Deep Image Compression Using Scene Text Quality Assessment

Shohei Uchigasaki^a, Tomo Miyazaki^{a,*}, Shinichiro Omachi^a

^aGraduate School of Engineering, Tohoku University, 6-6-05, Aoba Aramaki, Aoba, Sendai, 980-8579, Miyagi, Japan

Abstract

Image compression is a fundamental technology for Internet communication engineering. However, a high compression rate with general methods may degrade images, resulting in unreadable texts. In this paper, we propose an image compression method for maintaining text quality. We developed a scene text image quality assessment model to assess text quality in compressed images. The assessment model iteratively searches for the best-compressed image holding high-quality text. Objective and subjective results showed that the proposed method was superior to existing methods. Furthermore, the proposed assessment model outperformed other deep-learning regression models.

Keywords: image compression, scene text image, text quality, quality assessment, regression model

^{*}Corresponding author

Email addresses: uchiga@iic.ecei.tohoku.ac.jp (Shohei Uchigasaki), tomo@tohoku.ac.jp (Tomo Miyazaki), machi@ecei.tohoku.ac.jp (Shinichiro Omachi)



(a) Original Image

(b) Image reconstructed using JPEG

Figure 1: Results of high-compression image compression.

1. Introduction

The demand for higher image compression has increased due to the growth of communication traffic in recent years. Image compression is an essential technique for transmitting important content. One of the important contents in images is text, such as on signboards and book covers. However, traditional image compression methods such as JPEG [1] may generate block noise when the compression ratio is high, resulting in unreadable text. Generally, JPEG images at high compression lost the readability of the text significantly, as shown in Figure 1.

An image compression technique that can retain text quality even under high compression is required. Nevertheless, image compression focusing on text quality has considerable room for development. A fundamental technique to evaluate a compressed image's quality is assessing text quality assessment. It would be inefficient for a human to evaluate the readability of all compressed images. By contrast, only a human can assess the readability of text. Two types of images contain text: scene images, such as landscapes and buildings, and document images, such as books and documents. Models for assessing text quality in scene images are yet to be studied.

In this study, we propose a model for assessing text quality in scene images using deep learning. Using only the image features obtained from convolution networks makes text readability assessment difficult. Therefore, we use a text recognition model to obtain informative features for the assessment. A text recognition model is similar to the human eye in text understanding. The proposed method assesses text quality using the image features and the text probabilities predicted by a text recognition model. The text quality assessment model allows the evaluation of the text quality (i.e., readability) of a scene image without human intervention. In addition, we propose an image compression method to ensure text quality. We use a deep neural network to develop the proposed image compression method. The proposed method controls text quality by changing the image compression parameters based on the text quality assessed from the image compression results.

The main contributions of this paper are the following.

- We propose a quality assessment model for scene text images using image features and the text probabilities predicted by a text recognition model.
- We develop a deep image compression method that ensures text quality by applying the proposed assessment model.
- The proposed assessment and image compression models outperform traditional methods in objective measures.

The experimental results show that the proposed text quality assessment

model outperforms existing deep learning feature extractors. Furthermore, the image compression method that maintains text quality outperforms traditional image compression methods.

2. Related work

The proposed method has two modules: text quality assessment and image compression. Thus, we describe related works in the aspect of the two modules.

2.1. Text quality assessment

Image quality assessment is a fundamental technique in image processing and computer vision. There are representative works with excellent performance. A referenceless quality metric was developed for synthesized images produced by a depth image-based rendering [2]. The literature [3] discovered that saliency preservation effectively enhances automatic contrast. Gu et al. analyzed the parameters of an autoregressive model to develop a no-reference sharpness metric [4]. Also, an assessment model uses image features of picture complexity, content statics, global brightness quality, and sharpness [5].

Broadly, we can divide text images into two types: document images and scene text images. Various methods assess text quality in document images. For example, machine learning methods assess text quality from features extracted by filtering and edge detection [6]. In addition, deep learning algorithms, such as a convolutional neural network (CNN), have been used for text quality assessment. The first deep learning-based method divides a document image into small patches and assesses them using a CNN [7]. Variant methods extract patches from different images [8] and various scales [9]. Jemn et al. enhance text images using a generative adversarial network to improve character recognition performance [10].

Text quality assessment models for scene images have not yet been studied. Super-resolution (SR) methods for scene text images relate to text quality assessment since SR improves text quality. TextSR [11] uses a text recognition model as a loss function for SR. Likewise, TBSRN [12] uses a feature map extracted from the middle layer of a text recognition model. TextGestalt [13] calculates loss values using text strokes. TPGSR [14] includes a text recognition model in its network. TATT [15] combines a text recognition model and transformer. Inspired by [15], we use the transformer network to develop a quality assessment model for scene text images.

2.2. Image compression

Image compression methods include traditional methods, such as JPEG [1], JPEG2000 [16], and deep learning methods.

Firstly, we introduce conventional methods. JPEG and JPEG2000 are the most widespread image compression methods. Another method is sparse coding [17]. However, the conventional methods cannot use region-of-interest (RoI) compression and suffer from block noise.

Deep image compression methods can be divided into two types: variablerate image compression [18, 19] and ROI compression [20, 21]. The RoI methods use a binary mask to maintain high quality for specified regions while ignoring the quality of other regions. Song et al. [22] develop a method enabling variable-rate and RoI compression. The method [22] uses a quality



Figure 2: Proposed framework

map instead of a binary mask. The quality map can compress images in various bit rates and allows fine values of quality in ROI. In this study, we employed the deep image compression method [22] to control text quality.

3. Proposed method

We propose an image compression method for scene images. As shown in Figure 2, the proposed method compresses images three times by updating the quality map. Then, we produce the final compressed image by selecting the compressed image that achieves the best text quality. The main idea is to compress images by assessing text quality with the proposed scene text image quality assessment (STIQA) model. The proposed method combines variable-rate deep image compression [22] and the STIQA model. The proposed method defines the compression quality in each pixel. Thus, we can define a suitable quality map for scene text images. We update a quality map by assessing text quality using the STIQA model. This strategy provides a suitable text quality and compression for each text region.



Figure 3: Overall architecture of STIQA

3.1. Scene text image quality assessment

3.1.1. Overall framework of the STIQA model

An overall architecture of the STIQA model is shown in Figure 3. The STIQA model crops a text region from an image, resizing it to 32×128 . We extract features from the text region using two networks. Firstly, we predict a sequence of probabilities for character categories using a convolutional recurrent neural network (CRNN) [23]. The character categories include a–z, 0–9, and white space. Thus, the total number of character categories is 37. The second network extracts image features using a convolution layer with a window size of 7×7 . Then, we fuse the sequence and the image features using the transformer [24] to produce features. The final layer comprises adaptive average pooling, fully connected layers, and a sigmoid activation function. The output is a predicted text quality ranging from [0, 1]. Higher value is high quality.

3.1.2. Text image resizing

The text regions are cropped from a scene image and are resized to 32×128 . A standard rescale method, such as bilinear interpolation, distorts the

aspect ratio of the original text. Therefore, we resize the text by expanding backgrounds if the text size is less than 32×128 . We use the second-lowest pixel value among the four corner pixels, and Gaussian noise is then added to the pixel value. On the other hand, if the image size exceeds 32×128 , we randomly crop the text to 32×128 . As a result, the original text image can be input into the STIQA model.

3.1.3. Transformer

We adopted the transformer [24] to predict text quality. Figure 4 illustrates the architecture. The Transformer can find the global relationships between two different features and output them as features. Convolution cannot deal with scene text images of various shapes and angles owing to its narrow effective range. Therefore, the Transformer with a global correlation is appropriate. The Transformer has an encoder and decoder. The positional encoding processes, called FPE and RPE, are carried out separately since the input does not consider the position or semantic order. Therefore, positional information must be added in advance.

The encoder extracts correlations and features from the text strings. In addition, the encoder reshapes the size of the features. The encoder embeds the sequence of character probabilities and the image features into the same embedding space. Specifically, the self-attention operation is used to extract correlations between the probabilities. The feed-forward network performs feature refinement.

The decoder transforms the image features into those that contain text information. Cross-attention is used to extract mutual features for the encoder output and image features. In other words, the decoder associates the text



Figure 4: Detailed architecture of the Transformer in the STIQA

information output from the encoder with the corresponding positions in the image features. With the decoder, image features include text information, which strongly influences the text quality assessment.

3.1.4. Loss function

During training, the loss function uses the error between the ground truth and predicted labels. We use the L1 norm loss and the ϵ -insensitive loss function. The ϵ -insensitive loss function does not give a penalty if the error is within the allowed ϵ and gives a penalty if the error exceeds ϵ . The formula for the ϵ -insensitive loss function L_{ϵ} is defined in Eq. (1). *GT* represents the ground truth label and *Pred* the predicted value. In this experiment, $\epsilon = 0.1$. The reason for introducing the ϵ -insensitive loss function is to penalize the learning process further if the estimation of text quality is largely incorrect. The overall loss function L_{all} is defined by Eq. (2).



Figure 5: The architecture of deep image compression model using quality maps

$$L_{\epsilon} = \begin{cases} 0, & |GT - Pred| < \epsilon \\ |GT - Pred| - \epsilon, & |GT - Pred| \ge \epsilon \end{cases}$$
(1)

$$L_{all} = L_1 + L_{\epsilon} = |GT - Pred| + L_{\epsilon} \tag{2}$$

3.2. Image compression methods considering text quality

3.2.1. Image compression using quality maps

The proposed image compression method uses variable-rate deep image compression [22]. The network is illustrated in Figure 5. We can perform variable-rate image compression by using a pixel-by-pixel quality map.

The SFT layer [25] is used in the condition network, encoder, and decoder. The others consist of convolution and normalization layers. The condition network fuses the information from a quality map using image features. Pixel-by-pixel quality maps enable high image quality in various areas, such as text and significant areas. Therefore, it can be used for a wide variety of tasks. The quality map has values ranging from [0, 1] to allow fine quality control. Henceforth, the value assigned to the quality map is denoted by the



Figure 6: Quality map update method

weight k.

A previous work [22] assigns the same weight value k to all text regions in the image. However, the size and shape of the text vary significantly. Thus, the texts are different in readability. Therefore, we assigned different weights to each text region in this study.

3.2.2. Quality map update

We propose an image compression method using a quality map to preserve text quality assessed by STIQA. The flow of the proposed image compression method is illustrated in Figure 6. This quality map is weighted differently for each text region.

The initial quality map has weights of 0.5 on pixels of character edges and inside. Otherwise, the weight is 0.2 for the background. High weights are assigned to the edges because character edges are the most critical factors for readability. We use the Canny method to extract edges.

We compress images and assess the quality of each text in the compressed image. We update weights k_{new} by the distance between the assessment results and the target text quality. The update process is defined in Eq. (3). k_{old} denotes the value of the previous weights, λ is the hyperparameter that increases the update amount. $Score_{target}$ represents the target text quality. Score is the assessed text quality. We use $\lambda = 5$ and $Score_{target} = 0.90$ in this work.

$$k_{new} = k_{old} + \lambda(Score_{target} - Score) \tag{3}$$

We iteratively perform image compression using the updated quality map. In this paper, the repetition is three times. Consequently, we can obtain a quality map that achieves high compression while maintaining text quality.

4. Experiments

We conducted experiments to evaluate the proposed STIQA and the image compression method. Firstly, we described the datasets and training details. Then we showed the experimental results of the proposed STIQA and the image-compression methods.

4.1. Dataset

The experiments used two datasets: ICDAR2013 [26] and TextOCR [27]. The ICDAR2013 dataset is the focused scene text dataset of the Robust Reading Competition in ICDAR 2013, which is widely used in scene text analysis. TextOCR is a benchmark dataset for text recognition of arbitrarily shaped scene text. ICDAR2013 and TextOCR contain 229 and 21,778 training images, respectively.

We determined the quality of the text regions in 163 test images from the ICDAR2013 dataset by human annotation. In the human annotation task, 15 participants assessed the readability of the text in five grades. The answers were averaged and normalized to [0, 1]. We randomly selected 163 images from the test images of ICDAR2013 to alleviate annotation work.

We used a text recognition model, ABINet [28], to determine the quality of the text regions in the training images of ICDAR2013 and TextOCR since human annotation is infeasible for a large number of training images. We assumed that the human and the model are relevant if the model is accurate. Specifically, text quality q is determined by Eq. (4). We used two measurements: the prediction confidence c and accuracy a. ABINet predicted a set of probabilities $\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_n)$ for $n = |\mathbf{p}|$ characters in a text region. $\mathbf{p}_i = (p_1, \dots, p_{37})$ represents the output of a softmax function for the 37 classes. In addition, we measured the Levenshtein distance d_{lev} between the text of the ground truth t_g and the prediction t_p .

$$q = \frac{1}{2}(c+a),\tag{4}$$

$$c = \frac{1}{|\boldsymbol{P}|} \sum_{\boldsymbol{p} \in \boldsymbol{P}} \max(\boldsymbol{p}), \tag{5}$$

$$a = 1 - \frac{d_{\text{lev}}(t_g, t_p)}{n}.$$
(6)

4.2. Training details

The STIQA model was trained using text regions cropped from highly compressed JPEG images of ICDAR2013 training dataset. There were 844 text regions in 229 scene images. The text regions were split into training and validation sets at an 80:20 ratio. The scene images were converted into grayscale and resized to 192×128 pixels. The purpose of resizing the images was to verify the performance of the quality assessment, even for small text regions. We used a batch size of 16 and 500 epochs, a learning rate of 10^{-3} , and the Adam optimizer. The training images of the TextOCR dataset were used to train the proposed image compression method. Briefly, the same conditions as described in [22] were followed. The images of the TextOCR dataset were randomly cropped to 256×256 pixels.

4.3. Experiments on image compression

We compressed the test images and evaluated the quality of text regions. The proposed image compression method is compared with traditional image compression methods. We conducted a comparative experiment using objective and subjective evaluation metrics.

For comparison, we used the existing JPEG, JPEG2000, and a naive method using sparse coding [17]. We developed the naive method using optional weights in sparse coding. The naive method varied the weights and applied them to the text regions according to the text quality obtained by the STIQA. Therefore, the comparisons ensured fairness between the naive and the proposed methods. Furthermore, we used Qmap [22] to compress the images. All methods compressed the images at the same bit rate.

The evaluation metrics are PSNR, SSIM [29], and LPIPS [30]. In addition, the performance of a text detection model, TextFuseNet [31], was used to evaluate the compressed scene text images. We use precision, recall, and f-measure. The evaluation metrics were calculated for compressed images at almost identical bit rates.

4.3.1. Main results on image compression

The image compression results were presented in Table 1. The results were averaged over the test images in ICDAR2013. The results showed that

	Entire image quality			Text quality			Text detection		
	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	Р	R	F
JPEG	23.4	0.65	0.29	19.6	0.56	0.27	0.28	0.19	0.22
JPEG2000	26.9	0.77	0.25	22.5	0.72	0.18	0.60	0.51	0.54
Naive	25.6	0.71	0.23	<u>27.0</u>	<u>0.89</u>	<u>0.07</u>	<u>0.75</u>	<u>0.65</u>	<u>0.68</u>
$\operatorname{Qmap}[22]$	31.2	0.88	0.13	26.2	0.85	0.11	0.68	0.61	0.63
Proposed	<u>29.7</u>	<u>0.85</u>	0.17	28.7	0.93	0.05	0.79	0.74	0.76

 Table 1: Results on image compression (Bold and underline represent the best and the second best, respectively.)

the proposed method was superior for all the metrics. The proposed method obtained the second-best results for the entire image. The proposed image compression method outperformed traditional image compression methods, such as JPEG. Furthermore, the proposed method is the best in terms of text quality and text detection metrics. Therefore, the results verified the effectiveness of the proposed image compression method.

Qmap outperformed ours in overall image quality because Qmap was optimized for overall images, whereas our method focuses on text readability. Image compression is to determine a bitstream for a set of pixel values, as shown in Figure 5. Qmap determines a bitstream using uniform weights for all pixels, i.e., the overall image. Thus, the overall quality is favorable. However, text quality greatly deteriorates since the number of text pixels is smaller than the other pixels. In contrast, our method modifies weights by assigning more bits to text, resulting in a bitstream maintaining text readability. As shown in Figure 7, although overall image quality is reduced



Figure 7: Compressed images

slightly, text quality improved greatly.

Some of the compressed images are shown in Figure 7. JPEG and JPEG2000 were significantly degraded in both background and text quality. In addition, the naive method degrades only the background. By contrast, the proposed method achieves a higher image quality for text regions and a higher image quality for background regions compared to the other methods.

4.3.2. Subjective evaluation results

Subjective evaluation experiments were conducted using the proposed and naive methods. Fifteen participants chose better readability between two compressed images in terms of the entire image and the text region. Each participant had 50 pairs.

The results are presented in Table 2. The proposed method is superior

 Table 2: The ratios of Subjective evaluation on comparison of compressed images

 by the Naive and the proposed methods

Answer	Entire image	Text region
Same quality	0.3(%)	28.4 (%)
Naive is better	2.4 (%)	5.3~(%)
Proposed is better	97.3 (%)	66.3 (%)

for the entire image and text regions. The proposed method obtained more than 66% of the votes for the text region. These results demonstrate the capability of this method in maintaining text readability.

4.3.3. Hyper Parameter determination in image compression

The effect of varying the hyperparameter λ and the number of updates was also investigated. λ is defined in Eq. (3) to determine the amount of quality map update. Table 3 shows the bit rates (bpp) and the average score of the text qualities assessed using STIQA. The bit rates and STIQA scores increased according to λ and the number of updates. The best STIQA score was obtained when $\lambda = 5$. By considering the computing time, it was determined that three update cycles were sufficient.

4.4. Experiments on text quality assessment

We evaluated the STIQA model and other methods. The comparison methods used in this experiment were the ResNet-34 [32], VGG-13 [33], EfficientNet-B3 [34], and Document Image Quality Assessment (DIQA) [8]. We trained models to regress the text quality assessment.

	3 updates		10 updates		
λ	bpp	STIQA	bpp	STIQA	
10^{-1}	0.18	0.57	0.18	0.59	
10^{0}	0.19	0.59	0.22	0.62	
5×10^0	0.22	0.62	0.22	0.62	
10^{1}	0.23	0.62	0.23	0.62	
10^{2}	0.23	0.61	0.23	0.61	

Table 3: Bit rates and STIQA scores at various hyperparameters

4.4.1. Text quality assessment results

The results are presented in Table4. Three evaluation metrics were used: mean absolute error, Spearman's rank correlation coefficient, and Pearson's correlation coefficient. We used the indicators to show a relationship between humans and our model in text quality assessment. Specifically, MAE shows the gap between a model and humans. Pearson's correlation coefficient uses scores of text quality to evaluate a linear relationship. Spearman's rank correlation coefficient uses ranks among data to verify a monotonic relationship. Thus, we can comprehensively evaluate text quality assessment models using the three indicators.

The results showed that STIQA performed the best on the test dataset with text quality labels obtained from subjective evaluation experiments. The main difference between STIQA and other methods is the existence of a text recognition model, which may play a significant role in text quality assessment.

 Table 4: Results on text quality assessment

Model	MAE	Spearman	Pearson
EfficientNet [34]	0.31	0.27	0.24
VGG-13 [33]	0.30	0.25	0.28
Res- 34 [32]	0.28	0.37	0.34
DIQA [8]	0.22	0.43	0.47
STIQA (Proposed)	0.20	0.55	0.59

4.5. Ablation study on STIQA

We conducted an ablation study on the proposed Scene Text Image Quality Assessment model (STIQA). The STIQA consists of three modules: CRNN, Conv layer, and Transformer as shown in Figure 3. The CRNN extracts text features by predicting probabilities for character categories. Also, the Conv layer extracts image features. The Transformer fuses the text and image features to produce features. The final layer is a fully connected layer to predict text quality. The results are shown in Table 5. The CRNN model obtained relatively good performance. In contrast, the other models were superior. The full STIQA achieved the best performance. Comparing CRNN + Tans and the full STIQA, the Conv layer significantly improved performance. Therefore, the image features are essential to the proposed text quality assessment model.

4.5.1. Parameter combination results

We evaluated the impact of the combinations of the loss functions: the L1 loss and the ϵ -insensitive loss. The results are presented in Table 6. The L_{ϵ}

 Table 5: Ablation results of STIQA in text quality assessment

Model	MAE	Spearman	Pearson
CRNN	0.218	0.40	0.46
CRNN + Transformer	0.203	0.47	0.55
CRNN + Transformer + Conv (Full STIQA)	0.197	0.55	0.59

Table 6: Results on parameter combination in the proposed STIQA

L_1	L_{ϵ}	ϵ	MAE	Spearman	Pearson
\checkmark	-	NA	0.212	0.50	0.57
-	\checkmark	0.10	0.210	0.44	0.50
\checkmark	\checkmark	0.20	0.214	0.52	0.55
\checkmark	\checkmark	0.15	0.206	0.48	0.56
		0.10	0.197	0.55	0.59

loss was better than the L_1 loss. In addition, the best results were achieved by using the two loss functions simultaneously, and the combination with an ϵ of 0.10 was the best overall.

5. Conclusions

We proposed a method for text quality-preserving image compression using a text quality assessment model. The proposed image compression method can produce readable scene text images by assessing the quality of the text regions in the images. Text quality was assessed using a transformer to consider global and local features extracted from the entire image and the text regions. The experimental results showed that the proposed method is superior to the JPEG and naive methods using sparse coding compression. In addition, the proposed text quality assessment model outperformed other regression models. Both objective and subjective evaluations showed that the proposed method is capable of text quality preservation during image compression.

A promising application of the proposed method is video compression. We can efficiently preserve text information in a video conference, IoT, and sensor network. Besides, the proposed method can be applied to more applications if we extend the text quality assessment model to general objects.

Future work is to design an end-to-end model. The proposed model has two modules for image compression and text quality assessment. We trained the modules individually. Thus, we can improve the performance of the proposed method through end-to-end learning. The challenge is how to connect the two modules. A solution is text quality assessment using features encoded by the image compression model. This solution can facilitate learning the model and is beneficial to computational costs.

Constructing a practical system in combination with the text detection method is also essential for future work. Since it is difficult to evaluate the effectiveness of the proposed method when text detection errors occur, we conducted the experiments assuming that the text regions are known. The proposed STIQA can evaluate text readability. Thus, we believe it is possible to improve the accuracy of text detection by evaluating the readability of the regions detected by text detection methods.

Acknowledgments

This work was partially supported by JSPS KAKENHI under Grants JP20H04201, JP22K12729, and JP22H00540.

References

- W. B. Pennebaker, J. L. Mitchell, JPEG Still Image Data Compression Standard, 1st Edition, Kluwer Academic Publishers, USA, 1992.
- [2] K. Gu, V. Jakhetiya, J.-F. Qiao, X. Li, W. Lin, D. Thalmann, Modelbased referenceless quality metric of 3d synthesized images using local image description, IEEE Transactions on Image Processing 27 (1) (2018) 394–405. doi:10.1109/TIP.2017.2733164.
- K. Gu, G. Zhai, X. Yang, W. Zhang, C. W. Chen, Automatic contrast enhancement technology with saliency preservation, IEEE Transactions on Circuits and Systems for Video Technology 25 (9) (2015) 1480–1494. doi:10.1109/TCSVT.2014.2372392.
- [4] K. Gu, G. Zhai, W. Lin, X. Yang, W. Zhang, No-reference image sharpness assessment in autoregressive parameter space, IEEE Transactions on Image Processing 24 (10) (2015) 3218–3231. doi:10.1109/TIP. 2015.2439035.
- K. Gu, J. Zhou, J.-F. Qiao, G. Zhai, W. Lin, A. C. Bovik, No-reference quality assessment of screen content pictures, IEEE Transactions on Image Processing 26 (8) (2017) 4005–4018. doi:10.1109/TIP.2017. 2711279.

- [6] T. Obafemi-Ajayi, G. Agam, Character-based automated human perception quality assessment in document images, IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans 42 (3) (2012) 584–595. doi:10.1109/TSMCA.2011.2170417.
- [7] L. Kang, P. Ye, Y. Li, D. Doermann, A deep learning approach to document image quality assessment, in: 2014 IEEE International Conference on Image Processing (ICIP), 2014, pp. 2570–2574. doi:10.1109/ICIP. 2014.7025520.
- [8] X. Peng, C. Wang, Camera captured diqa with linearity and monotonicity constraints, in: X. Bai, D. Karatzas, D. Lopresti (Eds.), Document Analysis Systems, Springer International Publishing, Cham, 2020, pp. 168–181.
- [9] D. Rodin, V. Loginov, I. Zagaynov, N. Orlov, Document image quality assessment via explicit blur and text size estimation, in: J. Lladós, D. Lopresti, S. Uchida (Eds.), Document Analysis and Recognition – ICDAR 2021, Springer International Publishing, Cham, 2021, pp. 281– 292.
- S. Khamekhem Jemni, M. A. Souibgui, Y. Kessentini, A. Fornés, Enhance to read better: A multi-task adversarial network for handwritten document image enhancement, Pattern Recognition 123 (2022) 108370.
 doi:https://doi.org/10.1016/j.patcog.2021.108370.
 URL https://www.sciencedirect.com/science/article/pii/S0031320321005501

- [11] W. Wang, E. Xie, X. Liu, W. Wang, D. Liang, C. Shen, X. Bai, Scene text image super-resolution in the wild, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, Springer International Publishing, Cham, 2020, pp. 650–666.
- [12] J. Chen, B. Li, X. Xue, Scene text telescope: Text-focused scene image super-resolution, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 12021–12030. doi:10. 1109/CVPR46437.2021.01185.
- [13] J. Chen, H. Yu, J. Ma, B. Li, X. Xue, Text gestalt: Stroke-aware scene text image super-resolution, CoRR abs/2112.08171 (2021). arXiv: 2112.08171.
 URL https://arxiv.org/abs/2112.08171
- [14] J. Ma, S. Guo, L. Zhang, Text prior guided scene text image superresolution (2021). doi:10.48550/ARXIV.2106.15368.
- [15] J. Ma, Z. Liang, L. Zhang, A text attention network for spatial deformation robust scene text image super-resolution, in: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 5901–5910. doi:10.1109/CVPR52688.2022.00582.
- [16] A. Skodras, C. Christopoulos, T. Ebrahimi, The jpeg 2000 still image compression standard, IEEE Signal Processing Magazine 18 (5) (2001) 36-58. doi:10.1109/79.952804.
- [17] K. Skretting, K. Engan, Recursive least squares dictionary learning algo-

rithm, IEEE Transactions on Signal Processing 58 (4) (2010) 2121–2130. doi:10.1109/TSP.2010.2040671.

- [18] Z. Cui, J. Wang, S. Gao, T. Guo, Y. Feng, B. Bai, Asymmetric gained deep image compression with continuous rate adaptation, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 10527–10536. doi:10.1109/CVPR46437.2021.01039.
- [19] F. Yang, L. Herranz, J. v. d. Weijer, J. A. I. Guitián, A. M. López, M. G. Mozerov, Variable rate deep image compression with modulated autoencoder, IEEE Signal Processing Letters 27 (2020) 331–335. doi: 10.1109/LSP.2020.2970539.
- [20] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, L. Van Gool, Generative adversarial networks for extreme learned image compression, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 221–231. doi:10.1109/ICCV.2019.00031.
- [21] C. Cai, L. Chen, X. Zhang, Z. Gao, End-to-end optimized roi image compression, IEEE Transactions on Image Processing 29 (2020) 3442– 3457. doi:10.1109/TIP.2019.2960869.
- [22] M. Song, J. Choi, B. Han, Variable-rate deep image compression through spatially-adaptive feature transform, in: 2021 IEEE/CVF International Conference on Computer Vision (ICCV), 2021, pp. 2360–2369. doi: 10.1109/ICCV48922.2021.00238.

- [23] B. Shi, X. Bai, C. Yao, An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (11) (2017) 2298–2304. doi:10.1109/TPAMI.2016.2646371.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17, Curran Associates Inc., Red Hook, NY, USA, 2017, p. 6000–6010.
- [25] X. Wang, K. Yu, C. Dong, C. Change Loy, Recovering realistic texture in image super-resolution by deep spatial feature transform, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 606–615. doi:10.1109/CVPR.2018.00070.
- [26] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i. Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazàn, L. P. de las Heras, Icdar 2013 robust reading competition, in: 2013 12th International Conference on Document Analysis and Recognition, 2013, pp. 1484–1493. doi: 10.1109/ICDAR.2013.221.
- [27] O. Sidorov, R. Hu, M. Rohrbach, A. Singh, Textcaps: A dataset for image captioning with reading comprehension, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision ECCV 2020, Springer International Publishing, Cham, 2020, pp. 742–758.
- [28] S. Fang, H. Xie, Y. Wang, Z. Mao, Y. Zhang, Read like humans:

Autonomous, bidirectional and iterative language modeling for scene text recognition, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 7094–7103. doi: 10.1109/CVPR46437.2021.00702.

- [29] Z. Wang, A. Bovik, H. Sheikh, E. Simoncelli, Image quality assessment: from error visibility to structural similarity, IEEE Transactions on Image Processing 13 (4) (2004) 600–612. doi:10.1109/TIP.2003.819861.
- [30] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, O. Wang, The unreasonable effectiveness of deep features as a perceptual metric, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 586–595. doi:10.1109/CVPR.2018.00068.
- [31] J. Ye, Z. Chen, J. Liu, B. Du, Textfusenet: Scene text detection with richer fused features, in: C. Bessiere (Ed.), Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20, International Joint Conferences on Artificial Intelligence Organization, 2020, pp. 516–522, main track. doi:10.24963/ijcai.2020/72. URL https://doi.org/10.24963/ijcai.2020/72
- [32] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770–778. doi:10.1109/CVPR.2016.90.
- [33] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition (2014). doi:10.48550/ARXIV.1409.1556.
 URL https://arxiv.org/abs/1409.1556

[34] M. Tan, Q. Le, EfficientNet: Rethinking model scaling for convolutional neural networks, in: K. Chaudhuri, R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning, Vol. 97 of Proceedings of Machine Learning Research, PMLR, 2019, pp. 6105– 6114.

URL https://proceedings.mlr.press/v97/tan19a.html