al. [28]. See Table II for final results.

As we can see in Table II, the performances of majority of these thirteen languages improve compared with Shi et al. [28].

Especially, languages which share a subset of letters such as English, Russian and Greek still have a lot of room for improvement in Shi et al. [28], despite the use of discriminative clustering. Compared with results in Shi et al. [28], results of these languages make great progress by a large margin by explicitly modeling the spatial dependencies within text lines, proving the significance of it.

TABLE II  
ACCURACIES ON SIW-13 FOR ALL SCRIPT TYPES

	naive-CNN (see below)	Shi et al. [28]	our method
<b>Ara</b>	92.8	94	<b>96.2</b>
<b>Cam</b>	83.6	88	<b>93.4</b>
<b>Chi</b>	90.6	88	<b>94.0</b>
<b>Eng</b>	70.2	71	<b>83.6</b>
<b>Gre</b>	77.4	81	<b>89.4</b>
<b>Heb</b>	91.0	91	<b>93.8</b>
<b>Jap</b>	89.8	90	<b>91.8</b>
<b>Kan</b>	86.2	91	<b>91.8</b>
<b>Kor</b>	93.8	95	<b>95.6</b>
<b>Mon</b>	96.6	96	<b>97.0</b>
<b>Rus</b>	79.2	79	<b>87.0</b>
<b>Tha</b>	91.8	<b>94</b>	93.6
<b>Tib</b>	97.6	97	<b>98.6</b>
<b>Avg.</b>	87.7	89.4	<b>92.75</b>

TABLE III  
COMPARASION FOR CVSI2015, TASK 4

Method	CVSI-2015
<b>This work - CNN with stacked LSTM</b>	<b>94.2%</b>
Google [29]	<b>98.91%</b>
C-DAC [29]	84.66%
HUST [29], [32]	96.69%
CVC-1 [29]	95.88%
CVC-2 [29]	96.00%
CUK [29]	74.06%

### C. Experiments on CVSI2015

We also evaluate our method on a common script identification dataset, namely CVSI2015. CVSI2015 contains pre-segmented texts of ten types mainly extracted from various video sources while having few scene texts. 60% of texts are for training. 10% are for validation. The rest 30% are for testing. In our experiment, we do not use validation part.

CVSI2015 is relatively more idealised, with limited variation compared with SIW13, which limiting the power of our method, as our method can capture long spatial dependencies which is good at dealing with long text image. Despite this limit, our method still achieve a competitive result on it.

Table III compares our method with competitive methods on CVSI2015 for task 4. Refer Figure 4 for class level accuracy. As we can see, Google achieves the highest accuracy while our method also shows strong competence for the task. To be clear, Google uses binarization as image pre-processing method, which is available only for text with great background, limiting its ability for scene text identification. However, our method does not suffer from this problem. HUST [29] also achieves a fancy accuracy because the application of multiple features. In contrast, our model only applies one feature. We believe that with the help of more features, our model will achieve better performance.

### D. Improvement analysis

To understand each part's function in the network, we perform ablation experiment on SIW-13.

A strong baseline with the same convolutional layers and batch normalization is set up for comparison. The only dif-

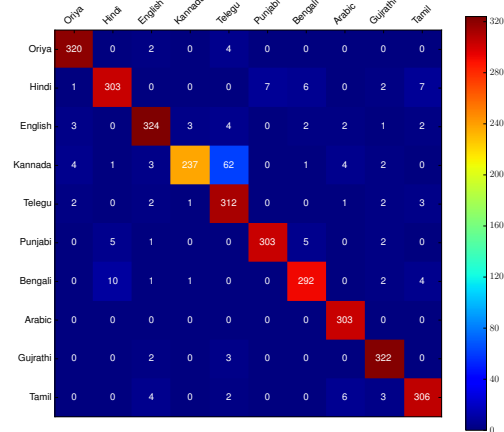


Fig. 4. Confusion matrix for CVSI2015, best view in color.

ference between the baseline and our proposed architecture is that the LSTM is replaced by a fully connected layer with dropout setting to 0.5.

This baseline architecture is holistic, so it suffers from the obstacles described above. It can only reach 87.74% accuracy. After adding a bidirectional LSTM to it, training with 32x128 images, the overall performance reaches 89.74% due to the application of spatial dependencies within text lines. As we expected, fine-tuning with varied length images attains another large improvement about 1.8% by avoiding deteriorate the key patch of the image which reveals the power of great image description. Finally, with the help of the stacked bidirectional LSTM for more powerful abstraction ability, we obtain an accuracy of 92.75%. The final results are listed in Table IV for comparison.

Our experiment reveals the great power of convolutional recurrent neural network against conventional holistic CNN, with a relevant improvement of 4.5%.

TABLE IV  
FINAL RESULT OF 4 MODELS

Model	Accuracy	Model size
Naive CNN	87.74%	6.5M
CNN with one LSTM (fix length)	89.74%	4M
CNN with one LSTM (vary length)	91.54%	4M
CNN with stacked LSTM	92.75%	5.2M

To be clear, the last two architectures have less parameters than the first one while beating it by a large margin, which reveals the power of the proposed method.

## IV. CONCLUSION

We have presented a novel approach for scene text script identification. The new approach exploits the image representation and spatial dependencies within text lines. The algorithm combines CNN and RNN into one end-to-end trainable network, while CNN generates rich image representations and

RNN effectively analyzes long-term spatial dependencies. The experiment on SIW-13 and CVSI2015 datasets demonstrates the power of the new algorithm.

This work explores the power of the image representation and spatial dependencies. Other successful strategies, for instance attention model [22], can be combined with these two factors. Using attention model over simple averaging is supposed to improve the result. This will be investigated in the future.

## V. ACKNOWLEDGE

This research was supported by National Natural Science Foundation of China (NSFC) (No. 61222308, No. 61573160) and Open Project Program of the National Key Laboratory of Digital Publishing Technology (No. F2016001).

## REFERENCES

- [1] A. L. Spitz, "Determination of the script and language content of document images," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 3, pp. 235–245, 1997.
- [2] U. Pal, S. Sinha, and B. B. Chaudhuri, "Multi-script line identification from indian document," in *ICDAR*, 2003, pp. 880–884.
- [3] J. Hochberg, P. Kelly, T. Thomas, and L. Kerns, "Automatic script identification from document images using cluster-based templates," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 2, pp. 176–181, 1997.
- [4] T. N. Tan, "Rotation invariant texture features and their use in automatic script identification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 7, pp. 751–756, 1998.
- [5] T. Q. Phan, P. Shivakumara, Z. Ding, S. Lu, and C. L. Tan, "Video script identification based on text lines," in *ICDAR*, 2011, pp. 1240–1244.
- [6] D. Zhao, P. Shivakumara, S. Lu, and C. L. Tan, "New spatial-gradient-features for video script identification," in *DAS workshop*, 2012, pp. 38–42.
- [7] L. G. i Bigorda, A. Nicolaou, and D. Karatzas, "Boosting patch-based scene text script identification with ensembles of conjoined networks," *CoRR*, vol. abs/1602.07480, 2016.
- [8] A. K. Singh, A. Mishra, P. Dabral, and C. V. Jawahar, "A simple and effective solution for script identification in the wild," in *DAS workshop*, 2016, pp. 428–433.
- [9] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [10] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *CVPR*, 2005, pp. 886–893.
- [11] Y. Bengio, "Learning deep architectures for AI," *Foundations and Trends in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient Based Learning Applied to Document Recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [13] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015, pp. 1–9.
- [14] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015.
- [16] R. B. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014, pp. 580–587.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [18] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015, pp. 3431–3440.
- [19] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," *CoRR*, vol. abs/1412.7062, 2014.
- [20] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *CoRR*, vol. abs/1409.0473, 2014.
- [21] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *EMNLP*, 2014, pp. 1724–1734.
- [22] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," *CoRR*, vol. abs/1506.07503, 2015.
- [23] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *CVPR*, 2015, pp. 3156–3164.
- [24] B. Shi, X. Bai, and C. Yao, "Script identification in the wild via discriminative convolutional neural network," *Pattern Recognition*, vol. 52, pp. 448–458, 2016.
- [25] S. Singh, A. Gupta, and A. A. Efros, "Unsupervised discovery of mid-level discriminative patches," in *ECCV*, 2012, pp. 73–86.
- [26] Y. Bengio, P. Y. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Trans. Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [27] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [28] B. Shi, X. Bai, and C. Yao, "An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition," *CoRR*, vol. abs/1507.05717, 2015.
- [29] N. Sharma, R. Mandal, R. Sharma, U. Pal, and M. Blumenstein, "ICDAR2015 competition on video script identification (CVSI 2015)," in *ICDAR*, 2015, pp. 1196–1200.
- [30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015, pp. 448–456.
- [31] M. D. Zeiler, "ADADELTA: an adaptive learning rate method," *CoRR*, vol. abs/1212.5701, 2012.
- [32] B. Shi, C. Yao, C. Zhang, X. Guo, F. Huang, and X. Bai, "Automatic script identification in the wild," in *ICDAR*, 2015, pp. 531–535.