

Memory Matters: Convolutional Recurrent Neural Network for Scene Text Recognition

Qiang Guo*, Dan Tu, Guohui Li and Jun Lei

Department of Information System and Management

National University of Defense Technology

Email: {guoqiang05*, tudan, guohuili, junlei}@nudt.edu.cn

Abstract—Text recognition in natural scene is a challenging problem due to the many factors affecting text appearance. In this paper, we presents a method that directly transcribes scene text images to text without needing of sophisticated character segmentation. We leverage recent advances of deep neural networks to model the appearance of scene text images with temporal dynamics. Specifically, we integrates convolutional neural network (CNN) and recurrent neural network (RNN) which is motivated by observing the complementary modeling capabilities of the two models. The main contribution of this work is investigating how temporal memory helps in an segmentation free fashion for this specific problem. By using long short-term memory (LSTM) blocks as hidden units, our model can retain long-term memory compared with HMMs which only maintain short-term state dependences. We conduct experiments on Street View House Number dataset containing highly variable number images. The results demonstrate the superiority of the proposed method over traditional HMM based methods.

I. INTRODUCTION

Text recognition in natural scene is an important problem in computer vision. However, due to the enormous appearance variations in natural images, e.g. different fonts, scales, rotations, illumination conditions, it is still quite challenging.

Identifying the position of a character and recognizing it are two interdependent problems. Straight-forward methods treat the task as separate character segmentation and recognition[1], [2]. This paradigm is fragile in unconstrained natural images for it's difficult to deal with low resolution, low contrast, blurring, large diversity of text fonts and highly complicated background clutters.

Due to the shortcoming of these methods, algorithms combining segmentation and recognition were proposed. GMM-HMMs are the mostly used models, especially in speech and handwriting communities[3], [4], [5], [6].

In this paradigm, scene text images are transformed to frame sequences by sliding window. GMMs are used for modeling frame appearance and HMMs are used to infer the target labels of the whole sequence.[7]

The merit of this method is avoiding the need of fragile character segmentation. However HMMs have several obvious shortcomings, e.g. lacking of long context consideration, improper independent hypothesis etc.

As the developments of deep neural networks (DNNs) flourishing, convolutional neural networks (CNNs) have been used to form the hybrid CNN-HMM model[7], replacing GMMs as

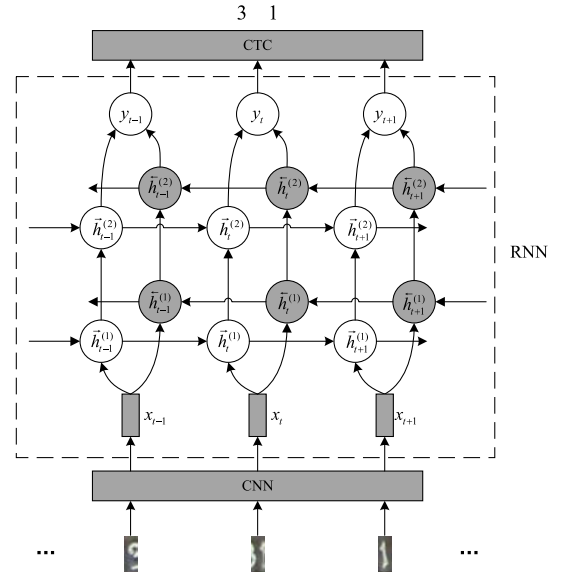


Fig. 1. The whole architecture of CRNN.

the observation model. The model generally performs better than the GMM-HMM model thanks to the strong representation capacity of CNN, however still doesn't eliminate the issues with HMM.

In this work, we address the issues of HMMs while keeping the algorithm free of segmentation. The novelty of our method is using Recurrent Neural Network (RNN), which has the ability of adaptively retaining long-term dynamic memory, as the sequence model. We combine CNN with RNN to utilize their representation abilities on different aspects.

RNN is a powerful connectionist model for sequences. Comparing with static feed-forward networks, it introduces recurrent connections enabling the network to maintain an internal state. It doesn't make any hypothesis on the independence of inputs, so each hidden unit can take into account more input information. Specifically, LSTM memory blocks are used enabling RNN to retain longer range of inter-dependences of input frames. Another virtue of using RNN

as the sequence model is the ease of build an end-to-end trainable system directly trained on the whole image without needing explicit segmentation. The main weakness of RNN is its feature extraction capability.

To alleviate the shortcomings of both HMM and RNN, we propose a novel end-to-end sequence recognition model named Convolutional Recurrent Neural Network (CRNN). The model is composed with hierarchical convolutional feature extraction layers and recurrent sequence modeling layers.

CNN is good at appearance modeling and RNNs have strong capacity for modeling sequences. The model is trained with Connectionist Temporal Classification (CTC)[8] object which enables the model directly learned from images without segmentation information.

Our idea is motivated by observing the complementary modeling capacities of CNN and RNN, and inspired by recent success applications of LSTM architectures to various sequential problems, such as handwriting[9], and speech recognition[10], image description[11], [12]. The whole architecture of our model is illustrated in [Figure 1](#).

II. RELATED WORK

In this section, we briefly survey methods that focus on sequence modeling without segmentation.

The paradigms of scene text recognition algorithms are similar with handwriting recognition. Straightforward methods[1], [2] are composed of two separated parts. A segmentation algorithm followed by a classifier to determine the category or each segment. Often, the classification results are post-processed to form the final results.

To eliminate the need of explicit character segmentation, many researchers use GMM-HMM for text recognition[13], [14]. GMM-HMM is a classical model widely used by the speech community. HMMs make it possible to model sequences without the necessity of segmentation. However, there are many shortcomings of GMM-HMM, which make it not widely used in scene text recognition. Firstly, GMM is a weak model for modeling characters in natural scene. Secondly, HMM has many limitations that are addressed in [section I](#).

To strengthen the representation capability of HMM based model, CNN is then used to replace GMM as the observation model which forms the hybrid CNN-HMM model[7], [15]. While improves the performance in comparison with GMM-HMM, it still doesn't eliminate the shortcomings of HMM.

Our idea is motivated by recent success of the RNN models applied to handwriting recognition[9], speech recognition[10] and image description[11], [12]. The main inspiration of our idea is observing the complementary modeling capacity of CNN and RNN. CNN can automatically learn hierarchical image features but only as a static model. RNN is good at sequence modeling while lacking the ability of feature extraction. We integrate the two models to form an end-to-end scene text recognition system.

Different with recent works[16] which use similar ideas, we investigate to use deep RNNs by stacking multiple recurrent hidden states on top of each other. Our experiment shows the improvements of the endeavor.

III. PROBLEM FORMULATION

We formulate scene text recognition as a sequence labeling problem by treating scene text as frame sequences. The label sequence is drawn from a fixed alphabet L . The length of the label sequence is not restricted to be equal to that of the frame sequence.

Each scene text image is treated as a sequence of frames denoted by $\mathbf{x} = (x_1, x_2, \dots, x_T)$. The target sequence is $\mathbf{z} = (z_1, z_2, \dots, z_U)$. We use bold typeface to denote sequences.

We constrain that $|\mathbf{z}| = U \leq |\mathbf{x}| = T$. The input space $\mathcal{X} = (\mathbb{R}^M)^*$ is the set of all sequences of M real valued vectors. The target space $\mathcal{Z} = L^*$ is the set of all sequences over the alphabet L of labels. We refer $\mathbf{z} \in L^*$ as a *labeling*.

Let S be a set of training samples drawn independently from a fixed distribution $\mathcal{D}_{\mathcal{X} \times \mathcal{Z}}$ composed of sequence pairs (\mathbf{x}, \mathbf{z}) . The task is to use S to train a sequence labeling algorithm $f : \mathcal{X} \mapsto \mathcal{Z}$ to label the sequences in a test set $S' \in \mathcal{D}_{\mathcal{X} \times \mathcal{Z}}$ as accurately as possible given the error criterion *label error rate* E^{lab} :

$$E^{lab}(h, S') = \frac{1}{Z} \sum_{(\mathbf{x}, \mathbf{z}) \in S'} ED(h(\mathbf{x}), \mathbf{z}) \quad (1)$$

where $ED(\mathbf{p}, \mathbf{q})$ is the *edit distance* between two sequences \mathbf{p} and \mathbf{q} .

IV. METHOD

A. The proposed model

The network architecture of our CRNN model is shown in [Figure 1](#). The model is mainly composed with two parts, a deep convolutional network for feature extraction and a bidirectional recurrent network for sequence modeling.

An scene text image is transformed into a sequence of frames which are fed into the CNN model to extract feature vectors.

The CNN model only map one frame feature to its corresponding output vector. The sequence of feature vectors are then used as the input of RNN which takes the whole sequence history into consideration.

The recurrent connections allow the network to retain previous inputs as memory in the internal states and discovery temporal correlations among time-steps even far from each other.

Given an input sequence \mathbf{x} , a standard RNN computes the hidden vector sequence $\mathbf{h} = (h_1, h_2, \dots, h_T)$ and output vector sequence $\mathbf{y} = (y_1, y_2, \dots, y_T)$ as following:

$$h_t = \mathcal{H}(W_{ih}x_t + W_{hh}h_{t-1} + b_h) \quad (2)$$

$$y_t = W_{ho}h_t + b_o \quad (3)$$

where the W terms denote weight matrices (e.g. W_{ih} is the input-hidden weight matrix), the b terms denote bias vectors (e.g. b_h is hidden bias vector), \mathcal{H} is the hidden layer activation function.

We stack a CTC layer on top of RNN. With the CTC layer, we can train the RNN model directly on the sequences' labellings with knowing frame-wise labels.

B. Feature extraction

CNNs[17], [18] have shown exceptionally powerful representation capability for images and have achieved state-of-the-art results in various vision problems. In this work, we build an CNN for feature extraction.

We use CNN as a transforming function $f(o)$ that takes an input image o and outputs an fixed dimensional vector x as the feature. The convolution and pooling operations in deep CNNs are specially designed to extract visual features hierarchically, from local low-level features to robust high-level ones. The hierarchically extracted features are robust to variable factors that characters facing in natural scene.

C. Dynamic modeling with Bidirectional RNN and LSTM

One shortcoming of conventional RNNs is that they only able to make use of previous context. However, it's reasonable to exploit both previous and future contexts. For scene text recognition, the left and right context are both useful for determining the category of a specific frame image.

In our model, we use Bidirectional RNN (BRNN)[19] to process sequential data from both directions with two separate hidden layers.

BRNN computes the *forward* hidden sequence \vec{h} , the *backward* hidden sequence \overleftarrow{h} . Each time-step y_t of the output sequence is computed by integrating both directional hidden states:

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_o \quad (4)$$

BRNN provides the output layer with complete past and future context for every time-step in the input sequence. During the forward pass, input sequence is fed to both directional hidden layers, and the output layer is not updated until both the two hidden layers have processed the entire input sequence. The backward pass of BPTT for BRNN is similar with unidirectional RNN, except that the output layer error δ is fed back to the both two directional hidden layers.

While in principle RNN is a simple and powerful model, in practice, it's unfortunately hard to train properly. RNN can be seen as a deep neural network unfolding in time. A critical problem when training deep networks is the *vanishing and exploding gradient* problem[20]. When error signal transmitting along the network's recurrent connections, it decays or blows up exponentially. Due to this, the range of context being accessed can be quite limited.

Long Short-term Memory (LSTM) block[21] is designed to control the input, output and transition of signal so as to retain useful information and discard useless one. LSTMs are used as memory blocks of RNN and can alleviate the vanishing and exploding gradient issues. They use a special structure to form memory cells. The LSTM updates for time-step t given inputs x_t , h_{t-1} and c_{t-1} are:

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \quad (5)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \quad (6)$$

$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (7)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \quad (8)$$

where σ is the logistic sigmoid function, i, f, o, c are respectively the *input gate*, *forget gate*, *output gate* and *cell* activation vectors, all of which are the same size as the hidden vector h .

The core part of LSTM block is the memory cell c that encodes the information of the inputs that have been observed up to that step. The gates determine whether the LSTM keeps the value from the gate or discards it. The input gate controls whether the LSTM considers its current input, the forget gate allows to forget its previous memory, and the output gate decides how much of the memory to transfer to the hidden state. Those features enable the LSTM architecture to learn complex long-term dependences.

D. Training

Both the CNN and RNN models are deep models. When stacking them together, it is difficult to train them together. Starting with an random initialization, the supervision information propagated from RNN to CNN are quite ambiguous. So we separately train the two parts.

The CNN model is trained with stochastic gradient descent. The samples for training CNN is got by performing forced-alignment on the scene text images with a GMM-HMM model[7].

For training the RNN model, we need a loss function that can directly compute the probability of the target labelling from the frame-wise outputs of RNN given the observations.

Connectionist Temporal Classification (CTC)[8] is an objective function designed for sequence labeling problem when the segmentation of data is unknown. It does not require pre-segmented training data, or post-processing to transform the network outputs into labelings. It trains the network to map directly from input sequences to the conditional probabilities of the possible labelings.

A CTC output layer contains one more unit than there are elements in the alphabet L , denoted as $L' = L \cup \{\text{null}\}$. The elements in L'^* are referred as *paths*. For an input sequence \mathbf{x} , the conditional probability of a path $\pi \in L'^T$ is given by

$$p(\pi|\mathbf{x}) = \prod_{t=1}^T y_{\pi_t}^t \quad (9)$$

where y_k^t is the activation of output unit k at time t . An operator \mathcal{B} is defined to merge the repeated labels and remove blanks. For example, $\mathcal{B}(-, 3, 3, -, -, 3, 2, 2)$ yields the labeling $(3, 3, 2)$. The conditional probability of a given labeling \mathbf{l} is the sum of the probabilities of all paths corresponding to it:

$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} p(\pi|\mathbf{x}). \quad (10)$$

We use CTC[8] as the objective function. A *forward-backward algorithm*[8] for CTC, which is similar to the forward-backward algorithm of HMM, is designed to effectively evaluate the probability.

The objective function for CTC is the negative log probability of the correct labelings for the entire training set:

$$O^{CTC} = - \sum_{(\mathbf{x}, \mathbf{z}) \in S} \ln(p(\mathbf{z}|\mathbf{x})). \quad (11)$$

Given the partial derivatives of some differential loss function \mathcal{L} with respect to the network outputs, we use *back propagation through time algorithm* (BPTT) to determine the derivatives with respect to the weights.

Like standard *back propagation*, BPTT follows the chain rule to calculate derivatives. The subtle difference is that the loss function depends on the activation of the hidden layer not only through its influence on the output layer, but also through the hidden layer of next time-step. So the error back propagation formula is

$$\delta_h^t = \theta'(a_h^t) \left(\sum_{k=1}^K \sigma_k^t w_{hk} + \sum_{h'=1}^H \sigma_{h'}^{t+1} w_{hh'} \right) \quad (12)$$

where

$$\delta_j^t \equiv \frac{\partial \mathcal{L}}{\partial a_j^t} \quad (13)$$

The same weights are reused at every timestep, so we sum the whole sequence to get the derivatives with respect to the network weights:

$$\frac{\partial \mathcal{L}}{\partial w_{ij}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial a_j^t} \frac{\partial a_j^t}{\partial w_{ij}} = \sum_{t=1}^T \delta_j^t b_j^t \quad (14)$$

E. Decoding

The *decoding* task is to find the most probable labeling \mathbf{l}^* given an input sequence \mathbf{x} :

$$\mathbf{l}^* = \arg \max_{\mathbf{l}} p(\mathbf{l}|\mathbf{x}). \quad (15)$$

We use a simple and effective approximation by choosing the most probable path, then get the labeling $\tilde{\mathbf{l}}^*$ corresponding to the path:

$$\tilde{\mathbf{l}}^* = \mathcal{B}(\arg \max_{\pi} p(\pi|\mathbf{x})). \quad (16)$$

V. EXPERIMENTS

This section presents the details of our experiments. We compare the CRNN model with 3 methods: (1) A baseline method by directly using CNN to predict the character at each time-step, with a post-processing procedure to merge repeated outputs; (2) The GMM-HMM model; (3) Hybrid CNN-HMM model.

A. Dataset

We explore the use of CRNN model on a challenging scene text dataset Street View House Numbers(SVHN)[2]. The dataset contains two versions of sample images. One contains only isolated 32×32 digits with 73257 for training and 26032 for testing. The other contains unsegmented full house number images containing variable number of unsegmented digits. The full number version is composed of 33402 training images and 13068 testing. House numbers in the dataset show quite large appearance variability, blurring and unnormalized layouts.

We use the isolated version of the dataset for training the CNN, which is then used to extract features for each sequence frame. The full house number version is used for HMM based methods and CRNN.

The training samples are randomly splitted out 10% as the validation set. The validation set is used only for tuning hyper parameters of different models. All the models use the same training, validation and testing set, which makes it fair to compare different models.

B. Implementation details

The full number images are normalized to the same height 32 while keeping the scale ratio, then transformed to frame sequences by 32×20 sliding window. Each frame o_t is fed into CNN to producing feature $x_t = f(o_t)$. We standardize the features by subtracting the global mean and dividing the standard deviation.

Our CNN model contains 3 convolution and pooling layer pairs and 2 fully connected layers. The activation neurons are all rectified linear units (ReLU)[22]. The output number sequence of the layers are 32, 32, 64, 128, 10. CNN is trained with stochastic gradient descent (SGD) under cross entropy loss with learning rate 10^{-3} and momentum 0.9.

We choose the output of CNN's first fully connected layer as frame features. The features are 128D and used for both HMM based models and CRNN.

We use 11 HMMs coinciding with the extended alphabet L' . All the HMMs are of 3-state left-to-right topology, except for the nul category which has 1 self-looped state. The GMM-HMM model is trained with Baum-Welch algorithm.

For the proposed CRNN model, we use a deep bidirectional RNN. We stack 2 RNN hidden layers, both of which are bidirectional. All hidden units are LSTM blocks. The CRNN model is trained with BPTT algorithm using SGD. We use a learning rate of 10^{-4} and a momentum of 0.9.

C. Results

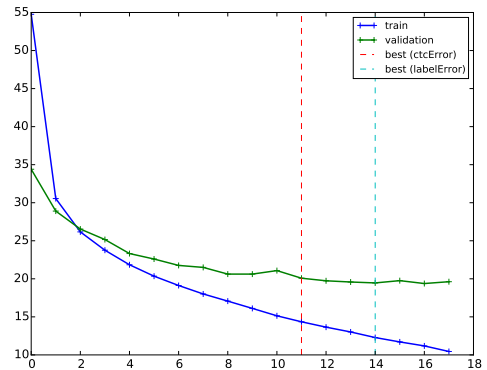


Fig. 2. Training curve during training of CRNN-2-layers.

a) *CNN based sequence labelling*: We train the CNN on the isolated version of SVHN, then use it to get the character predictions for each frame, choosing the most probable one as the result. After that, we merge consecutively repeated characters to get the final sequence labels. The recognition accuracy on the full number test set is only 0.23. Note that, the CNN model achieve an accuracy of 0.92 on the isolated test set. The correctly recognized house numbers by CNN

	CRNN-1-layer		CRNN-2-layers	
	epoch	accuracy	epoch	accuracy
CTC error	12	0.84	9	0.90
label error	15	0.86	10	0.91

TABLE I
COMPARISON OF DIFFERENT CRNN ARCHITECTURES.

are mostly contain only 1 or 2 digits. The simple experiment can gain us an intuition of how hard the problem is and how important the sequential information is.

b) *HMM based models*: We compare our method with 2 kinds of HMM based models. One is GMM-HMM model, the other is hybrid CNN-HMM model[7]. The number of mixture components is an important factor for GMM-HMM. We evaluate different number of Gaussian mixture components. The sequence accuracy stops improving at 800. The model tends to overfit when we continually increase the number of mixture components.

Hybrid CNN-HMM improves GMM-HMM by using CNN as the observation model. The training process is an iterative procedure, where network retraining is alternated with HMM re-alignment to generate more accurate state assignments.

c) *CRNN*: Recent developments of deep learning shows that *Deep* is an important factor for feed-forward models to gain stronger representation capability[23]. We evaluated two architectures of CRNN. One uses 1 hidden layer denoted as CRNN-1-layer, the other 2 hidden layers denoted as CRNN-2-layers. CRNN-1-layer has 128 LSTM memory cells, CRNN-2-layers has 128 for the first hidden layer and 32 for the second. **Figure 1** shows the architecture of CRNN-2-layers.

Experiment results are presented in **Table I**. *Epoch* column lists the epoch at which the best model reaches with respect to different error criterion. *Accuracy* column shows the sequence accuracy of the best model on test set.

As can be seen, the deeper architecture performs not only better but also with less training epochs. This is a surprising finding, as intuitively the deeper model has more parameters which makes it more difficult to train.

Model	Accuracy
CNN	0.23
GMM-HMM	0.56
Hybrid CNN-HMM	0.81
CRNN	0.91

TABLE II
PERFORMANCE COMPARISON OF DIFFERENT MODELS.

Performance comparison of CRNN with other models is represented in **Table II**. As shown by the experiments, CRNN outperforms CNN and both HMM based models.

VI. CONCLUSION

We have presented the Convolutional Recurrent Neural Network (CRNN) model for scene text recognition. It uses CNN to extract robust high-level features and RNN to learn sequence dependences. The model eliminates the need of character segmentation when doing scene text recognition. We apply our method on street view images and achieve

promising results. CRNN performs much better than HMM based methods. However, CNN is still trained separately. While the recognition process is segmentation-free, we still need cropped character samples for training the CNN. To eliminate the needing of cropped samples in training, we plan to investigate using forced alignment of GMM-HMM for bootstrapping of CRNN. Better method would be directly perform joint training of CNN and RNN from scratch. Another promising direction would be to investigate the potential of stacking more hidden LSTM layers of RNN.

ACKNOWLEDGMENT

This work was supported by Open Project Program of the State Key Laboratory of Mathematical Engineering and Advanced Computing(No.2015A04).

REFERENCES

- [1] A. Bissacco, M. Cummins, Y. Netzer, and H. Neven, "Photoocr: Reading text in uncontrolled conditions." in *ICCV*, 2013, pp. 785–792. **1, 2**
- [2] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," *NIPS workshop on deep learning and unsupervised feature learning*, vol. 2011, no. 2, p. 5, 2011. **1, 2, 4**
- [3] H. A. Bourlard and N. Morgan, *Connectionist speech recognition: a hybrid approach*. Springer Science & Business Media, 2012, vol. 247. **1**
- [4] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 20, no. 1, pp. 30–42, 2012. **1**
- [5] U.-V. Marti and H. Bunke, "Using a statistical language model to improve the performance of an hmm-based cursive handwriting recognition system," *International journal of Pattern Recognition and Artificial intelligence*, vol. 15, no. 01, pp. 65–90, 2001. **1**
- [6] S. Espana-Boquera, M. J. Castro-Bleda, J. Gorbe-Moya, and F. Zamora-Martinez, "Improving offline handwritten text recognition with hybrid hmm/ann models," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 33, no. 4, pp. 767–779, 2011. **1**
- [7] Q. Guo, D. Tu, J. Lei, and G. Li, "Hybrid cnn-hmm model for street view house number recognition," in *ACCV 2014 Workshops*, ser. Lecture Notes in Computer Science, C. Jawahar and S. Shan, Eds., 2015, vol. 9008, pp. 303–315. **1, 2, 3, 5**
- [8] A. Graves, S. Fernandez, F. J. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks." in *ICML*, 2006, pp. 369–376. **2, 3**
- [9] A. Graves and J. Schmidhuber, "Offline handwriting recognition with multidimensional recurrent neural networks." in *NIPS*, 2008, pp. 545–552. **2**
- [10] A. Graves and N. Jaitly, "Towards end-to-end speech recognition with recurrent neural networks." in *ICML*, 2014, pp. 1764–1772. **2**
- [11] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description." 2014. **2**
- [12] A. Karpathy and F.-F. Li, "Deep visual-semantic alignments for generating image descriptions." 2014. **2**
- [13] A. Vinciarelli, S. Bengio, and H. Bunke, "Offline recognition of unconstrained handwritten texts using hmms and statistical language models." 2004, pp. 709–720. **2**
- [14] M. Kozielski, P. Doetsch, and H. Ney, "Improvements in rwth's system for off-line handwriting recognition." in *ICDAR*, 2013, pp. 935–939. **2**
- [15] T. Bluiche, H. Ney, and C. Kermorvant, "Feature extraction with convolutional neural networks for handwritten word recognition." in *ICDAR*, 2013, pp. 285–289. **2**
- [16] K. Elagouni, C. Garcia, F. Mamalet, and P. Sbillot, "Text recognition in videos using a recurrent connectionist approach." in *ICANN (2)*, 2012, pp. 172–179. **2**
- [17] 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, Nov 1998. **3**

- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., 2012, pp. 1097–1105. 3
- [19] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *Signal Processing, IEEE Transactions on*, vol. 45, no. 11, pp. 2673–2681, Nov 1997. 3
- [20] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks." in *ICML (3)*, 2013, pp. 1310–1318. 3
- [21] S. Hochreiter and J. Schmidhuber, "Long short-term memory." 1997, pp. 1735–1780. 3
- [22] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines." in *ICML*, 2010, pp. 807–814. 4
- [23] R. Pascanu, G. Montufar, and Y. Bengio, "On the number of response regions of deep feed forward networks with piece-wise linear activations," in *International Conference on Learning Representations 2014 (ICLR 2014)*, Banff, Alberta, Canada, 2013. 5