

Visual attention models for scene text recognition

Suman K. Ghosh and Ernest Valveny

*Computer Vision Center, Dept. Ciències de la Computació
Universitat Autònoma de Barcelona
08193 Bellaterra (Barcelona), Spain
Email: sghosh,ernest@cvc.uab.es*

Andrew D. Bagdanov

*Media Integration and Communication Center (MICC)
Universita di Firenze, Firenze
Email: bagdanov@cvc.uab.es.*

Abstract—In this paper we propose an approach to lexicon-free recognition of text in scene images. Our approach relies on a LSTM-based soft visual attention model learned from convolutional features. A set of feature vectors are derived from an intermediate convolutional layer corresponding to different areas of the image. This permits encoding of spatial information into the image representation. In this way, the framework is able to learn how to selectively focus on different parts of the image. At every time step the recognizer emits one character using a weighted combination of the convolutional feature vectors according to the learned attention model. Training can be done end-to-end using only word level annotations. In addition, we show that modifying the beam search algorithm by integrating an explicit language model leads to significantly better recognition results. We validate the performance of our approach on standard SVT, ICDAR’03 and MS-COCO scene text datasets, showing state-of-the-art performance in unconstrained text recognition.

1. Introduction

The increasing ability to capture images in any condition and situation poses many challenges and opportunities for extracting visual information from images. One such challenge is the detection and recognition of text “in the wild”. Text in natural images is a high level semantic information that can aid automatic image understanding and retrieval.

However, robust reading of text in uncontrolled environments is very different from text recognition in document images and much more challenging due to multiple factors such as difficult acquisition conditions, low resolution, font variability, complex backgrounds, different lighting conditions, blur, etc. Therefore, OCR techniques used in document images do not generalize well to recognition of scene text.

The problem of end-to-end scene text recognition is usually divided in two different tasks: word detection and word recognition. The goal of the word detection stage is to generate bounding boxes around potential words in the images. Subsequently, the words in these bounding boxes are recognized in the word recognition stage. This paper is focused on this second stage, word recognition.

Existing word recognition methods can be broadly divided into dictionary-based methods, using some kind of predefined lexicon to guide the recognition, and unconstrained methods, able to recognize any word.

Dictionary-based scene text recognition. Traditionally, scene text recognition systems use character recognizers in a sequential way by localizing characters using a sliding window [9], [12], [17] and then grouping responses by arranging the character windows from left to right as words. A variety of techniques have been used to classify character bounding boxes, including random ferns [17], integer programming [14] and Convolutional Neural Networks (CNNs) [9]. These methods often use the lexical constraints imposed by a fixed lexicon while grouping the character hypotheses into words.

In contrast to sequential character recognizer models, holistic fixed-length representations have been proposed in [?], [1], [4], [8], [9], [13]. In [1], [4], [13], a holistic signature derived from a set of training images is used to learn a joint embedding space between images and words. The first attempt using CNN features was made by Jaderberg *et al.* in [9], where a sliding window over CNN features is used for robust scene text recognition using a fixed lexicon. Later, the same authors also proposed a fixed-length representation [?] using convolutional features trained on a synthetic dataset of 9 million images [8]

Unconstrained scene text recognition. Though most of the works in scene text recognition focus on fixed-lexicon recognition, a few attempts at unconstrained text recognition have also been made.

Biassco *et al.* in [3] rely on sequential character classifiers. They use a massive number of annotated character bounding boxes to learn character classifiers. Binarization and sliding window methods are used to generate character proposals followed by a text/background classifier. Finally, character probabilities given by character classifiers are used in a beam search to recognize words. They also integrate a static character n -gram language model in every step of the beam search to incorporate an underlying language model.

Though CNN models have achieved great success in lexicon-based text recognition, word recognition in unconstrained scenarios requires modeling the underlying

character-level language model. Jaderberg *et al.* in [6] proposed to use two separate CNNs, one modeling character unigram sequences and another n -gram language statistics. They additionally use a Conditional Random Field to model the interdependence of characters (n -grams). However, this significantly increases the computational complexity. In addition, to detect the presence of character n -grams in word images as neural activations, character n -grams are used as output nodes, leading to a huge (10k output units for $n=4$) output layer.

In contrast to the above strategies our approach neither recognizes individual characters in the word image nor uses any holistic representation to recognize the word. It rather uses a LSTM-based visual attention model on top of CNN features (based on [18]) to focus attention on relevant parts of the image at every step and infer a character present in the image (see figure 1). Thus, the system does not require explicit character segmentation and is able to recognize any word, without the help of any predefined dictionary. The visual attention model can be trained using only word bounding boxes and does not need explicit character bounding boxes at training time.

Recently, visual attention models have gained a lot of attention and have been used for machine translation [2], image captioning [18] and also text recognition [10]. In [18] the attention model is combined with an LSTM on top of CNN features. The LSTM outputs one caption word at every step focusing on a specific part of the image driven by the attention model. In our work, we mainly follow this attention model, adapted to the particular case of text recognition. Although the work of [10] also makes use of a soft attention model for text recognition in wild, there are significant differences with respect our work. Firstly, their model relies on Recursive CNN features to model the dependencies between characters. Instead we use traditional, much simpler CNN features and it is the visual attention model which learns to selectively attend to parts of the image and the dependencies between them. Secondly, Lee *et al.* [10] used the features from the fully connected layer, while we use features from an earlier convolutional layer, thus preserving the local spatial characteristics of the image and reducing the model complexity. This also allows the model to focus on a subset of features corresponding to certain area of the image and learn the underlying inter-dependencies. Thirdly, we used LSTM instead of RNN which has been shown to learn long term dependencies better than traditional RNNs.

Our contributions with respect to the state-of-the-art. In summary the contributions of our work are:

- We introduce a LSTM-based visual attention model on top of CNN features for unconstrained scene text recognition. This model is able to selectively attend to specific parts of word images, allowing it to model inter-character dependencies as needed and thus to *implicitly* model the underlying language.
- We show that weak explicit language models (in the form of prefix probabilities) can significantly boost

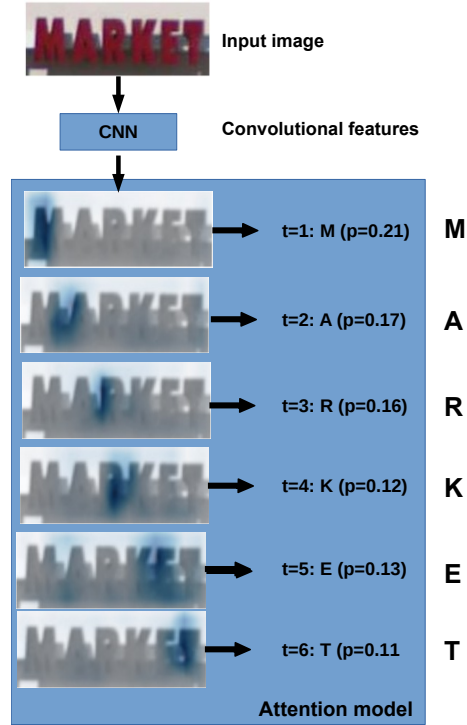


Figure 1. Overall scheme of the proposed recognition framework. Given a cropped word image, a set of spatially localized features are obtained using a CNN. Then, an LSTM decoder is combined with an attention model to generate the sequence of characters. At every time step the attention model weights the set of feature vectors to make the LSTM focus on a specific part of the image.

the final recognition result without having to resort to a fixed lexicon. For that, We modify the beam search to take into account the language model. Additionally, the beam search can also incorporate a lexicon whenever it is available.

- We experimentally validate that our approach with weak language modeling outperforms the state-of-the-art in unconstrained scene text recognition and performs comparably to lexicon-based approaches with a model complexity lower than similar approaches.

The rest of the paper is organized as follows. In Section 2, we present our attention-based recognition approach and 3. In Section 4 we experimentally validate the model on a variety of standard and public benchmark datasets. We conclude in Section 5 with a summary of our contributions and a discussion of future research directions.

2. Visual attention for scene text recognition

Our recognition approach is based on an encoder-decoder framework for sequence to sequence learning. An overall scheme of the framework is illustrated in figure 2. The encoder takes an image of a cropped word as input and encodes this image as a sequence of convolutional features. The attention model in between the encoder and the decoder drives, at every step, the focus of attention of the decoder

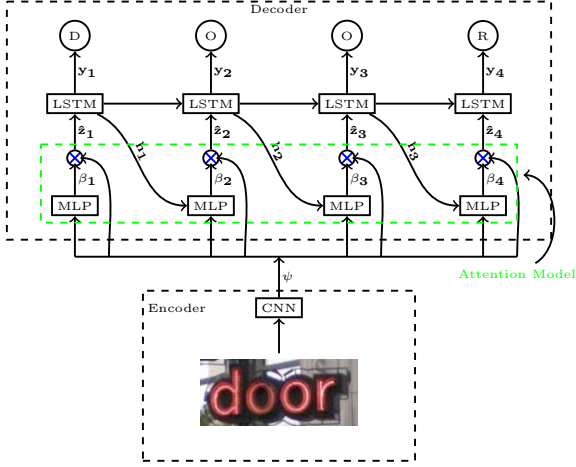


Figure 2. The proposed Encoder-decoder framework with attention model. The encoder uses a convolutional neural network to extract a set of features from the image. Specifically, we make use of the CNN model proposed by Jaderberg et al. [?] for scene text recognition – however we do not use the fully connected layer as a fixed-length representation as it is common in previous works. Instead, we take the features produced by the last convolutional layer. In this way we can produce a set of feature vectors, each of them linked to a specific spatial location of the image through its corresponding receptive field. This preserves spatial information about the image and reduces model complexity. Through the attention model, the decoder is able to use this spatial information to selectively focus on the most relevant parts of the image at every step.

Encoder: The encoder uses a convolutional neural network to extract a set of features from the image. Specifically, we make use of the CNN model proposed by Jaderberg et al. [?] for scene text recognition – however we do not use the fully connected layer as a fixed-length representation as it is common in previous works. Instead, we take the features produced by the last convolutional layer. In this way we can produce a set of feature vectors, each of them linked to a specific spatial location of the image through its corresponding receptive field. This preserves spatial information about the image and reduces model complexity. Through the attention model, the decoder is able to use this spatial information to selectively focus on the most relevant parts of the image at every step.

Thus, given an input image of a cropped word, the encoder generates a set of feature vectors:

$$\Psi = \{x_i : i = 1 \dots K\}, \quad (1)$$

where x_i denotes the feature vector corresponding to i^{th} part of the image. Each x_i corresponds to a spatial location in the image and contains the activations of all feature maps at that location in the last convolutional layer of the CNN.

Attention model: For the attention model, we adapt the soft attention model of [18] for image captioning, originally introduced by [2] for neural machine translation. In [18] slightly better results are obtained using the hard version of the model that focuses, at every time step, on a single feature vector. However, we argue that, in the case of text recognition, the soft version is more appropriate since a single character will usually span more than one spatial cell of the image corresponding to each of the feature vectors. The soft version of the model can combine several feature vectors with different weights into the final representation.

As shown in figure 2, the attention model generates, at every time step t , a vector \hat{z}_t that will be the input to the LSTM decoder. This vector \hat{z}_t can be expressed as a weighted combination of the set Ψ of feature vectors x_i extracted from the image:

$$\hat{z}_t = \sum_{i=1}^K \beta_{t,i} x_i \quad (2)$$

Thus, the vector \hat{z}_t encodes the relative importance of each part of the image in order to predict the next character for the underlying word. At every time step t , and for each location i a positive weight $\beta_{t,i}$ is assigned such that $\sum(\beta_i) = 1$. These weights are obtained as the softmax output of a Multi Layer Perpectron (denoted as Φ) using the set of feature vectors Ψ and the hidden state of the LSTM decoder at the previous time step, h_{t-1} . More formally:

$$\alpha_{ti} = \Phi(x_i, h_{t-1}) \quad (3)$$

$$\beta_{ti} = \frac{\exp(\alpha_{ti})}{\sum_{j=1}^K \exp(\alpha_{t,j})} \quad (4)$$

This model is smooth and differentiable and thus it can be learned using standard back propagation.

Decoder: Our decoder is a Long Short Term Memory (LSTM) network [5] which produces one symbol from the given symbol set L , at every time step. The output of the LSTM is a vector y_t of $|L|$ character probabilities which represents the probability of emitting each of the characters in the symbol set L at time t . It depends on the output vector of the soft attention model \hat{z}_t , the hidden state at previous step h_{t-1} and the output of the LSTM at previous step y_{t-1} . We follow the notation introduced in [18] where the network is described by:

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} T \begin{pmatrix} E y_{t-1} \\ h_{t-1} \\ \hat{z}_t \end{pmatrix} \quad (5)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (6)$$

$$h_t = o_t \odot \tanh(c_t), \quad (7)$$

where T is the matrix of weights learned by the network and i_t, f_t, c_t, o_t , and h_t are the input, forget, memory, output and hidden state of the LSTM, respectively. In the above definition, \odot denotes the element-wise multiplication and E is an embedding of the output character probabilities that is also learned by the network. σ and \tanh denote the activation functions that are applied after the multiplication by the matrix of weights

Finally, to compute the output character probability y_t , a deep output layer is added that takes as input the character probability at the previous step, the current LSTM hidden state, and the current feature vector. The output character probability is:

$$P(y_t | \Psi, y_{t-1}) \sim \exp(L_0(E y_{t-1} + L_h h_t + L_z \hat{z}_t)) \quad (8)$$

where L_0, L_h and L_z are the parameters of the deep output layer that are learned using back-propagation.

3. Inference

We use beam search over LSTM outputs to perform word inference. We first introduce the basic procedure, and then describe how we extend it to incorporate language models.

3.1. The basic inference procedure

Once the model is trained, we use a beam search to approximately maximize the following score function over every possible word: $\mathbf{w} = [c_1, \dots, c_n]$:

$$S(\mathbf{w}, x) = \sum_{t=1}^N \log(P(c_t|c_{t-1})), \quad (9)$$

where c_n is a special symbol signifying the end of a word, which immediately stops the beam search.

The beam search keeps track at every step of the top N most probable sequences of characters. For every active branch of the beam search, given the previous character of the sequence, c_{t-1} , the output character probability y_t of the LSTM is used to obtain $P(c_t|c_{t-1})$ for all characters c_t in the symbol set L .

3.2. Incorporating language models and Lexicon

Text is a strongly contextual. There are some strict constraints imposed by the grammar of the language. For example any word in English cannot carry more than two consecutive occurrences of any alphabet letter. Leveraging such knowledge can positively impact the final recognition output. Although the LSTM implicitly learns some dependences between consecutive characters, we show that adding an explicit language model that takes into account longer dependencies gives a significant boost to recognition accuracy.

In this work we use a standard n -gram based language model during inference to leverage the language prior. The character n -gram model gives probability of a character conditioned on k previous characters, where k is a parameter of the model:

$$\Theta(c_k|c_{k-1}, c_{k-2}, \dots, c_1) = \frac{\#(c_1 c_2 \dots c_{k-1})}{\#(c_k c_k \dots c_k)}, \quad (10)$$

where, $\#(c_1, \dots, c_n)$ is the number of occurrences of a particular substring in a training corpus.

Finally, the score function in equation 9 can be modified to take the n -gram language model into account as:

$$S(\mathbf{w}, x) = \sum_{t=1}^N \log(P(c_t|c_{t-1})) + \alpha \log \Theta(w_t|w_{t-1}, w_{t-2}, \dots, w_1) \quad (11)$$

At every step we fix the parameter k of the language model to the number of previously generated characters in order to take into account the longest possible sequence.

Although our method is originally designed for unconstrained text recognition, it can also leverage a lexicon whenever available. The use of a lexicon D can be integrated by modifying the beam search so that all active sequences

that do not correspond to any valid word are automatically removed from the beam.

This can be efficiently implemented by storing the lexicon in a trie structure and automatically removing from the beam search any alternative that do not correspond to any partial branch of the trie.

4. Experimental Results

4.1. Datasets and experimental protocols

We evaluate the performance of the proposed method using the following standard datasets.

Street View Text (SVT) dataset: this dataset contains 647 cropped word images downloaded from Google Street View. Results using the predefined lexicons defined by Wang *et al.* in [17] of 50 words for each image refereed as SVT-50.

ICDAR'03 text dataset: this dataset dataset contains 251 full images and 860 cropped word images [15]. We used the same protocol as [1], [10], [17] and evaluate cropped word images for which the groundtruth text contains only alphanumeric characters and contains at least three characters.

MSCOCO [16] text Dataset: This is a recently published dataset. This dataset is also challenging as none of the images are captured specifically with text recognition in mind. Also this dataset is much bigger than previous scene text datasets.

Synth90k text dataset: this dataset is used only for training [8]. It contains 9 million synthetically-generated text images. We use the official partition for training as in other works like [10].

Evaluation protocol: We use the standard evaluation protocol adopted in most previous work on text recognition in scene images [6], [10], [17]. The accepted metric is word level accuracy in percentage. SVT and ICDAR'03 are used for evaluation. For lexicon-based recognition, we used the same set of 50 for all images in for SVT and ICDAR'03 dataset, as proposed buy Wang *et al.* [17].

Implementation details: The CNN encoder used in this work is the Dictnet model by Jaderberg *et al.* [8]. Their deep convolutional network consists of four convolutional layers and two fully connected layers. In this work we used features from the last convolutional layer. Thus, the feature map used is of size 4×13 and therefore, the LSTM takes input in the form of 52×512 .

For lexicon-based recognition when we do not use the lexicon-based inference explained in section 3.2. Instead, we take the output of unconstrained recognition and find the closest word in the lexicon using the Levenshtein edit distance. For lexicon-based inference in unconstrained datasets (SVT and ICDAR'03) we use the 90k-words lexicon provided by Jaderberg *et al.* in [8]. The explicit language model is also learned using this 90k word lexicon.

The parameter α (see equation 11) to weight the language model with respect to LSTM character probability is

empirically established. In our experiments we found the best results with α between 0.25 to 0.3

4.2. Baseline performance analysis

In this section we analyze the impact on performance of all the components of the proposed model. We start with a baseline that consists of a simple one layer LSTM network as decoder, without any attention or explicit language model. As we are interested mainly in the impact of the attention model, we use a simple version in which CNN features from the encoder are fed to the LSTM only at the first time step. At every step the output character is determined based on the output of the previous step and the previous hidden state.

In an effort to evaluate each of our contributions, we trained the baseline system and our model with exactly the same training data. For this purpose we randomly sampled one million training samples from the Synth90k [8] dataset. For validation we used 300,000 samples randomly taken from the same synth90K dataset.

We present the results for each of the component of the framework as described above in Table 1. The attention model outperforms the baseline by a significant margin (around 7%). Also these results confirm the advantage of using an explicit language model in addition to the implicit conditional character probabilities learned by the LSTM model. Using the language model improves accuracy in another 7%. We also see that further constraining the inference with a dictionary does not improve the result much, probably because the language model is learned from the same 90K dictionary proposed by Jaderberg *et al.* in [8].

In comparison with other related works on unconstrained text recognition, it is noteworthy that with only one million training samples our complete framework can learn a better model than Jaderberg *et al.* [6] and obtain results that are close to other state-of-the-art methods that are using the whole 9 million sample training dataset (see table 2).

Methods	SVT
Baseline (LSTM-no attention)	61.7
Proposed (LSTM + attention model)	68.16
Proposed (LSTM + attention model + LM)	75.57
Proposed (LSTM + attention model+LM+dict)	76.04

TABLE 1. IMPACT OF THE DIFFERENT COMPONENTS OF OUR FRAMEWORK WITH RESPECT TO THE BASELINE. WE COMPARE THE BASELINE (LSTM WITH NO ATTENTION MODEL) WITH ALL THE VARIANTS OF THE PROPOSED METHOD

4.3. Comparison with state of the art

In this section we will compare our result with other related works on scene text recognition. The results of this comparison are shown in table 2 for SVT and ICDAR’03 and 3 for COCO dataset. First, we will discuss results on unconstrained text recognition which is the main focus of our work. Then, we will analyze results for lexicon-based recognition.

Unconstrained text recognition: apart from our method Jaderberg *et al.* [6], Lee *et al.* [10] and Bissacco *et al.* [3] are the only methods which are capable of performing totally unconstrained recognition of scene text. Among these methods, our visual attention based model performs significantly better than Bissacco *et al.* [3] and Jaderberg *et al.* [6] in both SVT and ICDAR’03 datasets. Our model also performs as good as Lee *et al.* [10] in SVT dataset and outperforms them by 3% in ICDAR’03 dataset, which is significant given the high recognition rates.

If we further compare our model with that of Lee *et al.* [10], that also uses different variants of RNN architectures and an attention model on top of CNN features, we find that they use recursive CNN features. They report that this gives an 8% increase in accuracy over the baseline. This success is due to the recurrent nature of the CNN feature which implicitly model the conditional probability of character sequences.using recursive CNN performs better than the traditional convolutional feature. However, the RNN architecture they use improves only 4% over the baseline. In contrast our method rely on traditional CNN features (which can possibly encodes the presence of individual characters as shown in [6] from lower convolutional layer preserving local spatial characteristics, which reduces the complexity of the model. In addition, as reported in table 1, our combination of LSTM and soft attention model achieves a much larger margin, 14%, over the baseline. Theses results show that a combination of local convolutional features using the context based attention attention performs better or comparable to the previous state-of-the- art results.

We provide the results on the COCO text dataset in Table 3: Being this the most recently released dataset in this domain, there are no published results that could be comparable with our work. To make a valid comparison we used two neural network based approaches by M. Jaderberg *et. al.* [9] as they have made their models available online. We also fine-tuned the models on COCO dataset which leads to significant improvement (last row in Table 3). We can see that our simplest model is comparable to Jaderberg’s results while including the explicit language model leads to a significant improvement by a large margin.

Lexicon-based recognition For SVT-50 we can observe that our method obtain a similar result than the best of the methods [8] specifically designed to work in a lexicon-based scenario. Comparing with methods for unconstrained text recognition, only the method of Lee *et al.* [10] outperforms our best setting. But as we have already discussed, part of this better performance can be explained by the use of the more complex recursive CNN features.

Concerning ICDAR’03-50 and ICDAR’03-full, our results, although do not beat current state of the art are very competitive and comparable to the best performing methods.

5. Conclusions

In this paper we proposed an LSTM-based visual attention model for scene text recognition. The model uses

	Methods	SVT-50	SVT	ICDAR'03-50	ICDAR'03-full	ICDAR'03
Lexicon-based	Almazan <i>et al.</i> [1]	89.2	-	-	-	-
	Lee <i>et al.</i> [11]	80.0	-	88.0	76.0	-
	Yao <i>et al.</i> [19]	75.9	-	88.5	80.3	-
	Rodriguez-Serrano <i>et al.</i> [13]	70.0	-	-	-	-
	Jaderberg <i>et al.</i> [7]	86.1	-	96.2	91.5	-
	Su and Luet <i>et al.</i> []	83.0	-	92.0	82.0	-
	Gordo <i>et al.</i> [4]	90.7	-	-	-	-
	*DICT Jaderberg <i>et al.</i> [8]	95.4	80.7	98.7	98.6	93.1
Unconstrained	Bissacco <i>et al.</i> [3]	90.4	78.0	-	-	-
	Jaderberget al. [6]	93.2	71.7	97.8	97.0	89.6
	Lee <i>et al.</i> [10]	96.3	80.7	97.9	97.0	88.7
	Proposed (LSTM + attention model)	91.7	75.1	93.4	91.0	89.3
	Proposed (LSTM + attention model + LM)	95.2	80.4	95.7	94.1	92.6
	Proposed (LSTM + attention model+LM+dict)	95.4	-	96.2	95.7	-

TABLE 2. SCENE TEXT RECOGNITION ACCURACY. “50” AND “FULL” DENOTE THE LEXICON SIZE USED FOR CONSTRAINED TEXT RECOGNITION AS DEFINED IN [17]. RESULTS ARE DIVIDED INTO LEXICON-BASED AND UNCONSTRAINED (LEXICON-FREE) APPROACHES. *DICT [8] IS NOT LEXICON-FREE DUE TO INCORPORATING GROUND-TRUTH LABELS DURING TRAINING.

Methods	Accuracy
Charnet [9]	24.72
Dictnet [9]	26.79
Proposed (LSTM + attention model)	24.11
Proposed (LSTM + attention model+LM)	33.67
Proposed (LSTM + attention model+LM+FT)	43.86

TABLE 3. PERFORMANCE OF OUR METHODS ON RECENTLY RELEASED COCO-TEXT DATASET. WE COMPARE DIFFERENT VARIANTS OF OUR METHOD USING ONLY THE ATTENTION MODEL, INTEGRATING EXPLICIT LANGUAGE MODEL AND ALSO FINE TUNING THE MODEL ON COCO TEXT DATASET).

convolutional features from a standard CNN as input to an LSTM network that selectively attends to parts of the image at each time step in order to recognize words without resorting to a fixed lexicon. We also propose a modified beam search strategy that is able to incorporate weak language models (n -grams) to improve recognition accuracy. Experimental results demonstrate that our approach outperforms or performs comparably to state-of-the-art approaches that use lexicons to constrain inferred output words. Experimental results shows that context plays a important part in case of real data, thus using a explicit language model always helps to improve the result.

In future we can extend the attention model for the text detection task, which will lead to an end-to-end framework for text recognition from images. Moreover, in our current framework convolutional features are taken from one single layer, which can lead to poorer results when the text is either too big or too small. This can be dealt with combining features from multiple layers.

References

- [1] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, “Word spotting and recognition with embedded attributes,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 12, pp. 2552–2566, 2014.
- [2] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *CoRR*, vol. abs/1409.0473, 2014. [Online]. Available: <http://arxiv.org/abs/1409.0473>
- [3] A. Bissacco, M. Cummins, Y. Netzer, and H. Neven, “Photoocr: Reading text in uncontrolled conditions,” in *2013 IEEE International Conference on Computer Vision*, Dec 2013, pp. 785–792.
- [4] A. Gordo, “Supervised mid-level features for word image representation,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [5] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [6] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep structured output learning for unconstrained text recognition,” *ICLR 2015*, 2014.
- [7] —, “Reading text in the wild with convolutional neural networks,” 2014.
- [8] —, “Synthetic data and artificial neural networks for natural scene text recognition,” in *NIPS Deep Learning Workshop 2014*, 2014.
- [9] M. Jaderberg, A. Vedaldi, and A. Zisserman, “Deep features for text spotting,” in *ECCV 2014*, 2014.
- [10] C. Lee and S. Osindero, “Recursive recurrent nets with attention modeling for OCR in the wild,” *CoRR*, vol. abs/1603.03101, 2016. [Online]. Available: <http://arxiv.org/abs/1603.03101>
- [11] S. Lee, M. S. Cho, K. Jung, and J. H. Kim, “Scene text extraction with edge constraint and text collinearity,” in *Proc. ICPR*, 2010.
- [12] L. Neumann and J. Matas, “Real-time scene text localization and recognition,” in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 3538–3545.
- [13] J. A. Rodriguez-Serrano, A. Gordo, and F. Perronnin, “Label embedding: A frugal baseline for text recognition,” *International Journal of Computer Vision*, vol. 113, no. 3, pp. 193–207, 2015.
- [14] D. L. Smith, J. Field, and E. Learned-Miller, “Enforcing similarity constraints with integer programming for better scene text recognition,” in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 73–80.
- [15] L. P. Sosa, S. M. Lucas, A. Panaretos, L. Sosa, A. Tang, S. Wong, and R. Young, “Icdar 2003 robust reading competitions,” in *In Proceedings of the Seventh International Conference on Document Analysis and Recognition*. IEEE Press, 2003, pp. 682–687.
- [16] A. Veit, T. Matera, L. Neumann, J. Matas, and S. Belongie, “Cocotext: Dataset and benchmark for text detection and recognition in natural images.”
- [17] K. Wang, B. Babenko, and S. Belongie, “End-to-end scene text recognition,” in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 1457–1464.
- [18] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” *CoRR*, vol. abs/1502.03044, 2015. [Online]. Available: <http://arxiv.org/abs/1502.03044>
- [19] C. Yao, X. Bai, B. Shi, and W. Liu, “Strokelets: A learned multi-scale representation for scene text recognition,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.