

# SHARP ATTENTION FOR SEQUENCE TO SEQUENCE LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Attention mechanism has been widely applied to tasks that output some sequence  
2 from an input image. Its success comes from the ability to align relevant parts of  
3 the encoded image with the target output. However, most of the existing methods  
4 fail to build clear alignment because the aligned parts are unable to well represent  
5 the target. In this paper we seek clear alignment in attention mechanism through  
6 a *sharpen* module. Since it deliberately locates the target in an image region  
7 and refines representation to be target-specific, the alignment and interpretability  
8 of attention can be significantly improved. Experiments on synthetic handwritten  
9 digit as well as real-world scene text recognition datasets show that our approach  
10 outperforms the mainstream ones such as soft and hard attention.

## 11 1 INTRODUCTION

12 In modern sequence to sequence learning, attention mechanism has become a key building block,  
13 because it helps identify relevant parts of the input sequence and align them with the target output at  
14 each time step. Such alignment resembles fixation in human vision, where only the object of interest  
15 falls on the fovea in the retina, leading the visual scene outside to be largely ignored. Thanks to such  
16 ability to select perceptual information, attention mechanism has been successfully applied to many  
17 visual tasks, such as scene text recognition (Shi et al., 2019) and image captioning (Xu et al., 2015).

18 Although a variety of attention mechanisms have been proposed to build alignment, most of them  
19 fail to achieve clear alignment. Soft attention (Bahdanau et al., 2015; Luong et al., 2015), the most  
20 popular one, aligns a weighted average of the input sequence with the target output throughout the  
21 time. Since the weights are never zeros, irrelevant parts are inevitably involved in the alignment and  
22 may introduce distraction. For distinct alignment, hard attention (Xu et al., 2015) enforces exactly  
23 one input part is employed, regardless of whether it represents the target or not, and thus may still  
24 suffer from irrelevant parts. Besides, existing attention mechanisms often regard the input sequence  
25 as fixed during alignment establishment. If the target representation given by the selected part(s)  
26 is poor, there is no way to fix it. This is especially the case in visual sequence learning, where  
27 features are precomputed by a convolutional neural network (CNN) before being fed into attention  
28 mechanisms. As each feature only characterises a local fixed image region (i.e. the receptive field), it  
29 hardly covers the appearance of the target exactly, thus leading to noisy representation. See Fig. 1(b)  
30 for an example.

31 In this work we address the construction of clear alignment in attention mechanism. This is achieved  
32 by aligning the target output with image regions instead of features, which is a more natural approach  
33 to alignment for visual sequence learning. A *sharpen* module is then used to make the aligned  
34 region as specific as possible to the target, essentially a clear alignment. While it can take any form,  
35 the module explored in this paper consists of a localiser and an encoder. The former locates the  
36 target in the region, while the latter extracts features from the result for alignment. It is such accurate  
37 and specific representation that makes attention mechanism able to pay close attention to the target,  
38 leading to improved alignment and thus interpretability. The sharpener can be trained along with  
39 any sequence-to-sequence (Seq2Seq) model through back-propagation without extra supervision.  
40 Nonetheless, it is also possible to guide its training to further improve alignment quality if auxiliary  
41 information is available (see Sect. 4.1). Therefore, the sharpener naturally lends itself to direct  
42 attention manipulation, which is yet not available in most of the existing attention mechanisms.

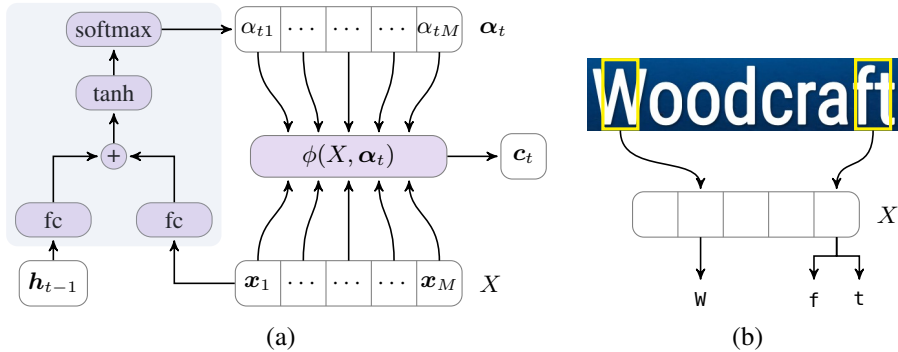


Figure 1: (a) An outline of the attention mechanism. A query vector  $h_{t-1}$  is compared with an input sequence  $X$  to figure out where to look via a multilayer perceptron (see the grey box), leading to a set of weights  $\alpha_t$ . A context vector  $c_t$  is then computed for alignment by a function  $\phi$ , where attention mechanisms generally differ. Soft attention computes  $c_t$  as a weighted average of  $X$  while hard attention does it by randomly sampling an element from  $X$ . (b) Poor target representation in attention mechanism. Each element of  $X$  only represents a fixed region defined by the receptive field (see yellow boxes). It is hard to accurately represent the target without the distraction of redundant or missing parts (see letters ‘f’ and ‘W’ for examples).

## 43 2 SHARP ATTENTION

### 44 2.1 ALIGNMENT

45 Let  $X = \{x_1, \dots, x_M\}$  be an input sequence of length  $M$ , where  $x_i \in \mathbb{R}^D$  is a feature vector  
 46 representing a region of an input image. Usually,  $X$  comes from the last convolutional layer of a  
 47 backbone network, and  $x_i$  delineates a region of fixed size given by the receptive field of that layer.  
 48 Similarly, we denote the output sequence of length  $N$  as  $Y = \{y_1, \dots, y_N\}$ , where  $y_t \in \mathbb{R}^K$  is  
 49 a one-of- $K$  encoded vector indicating a discrete token in a vocabulary of size  $K$ . Our goal is to  
 50 learn a model that can accurately predict  $y_t$  by choosing appropriate  $x_i$  in a sequential process.  
 51 This implicitly asks for building an alignment between  $X$  and  $Y$  at each time step. To facilitate the  
 52 construction, we introduce a latent variable  $A = \{a_1, \dots, a_N\}$ , where  $a_t \in \mathbb{R}^M$  is a one-of- $M$   
 53 encoded vector indicating the index of the selected feature vector at some time. For example,  $a_{ti} = 1$   
 54 refers to the  $i$ -th element of  $X$  (i.e.  $x_i$ ) being chosen to predict  $y_t$  at time  $t$ . Usually, the selected  
 55 feature is called context vector and denoted as  $c_t$ .

56 To learn the model, we maximise a conditional probability

$$\begin{aligned}
 p(Y|X) &= \sum_A p(Y, A|X) = \sum_A \prod_{t=1}^N p(y_t, A | \underbrace{y_1, \dots, y_{t-1}}_{\mathbf{y}_{<t}}, X) = \prod_{t=1}^N \sum_{\mathbf{a}_t} p(y_t, \mathbf{a}_t | \mathbf{y}_{<t}, X) \quad (1) \\
 &= \prod_{t=1}^N \sum_{\mathbf{a}_t} p(\mathbf{a}_t | \mathbf{y}_{<t}, X) p(y_t | \mathbf{y}_{<t}, X, \mathbf{a}_t) \equiv \prod_{t=1}^N \sum_{i=1}^M p(a_{ti} = 1 | \mathbf{y}_{<t}, X) p(y_t | \mathbf{y}_{<t}, x_i) \quad (2)
 \end{aligned}$$

57 where the last two terms in (1) are obtained by applying chain rule on  $Y$ , and by using the assumption  
 58 that  $y_t$  only depends on  $\mathbf{a}_t$  at time  $t$ , respectively. Equation (2) clearly defines two major components  
 59 to compute  $p(Y|X)$ . One is the chance of selecting each element of  $X$ , and the other is the likelihood  
 60 of the target token given the selection. However, the computation of the latter is impractical when  $M$   
 61 is large because every element of  $X$  has to be considered. Two typical approximations to (2) are thus  
 62 proposed, leading to the soft and hard attention mechanisms.

63 By using the first order Taylor expansion,<sup>1</sup> we obtain the loss function for soft attention,

$$\log p(Y|X) = \sum_{t=1}^N \log \left( \sum_{i=1}^M \alpha_{ti} p(y_t | \mathbf{y}_{<t}, x_i) \right) \approx \sum_{t=1}^N \log p \left( y_t | \mathbf{y}_{<t}, \sum_{i=1}^M \alpha_{ti} x_i \right), \quad (3)$$

<sup>1</sup>Let  $f(\cdot)$  be a function of some random variable. By using Taylor’s theorem, the first-order approximation to the expectation  $\mathbb{E}[f(\cdot)]$  is given by  $\mathbb{E}[f(\cdot)] \approx f(\mathbb{E}[\cdot])$ .

64 where we define  $\alpha_{ti} \equiv p(a_{ti} = 1 | \mathbf{y}_{<t}, X)$  to simplify the notation. In (3), the context vector  $\mathbf{c}_t$  is  
 65 given by  $\sum_{i=1}^M \alpha_{ti} \mathbf{x}_i$ , which means that  $\mathbf{y}_t$  is no longer predicted by a single element of  $X$  but rather  
 66 a weighted average of  $X$ , thus leading to the break in alignment. To pursue the alignment such that  
 67 each target token only depends on one element of  $X$ , hard attention instead estimates a variational  
 68 lower bound on  $\log p(Y|X)$  using Jensen’s inequality,

$$\log p(Y|X) = \sum_{t=1}^N \log \left( \sum_{i=1}^M \alpha_{ti} p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}_i) \right) \geq \sum_{t=1}^N \sum_{i=1}^M \alpha_{ti} \log p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}_i). \quad (4)$$

69 Now the optimisation of  $\log p(Y|X)$  can be thought as repeating the following steps until happy:  
 70 (i) estimating  $p(\mathbf{a}_t | \mathbf{y}_{<t}, X)$  by fixing all model parameters; (ii) modifying the parameters to maximise  
 71  $\log p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}_i)$  using each  $\mathbf{x}_i$ . The underlying idea is to increase  $\log p(Y|X)$  by iteratively raising  
 72 the lower bound. Since it is infeasible to consider every  $\mathbf{x}_i$  as aforementioned, approximation is often  
 73 adopted and achieved by Monte Carlo sampling (See Sect. 2.4 for details), thus leading  $\mathbf{c}_t$  to be the  
 74 sampled feature vector for hard attention.

## 75 2.2 LOSS FUNCTION

76 Intuitively, if each  $\mathbf{x}_i$  is an accurate representation of the target token when optimising (4), particularly  
 77 in the second step, there would be a tight gap between  $\log p(Y|X)$  and its lower bound. In contrast,  
 78 if any poor  $\mathbf{x}_i$  occurs, the gap may become large and thus result in performance degradation. Subject  
 79 to the fixed local region, it is unlikely for  $\mathbf{x}_i$  to well characterise the target token, which makes hard  
 80 attention hardly achieve clear alignment (see Fig. 1(b)). This motivates us to reformulate (4) for more  
 81 flexible representation of the target tokens.

82 Suppose we can break down the input image into a set of  $M$  local regions, each of which is sufficiently  
 83 large to cover objects of interest. We would like to maximise the marginal log-likelihood  $\log p(Y|R)$ ,  
 84 where  $R = \{r_1, \dots, r_M\}$  is the set of regions. Similarly, its lower bound  $\ell$  is given by

$$\log p(Y|R) \geq \sum_{t=1}^N \sum_{\mathbf{a}_t} p(\mathbf{a}_t | \mathbf{y}_{<t}, R) \log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t) \equiv \ell, \quad (5)$$

85 where  $\mathbf{a}_t$  is still a one-of- $M$  encoded vector but now refers to the index of the selected region at time  
 86  $t$ . By working with (5), we are not restricted to the representation given by  $X$  any more. Consider  
 87  $\mathbf{x}$  to be a function of  $r$  parameterised by the backbone network, e.g.,  $\mathbf{x} = f_g(r; \boldsymbol{\theta}_g)$ , where  $\boldsymbol{\theta}_g$   
 88 denotes all the weights in the network. By plugging  $X = \{f_g(r_1; \boldsymbol{\theta}_g), \dots, f_g(r_M; \boldsymbol{\theta}_g)\}$  into (4), it  
 89 is easy to see that the lower bound of  $\log p(Y|X)$  is equivalent to that of  $\log p(Y|R)$  when partially  
 90 parameterising the two terms  $p(\mathbf{a}_t | \mathbf{y}_{<t}, R)$  and  $\log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t)$  using  $\boldsymbol{\theta}_g$ . As  $\ell$  is valid for all  
 91 model parameters, it does not rely on any specific modelling. This allows us to separately model the  
 92 use of  $R$  in each term. For example, we may leverage the same backbone network for the first term to  
 93 get a rough idea on where to look, while deliberately design a network for the second term to sharpen  
 94 the focus. It is such flexible parameterisation that makes the construction of clear alignment possible.  
 95 Below we elaborate the modelling of each term in (5).

## 96 2.3 MODELLING

97 We use a variant of VGG (Shi et al., 2017) as the backbone network to not only create the set of  
 98 regions but also process it in  $p(\mathbf{a}_t | \mathbf{y}_{<t}, R)$ . When the backbone network is a CNN, we may take  
 99 advantage of the implicitly defined sliding widow for region generation. For example, our backbone  
 100 network effectively divides an input image of size  $100 \times 32$  into 24  $86 \times 46$  regions when extracting  
 101 features from the final layer (Araujo et al., 2019). Note that the generation of  $R$  is arbitrary and we  
 102 just use the sliding window for simplicity. To emphasise clear alignment with  $\log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t)$ ,  
 103 we leverage a sharpener module that consists of a localiser and an encoder. The former seeks the  
 104 target in the selected region while the latter extracts features from the result. While a natural choice  
 105 of the localiser is object detectors, we instead resort to spatial transformer networks (STNs, Jaderberg  
 106 et al. (2015)) for both computational and labelling efficiency. STN is a lightweight CNN that is able  
 107 to crop and transform an image region. Its training does not need expensive annotations such as  
 108 bounding boxes, which are usually unavailable in sequential learning tasks. The encoder can take  
 109 any form and we use the same CNN to the backbone to simplify the implementation. Let  $Z$  be the

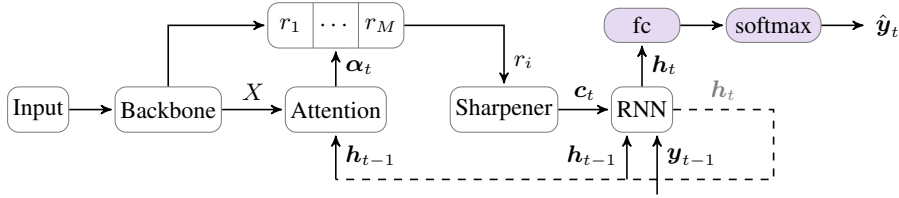


Figure 2: A generic Seq2Seq learning architecture with sharp attention. Given an input image, a backbone network breaks down it into a set of regions and extracts features from each region, leading to a sequence of feature vectors  $X$ . An attention mechanism (see Fig. 1(a)) uses  $X$  and a hidden state vector  $\mathbf{h}_{t-1}$  to compute a categorical distribution  $\alpha_t$ , based on which a region is randomly chosen and fed into a sharpener module to compute the context vector  $\mathbf{c}_t$  for clear alignment. A recurrent neural network (RNN) takes in  $\mathbf{h}_{t-1}$ ,  $\mathbf{c}_t$  and  $\mathbf{y}_{t-1}$  (the token at previous time step) to update its internal state, and then outputs current  $\mathbf{h}_t$  for token prediction and next iteration (dashed line).

110 encoder output, which is also a sequence of feature vectors similar to  $X$  with length dependent on the  
 111 size of the localisation result. The output of the sharpener is the context vector, whose computation  
 112 will be detailed in Sect. 2.5.

113 The dependency of the current token  $\mathbf{y}_t$  on all previous ones  $\mathbf{y}_{<t}$  over the time is often modelled by  
 114 an RNN, e.g., long short-term memory (LSTM, Hochreiter & Schmidhuber (1997)). Specifically, it is  
 115 defined by<sup>2</sup>

$$\mathbf{h}_t = f_r(\mathbf{h}_{t-1}, \mathbf{y}_{t-1}, \mathbf{c}_t; \boldsymbol{\theta}_r), \quad \mathbf{h}_0 \equiv f_i \left( \frac{1}{M} \sum_{i=1}^M \mathbf{x}_i; \boldsymbol{\theta}_i \right), \quad \mathbf{y}_0 \equiv \mathbf{0}, \quad (6)$$

116 where  $f_r(\cdot)$  refers to the non-linear function defined by LSTM with parameters  $\boldsymbol{\theta}_r$ ,  $\mathbf{h}_t$  is a hidden  
 117 state vector at time  $t$  that summarises the history tokens  $\mathbf{y}_{<t}$  and is initialised by a fully connected  
 118 layer  $f_i(\cdot)$  that takes an average of  $X$  as the input, and a zero vector is used as the initial token  $\mathbf{y}_0$ .

119 Now we are ready to define the two terms in (5). As we use the backbone network to process  $R$  in  
 120  $p(\mathbf{a}_t | \mathbf{y}_{<t}, R)$ , the probability of selecting a particular region is given by

$$\alpha_{ti} \equiv p(a_{ti} = 1 | \mathbf{y}_{<t}, R) \equiv p(a_{ti} = 1 | \mathbf{h}_{t-1}, X) = \text{softmax}(f_a(\mathbf{x}_i, \mathbf{h}_{t-1}; \boldsymbol{\theta}_a)), \quad (7)$$

121 where  $f_a(\cdot)$  is the attention function as described in Bahdanau et al. (2015) (Fig. 1(a)). The probability  
 122 of the output token  $\hat{\mathbf{y}}_t$  given all previous ones as well as the selected region is computed by

$$p(\mathbf{y}_t = \hat{\mathbf{y}}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t) \equiv p(\mathbf{y}_t = \hat{\mathbf{y}}_t | \mathbf{h}_{t-1}, \mathbf{c}_t) = \text{softmax}(f_e(\mathbf{y}_{t-1}, \mathbf{h}_t; \boldsymbol{\theta}_e)), \quad (8)$$

123 where  $f_e(\cdot)$  is a fully connected layer. An overview of the whole architecture is given in Fig. 2.

## 124 2.4 OPTIMISATION

125 The differentiation of  $\ell$  w.r.t. all model parameters yields the following learning rule<sup>3</sup>

$$\frac{\partial \ell}{\partial \theta} = \sum_{t=1}^N \sum_{\mathbf{a}_t} p(\mathbf{a}_t | \mathbf{y}_{<t}, R) \left[ \frac{\partial \log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t)}{\partial \theta} + \log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \mathbf{a}_t) \frac{\partial \log p(\mathbf{a}_t | \mathbf{y}_{<t}, R)}{\partial \theta} \right],$$

126 where  $\theta$  is a collection of model parameters, e.g.,  $\theta = \{\boldsymbol{\theta}_g, \boldsymbol{\theta}_s, \boldsymbol{\theta}_r, \boldsymbol{\theta}_i, \boldsymbol{\theta}_a, \boldsymbol{\theta}_e\}$  ( $\boldsymbol{\theta}_s$  for the sharpener  
 127 module). To reduce the computational cost as explained in Sect. 2.1, the derivative at time  $t$  is often  
 128 numerically approximated by the Monte Carlo method as follows

$$\frac{\partial \ell_t}{\partial \theta} \approx \frac{1}{S} \sum_{s=1}^S p(\hat{\mathbf{a}}_t^s | \mathbf{y}_{<t}, R) \left[ \frac{\partial \log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s)}{\partial \theta} + \log p(\mathbf{y}_t | \mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s) \frac{\partial \log p(\hat{\mathbf{a}}_t^s | \mathbf{y}_{<t}, R)}{\partial \theta} \right],$$

<sup>2</sup>Strictly speaking, the alignment in this modelling is no longer conditionally independent throughout the time as assumed in Sect. 2.1. In fact, this is the modelling used in both soft and hard attention mechanisms (Bahdanau et al., 2015; Xu et al., 2015). However, the discussion on why these mechanisms fail to achieve clear alignment still applies.

<sup>3</sup>We use the trick  $\nabla_{\theta} p(x; \theta) = p(x; \theta) \nabla_{\theta} \log p(x; \theta)$  in the derivation.

129 where  $S$  is the number of samples  $\hat{\mathbf{a}}_t^s$  drawn from a categorical distribution defined by  $p(\mathbf{a}_t|\mathbf{y}_{<t}, R)$   
 130 (Xu et al., 2015). Similar to Mnih et al. (2014) and Ba et al. (2015), this approximation yields a hybrid  
 131 loss function asking for different optimisation strategies for the two terms in the square brackets.  
 132 The former is optimised by a cross entropy loss together with the ground-truth token at time  $t$  and  
 133 gradient back-propagation. By regarding the accumulated sum of  $\log p(\mathbf{y}_t|\mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s)$  over the time  
 134 as a reward, the latter is achieved by the REINFORCE algorithm (Williams, 1992). To reduce the  
 135 high variance in gradient estimate caused by the unbounded  $\log p(\mathbf{y}_t|\mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s)$  (Ba et al., 2015),  
 136 we follow Xu et al. (2015) to introduce a moving average baseline

$$b_j = 0.9 \times b_{j-1} + 0.1 \times \frac{1}{NS} \sum_{s=1}^S \sum_{t=1}^N \log p(\mathbf{y}_t|\mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s),$$

137 where  $j$  is the index of the mini-batch. Finally, we use the following learning rule for optimisation

$$\frac{\partial \ell}{\partial \theta} \approx \frac{1}{S} \sum_{s=1}^S \sum_{t=1}^N p(\hat{\mathbf{a}}_t^s|\mathbf{y}_{<t}, R) \left[ \frac{\partial \log p(\mathbf{y}_t|\mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s)}{\partial \theta} + \lambda_r (\log p(\mathbf{y}_t|\mathbf{y}_{<t}, R, \hat{\mathbf{a}}_t^s) - b) \frac{\partial \log p(\hat{\mathbf{a}}_t^s|\mathbf{y}_{<t}, R)}{\partial \theta} \right], \quad (9)$$

138 where  $\lambda_r$  is a learning hyper-parameter. We do not add entropy  $H[\mathbf{a}_t]$  to (9) to further reduce gradient  
 139 variance as in Xu et al. (2015), because it encourages a uniform distribution and breaks the alignment.

## 140 2.5 CONTEXT VECTOR

141 Rather than compute  $\mathbf{c}$  from  $X$  like both soft and hard attention do, the proposed sharp attention  
 142 leverages the encoder output. Specifically, we explore three ways to compute  $\mathbf{c}_t$  given  $Z_t$ , the encoder  
 143 output at time  $t$ . The first one is *pooling*, where an average pooling of  $Z_t$  is used for  $\mathbf{c}_t$ . Inspired  
 144 by the glimpse idea (Mnih et al., 2014), we combine the result from pooling with  $\hat{\mathbf{x}}_t^s$ , the feature  
 145 vector associated with the sampled region, to incorporate both fine and coarse representation of the  
 146 target token, leading to the second way—*chain*. Alternatively, we may compute a weighted average  
 147 of the feature set  $\{Z_t, \hat{\mathbf{x}}_t^s\}$  for better mixed representation. With  $\mathbf{h}_{t-1}$ , the weights can be learned in  
 148 a similar way to the soft attention as shown in Fig. 1(a). We call this last approach to  $\mathbf{c}_t$  *weighting*.

## 149 3 RELATED WORK

150 In Xu et al. (2015), the alignment is treated as a latent variable to help define a loss function, which is  
 151 optimised by gradually increasing a variational lower bound on marginal log-likelihood of the target  
 152 output given the input sequence. Later, a variety of variational inference techniques (Lawson et al.,  
 153 2018; Deng et al., 2018; Bahuleyan et al., 2018) are proposed to further reduce the gap between the  
 154 lower bound and the marginal log-likelihood. Alternative approaches to approximating the marginal  
 155 log-likelihood can be found in Shankar et al. (2018) and Shankar & Sarawagi (2019). The former  
 156 cherry-picks a set of alignments, computes the log-likelihood conditioned on each alignment and  
 157 averages the results, while the latter extends the idea by enforcing Markov property on adjacent  
 158 alignments. Instead of approximation, Wu et al. (2018) attempted to compute exact marginal log-  
 159 likelihood by assuming that each alignment is conditionally independent across time steps. All of the  
 160 above methods regard the input sequence as fixed in optimisation and thus cannot tailor the input  
 161 part(s) to clear alignment.

162 The idea of using relevant parts to improve attention has been explored in various tasks. Mei et al.  
 163 (2016) adjusted the weights of the input parts resulting from soft attention to highlight the most  
 164 relevant ones for selective generation. A similar work can be found in Nallapati et al. (2016), where  
 165 keywords are interwoven with the sentences in which they lie for text summarisation by applying soft  
 166 attention to sentences and words respectively and rescaling word weights. Instead of reweighting,  
 167 Cheng et al. (2017) used character-level masks to guide the selection of useful parts for scene text  
 168 recognition. None of these methods build clear alignment due to the use of soft attention.

169 Our work is closely related to Xu et al. (2015) and Ba et al. (2015). We generalise the former’s  
 170 mathematical formulation on hard attention by introducing flexible representation of objects of interest  
 171 via a sharpener module. While the generalisation appears similar to the latter, we use it to tackle the  
 172 alignment issue in attention mechanisms rather than develop a new Seq2Seq model. Our modelling  
 173 also differs. Instead of seeking desired objects within the whole image, a divide-and-conquer scheme

174 is used to gradually narrow the search range for accurate localisation. Another difference in modelling  
 175 is that in Ba et al. (2015) prediction will not happen until a series of localisation across predefined  
 176 time steps whereas in our work that immediately follows localisation at each time.

## 177 4 EXPERIMENTS

178 We demonstrate the efficacy of the proposed method in two different scenarios of increasing difficulty:  
 179 (i) synthetic handwritten digit recognition and (ii) real-world scene text recognition. Rather than strive  
 180 for state-of-the-art results, the focus here is to highlight (i) the performance of Seq2Seq models can  
 181 be boosted if attention mechanism really yields clear alignment, that is, paying attention to the target  
 182 object, and (ii) the proposed sharp attention is an effective approach to reaching the goal. Therefore,  
 183 our vanilla system is built upon off-the-shelf modules and was trained without sophisticated parameter  
 184 tuning schemes. Below we describe some common choices for all scenarios.

185 **Implementation** All images are converted to grey scale and resized to  $100 \times 32$ . A variant of VGG  
 186 (Shi et al., 2017) is then used as the backbone to extract a  $24 \times 1$  feature map from each resized  
 187 image as well as create the set of regions, whose height is clipped to 32. The input sequence  $X$   
 188 is obtained by splitting the feature map along its width, leading the dimensionality of each  $x_i$  to  
 189 be the feature map depth (i.e. 512). Before feeding it into the attention mechanism, we follow  
 190 Shi et al. (2017) to further process  $X$  to capture long-range contextual information with a 2-layer  
 191 bidirectional LSTM, where each layer has a forward LSTM and a backward LSTM, each having 256  
 192 hidden units. The depth of the attention mechanism is set to 256 and so is the number of the hidden  
 193 units of the associated LSTM, which runs over some time to predict the output sequence  $Y$ . The  
 194 number of time steps is set to the maximum transcription length in each scenario. As in Sutskever  
 195 et al. (2014), an end-of-sequence token is used to indicate the finish of prediction. The STN in the  
 196 sharpener is composed of a localisation network, a grid generator and a sampler. Given a region,  
 197 the localisation network, achieved by the one described in Liu et al. (2016), uses it to estimate an  
 198 affine transformation, which is then used by the grid generator to place a set of control points on the  
 199 region. By sampling the intensity value at each control point in a way similar to Shi et al. (2019), the  
 200 sampler produces a patch of given size as the STN output. When multiple STNs are used, the output  
 201 of the previous one is used as the input to the next one. The output of the last STN is plugged into  
 202 the encoder in the sharpener to compute  $Z$ . The whole system was implemented using TensorFlow  
 203 (Abadi et al., 2016) and the code will be released in the near future.

204 **Training** Three kinds of Seq2Seq models were trained from scratch in terms of the attention  
 205 mechanism used. Specifically, the soft and hard models were learned via corresponding attention  
 206 mechanisms respectively. Unlike the previous two baseline models, the sharp models were obtained  
 207 by the sharpener with the context vector schemes described in Sect. 2.5. A stochastic gradient descent  
 208 method, ADADELTA (Zeiler, 2012), was used to learn the model parameters until certain number  
 209 of iterations in different scenarios. The learning rate was constant and set to 1.0 and the decay rate  
 210 was 0.95. In addition, all model parameters were regularised by an  $L_2$  norm with a weight decay  
 211 of  $4 \times 10^{-5}$ . All experiments were done with a batch size of 192 (per GPU) on a workstation of 4  
 212 NVIDIA GEFORCE RTX 2080 Ti GPUs. The number of samples  $S$  was set to the batch size.

213 **Evaluation** A prediction is correct if the predicted transcription matches ground truth. We reported  
 214 the proportion of correct predictions on each testing dataset. As in Shi et al. (2017), all transcriptions  
 215 were converted to lower cases and had punctuations ruled out before evaluation if applicable.

### 216 4.1 HANDWRITTEN DIGIT RECOGNITION

217 We randomly chose  $l$  images from the MNIST dataset (Lecun et al., 1998), resized them to  $32 \times 32$   
 218 and concatenated them horizontally, leading to an image of  $l$  handwritten digits. For each  $l$  in  $\{5,$   
 219  $7, 9, 11, 13\}$ , we created 20,000 images for training and 10,000 for testing by selecting from the  
 220 MNIST training and testing datasets respectively, leading to a normal handwritten digit dataset. To  
 221 introduce some distortion, we repeated the above procedure for a rotated dataset by randomly turning  
 222 the selected images around  $y$ -axis within  $[-30^\circ, 30^\circ]$  before concatenation. Examples of the generated  
 223 images are given in Fig. 3. We trained all models with the resulting datasets for 30,000 iterations  
 224 by setting  $\lambda_r = 1.0$  when applicable. For better localisation, we upsampled the set of regions along

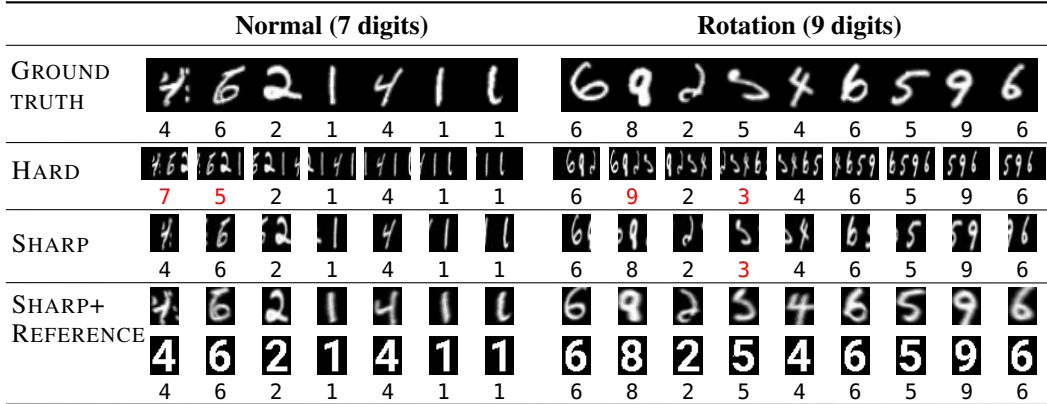


Figure 3: From top to bottom, we show examples of the created handwritten digit images and transcriptions (row 2), and patches for context vector computation as well as predicted digits (failures shown in red) at each time step (rows 3-5). Note that the patches are the receptive fields for the hard model and the output of the STN for the pooling based sharp models. The reference images are given below the patches in the last row.

Table 1: Recognition accuracy of all models for the handwritten digit datasets.

Model	Normal					Rotation					Mean Acc.
	5	7	9	11	13	5	7	9	11	13	
<b>Baseline</b>											
SOFT	97.2	96.6	95.0	91.9	82.2	96.3	95.4	93.6	89.1	78.3	91.6
HARD	97.2	96.4	94.1	91.6	81.7	96.6	95.4	93.5	89.2	78.6	91.4
<b>Sharp</b>											
POOLING	97.8	97.5	96.6	94.9	92.8	97.1	96.4	95.0	92.8	89.6	95.1
CHAIN	97.5	97.3	96.3	95.0	93.7	97.3	96.7	95.5	94.1	91.7	95.5
WEIGHTING	95.0	95.4	95.0	92.6	90.2	94.3	94.9	94.6	91.7	87.9	93.2
<b>Pooling-based Sharp+Reference</b>											
AFFINE	<b>98.1</b>	<b>97.7</b>	<b>96.8</b>	<b>96.0</b>	<b>94.4</b>	<b>98.0</b>	<b>97.1</b>	<b>96.1</b>	<b>95.2</b>	<b>93.0</b>	<b>96.2</b>

225 the width with a scale factor of 1.8 before plugging them into the sharpener, which was efficiently  
 226 achieved by running region generation with  $180 \times 32$  images. The patch output by the STN had a size  
 227 of  $24 \times 32$ .

228 Table 1 shows that all sharp models significantly outperform the baseline, demonstrating the efficacy  
 229 of clear alignment in attention mechanism. This can be easily seen from the pooling based model  
 230 whose context vector purely results from the sharpener. Figure 3 also illustrates how attention can  
 231 benefit from the sharpener. Take the rotation case for example. To predict digit ‘8’ (the second  
 232 column from left), hard attention chose a feature corresponding to a region filled with four digits. Due  
 233 to the distraction of irrelevant digits (e.g., ‘6’, ‘2’ and ‘5’), the feature failed to precisely represent  
 234 ‘8’, thus giving wrong result. Although sharp attention selected the same region, it avoided most  
 235 of the distraction by deliberately locating ‘8’ in that region, thus leading to more accurate and  
 236 specific representation as well as correct prediction. Besides, the resulting attention also has better  
 237 interpretability since it is more focused. The above results testify that the performance of Seq2Seq  
 238 models can be largely improved when attention mechanism is really focused on the object of interest.

239 To show that the sharpener allows for external supervision, a set of  $24 \times 32$  reference images for each  
 240 digit (see Fig. 3) was created with the Roboto Bold typeface of font size 36.<sup>4</sup> The supervision was  
 241 achieved by introducing an  $L_2$  image similarity loss of a weight of 1.0 to (9) to minimise the intensity  
 242 difference between the patch and the reference. By showing what a desired digit would look like, the

<sup>4</sup>The font is available at <https://fonts.google.com/specimen/Roboto>. We used Pygame for rendering.

Table 2: Recognition accuracy of all models for the scene text recognition datasets.

Model	IIIT	SVT	IC03	IC13	IC15	SP	Mean Acc.
	3000	647	867	857	1811	645	
<b>Baseline</b>							
SOFT	77.9	78.8	87.8	86.1	61.4	62.5	75.8
HARD	77.6	78.8	<b>89.2</b>	86.0	59.9	63.9	75.9
<b>Sharp+One STN</b>							
CHAIN	76.9	78.1	88.7	<b>88.8</b>	61.1	<b>65.6</b>	76.5
<b>Sharp+Two STNs</b>							
POOLING	73.9	72.6	83.6	83.3	54.2	60.2	71.3
CHAIN	77.9	<b>80.1</b>	89.0	87.7	<b>62.0</b>	65.1	<b>77.0</b>
WEIGHTING	77.7	78.5	89.0	87.4	61.5	64.0	76.4

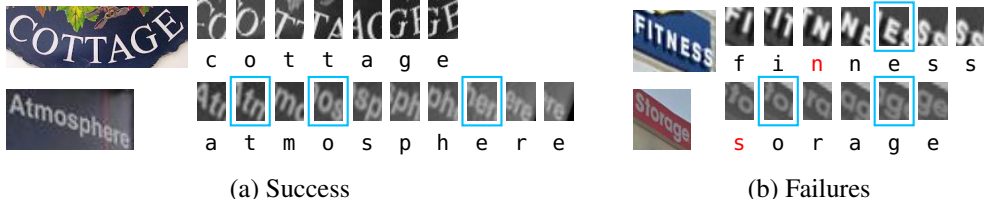


Figure 4: Examples of recognition results for scene text recognition. All real-world testing images are shown as is without rescaling and grey scale conversion. We leverage the chain based sharp model to generate the patches localised by the sharpener and predicted tokens (failures highlighted in red) across time steps. The results from other sharp models look similar to what has been shown here.

243 pooling based sharpener gives better localisation where all digit shape is well preserved with little  
 244 distraction, compared to its counterpart without such guidance (see Fig. 3). This, in turn, improves  
 245 the representation of the context vector for alignment, leading to a boost of the model performance.  
 246 Table 1 shows that such improvement is the most remarkable when handling images of more digits  
 247 (e.g., 13). This is not surprising because the receptive fields in such case have more digits filled and  
 248 thus make it difficult for the baseline models to decide where to look. Even if compared with the best  
 249 one amongst all sharp models trained without external supervision, i.e. the chain based model, the  
 250 improvement is still noticeable, further demonstrating the importance of accurate and target-specific  
 251 representation in attention mechanism.

## 252 4.2 SCENE TEXT RECOGNITION

253 To further show the effectiveness of the proposed sharp attention, we applied it to a real-world visual  
 254 sequence learning task, scene text recognition. The training datasets include MJSynth (Jaderberg et al.,  
 255 2014) and SynthText (Jaderberg et al., 2016), while the testing ones consist of IIIT5K-Words (Mishra  
 256 et al., 2012), Street View Text (Wang et al., 2011), ICDAR2003 (Lucas et al., 2003), ICDAR2013  
 257 (Karatzas et al., 2013), ICDAR2015 (Karatzas et al., 2015) and SVT Perspective (Phan et al., 2013),  
 258 which are short for IIIT, SVT, IC03, IC13, IC15 and SP respectively. There are 8.9 million images  
 259 in the MJSynth dataset and 5.5 million in SynthText by cropping text regions and ruling out non-  
 260 alphanumeric characters. The number of images in each testing dataset is detailed in Table 2, where  
 261 the three datasets, IC03, IC13 and IC15, were prepared by following the protocol in Baek et al. (2019).  
 262 The total number of testing images is 7, 827. Note that some of these datasets (e.g., IC15 and SP) are  
 263 quite challenging due to various nuisance factors, such as poor lighting and geometry change (Fig. 4).

264 For fast evaluation of various model configurations, we randomly sampled 2 million labelled images  
 265 from the MJSynth dataset. The sampling was done by following the distribution (i.e. histogram of  
 266 the transcription length) of the original dataset. We set the maximum transcription length to 16 to  
 267 enable a large batch size (i.e. 192) for robust training. For sharp models, we upsampled the regions  
 268 along the width with a scale factor of 2.0 for better localisation. To explore the effects of using  
 269 multiple STNs for localisation, we first used two STNs to train all sharp models, and then just used



270 the second one to train a chain based model for comparison. To see whether localisation benefits  
 271 from a coarse-to-fine search strategy, the first STN was designed to estimate a simple transformation  
 272 (i.e.  $x$ -direction translation) but output a large patch (i.e.  $64 \times 32$ ), while the second STN had the same  
 273 configuration as used in Sect. 4.1. All models were trained for 400,000 iterations with the sampled  
 274 dataset, and no reference images were used. To train the hard and sharp models we used  $\lambda_r = 0.1$ .

Table 3: Mean entropy of different attention mechanisms for the scene text recognition datasets.

Model	IIT	SVT	IC03	IC13	IC15	SP	Mean
SOFT	1.06	1.01	1.06	1.00	1.15	1.12	1.07
HARD	0.63	0.61	0.67	0.63	0.63	0.62	0.63
SHARP	<b>0.42</b>	<b>0.36</b>	<b>0.42</b>	<b>0.39</b>	<b>0.38</b>	<b>0.38</b>	<b>0.39</b>

Table 4: Sharp attention vs. soft attention in a published scene text recognition work.

Model	IIT	SVT	IC03	IC13	IC15	SP	Mean Acc.
SOFT(Baek et al., 2019)	<b>84.3</b>	83.8	93.1	91.9	70.8	71.9	82.6
SHARP	83.9	<b>85.8</b>	<b>94.3</b>	<b>92.2</b>	<b>71.3</b>	<b>73.6</b>	<b>83.5</b>

275 From Table 2, we see that the chain based one again works the best amongst all sharp models of two  
 276 STNs. It beats the baseline models remarkably on most of the datasets, whereas the other two either  
 277 moderately outperform or lag behind the baseline. It also beats its counterpart of single STN, even  
 278 though the latter is slightly better than the baseline as well. The result from the winning sharp model  
 279 further testifies our hypothesis on clear alignment in large-scale real-world datasets, which is also  
 280 revealed by Fig. 4. Whenever prediction is successfully performed, there is sensible localisation of  
 281 the target token. In fact, the sharpener attempts to highlight the token by placing it in the patch centre  
 282 (see Fig. 4 for examples surrounded by blue boxes). This is an encouraging result given that the  
 283 sharpener was trained in a data-driven manner with merely sequential labelling (i.e. transcriptions).  
 284 However, the localisation is by no means satisfactory since all patches have some sort of distractions,  
 285 such as skewed target tokens and irrelevant parts of adjacent tokens. Both this observation and the  
 286 benefit of multi-STN suggest a potential increase in accuracy if the sharpener is properly designed  
 287 such that it can produce good localisation as shown in Fig. 3, which is beyond the scope of this paper.

288 In Fig. 3, we have shown sharp attention is more focused and yields better interpretability. This can  
 289 be evaluated by  $H[\mathbf{a}_t]$ , entropy of the categorical distribution defined by  $\alpha_t$  (Shankar & Sarawagi,  
 290 2019). It measures attention uncertainty, that is, the lower entropy, the better alignment and thus  
 291 interpretability. Averaging  $H[\mathbf{a}_t]$  across all valid time steps leads to the entropy for an image. We  
 292 reported the mean of such entropy on each testing dataset for various attention mechanisms in  
 293 Table 3. The results clearly show that sharp attention indeed boosts interpretability since entropy is a  
 294 logarithmic metric. We only reported the entropy from the best sharp model in Table 2.

295 Finally, we used the full datasets (i.e. MJSynth & SynthText) to train a sharp model with the same  
 296 configuration to the best one in Table 2. To fairly compare the proposed attention with other attention  
 297 mechanisms in existing scene text recognition works, we reported the performance of the sharp model  
 298 and a model (i.e. VGG+BiLSTM+Attn) based on soft attention trained with the same datasets by  
 299 Baek et al. (2019) in Table 4. The two models share the same backbone and RNN decoder. They only  
 300 differ in attention mechanism. Table 4 further shows the superiority of our method.

## 301 5 CONCLUSION

302 We have described a novel attention mechanism that is able to build clear alignment between relevant  
 303 regions in the input image and the target output. This is achieved by a generic sharpener module that  
 304 computes accurate representation of the targets across time steps. Experimental results show that a  
 305 vanilla implementation of our method can significantly beat soft and hard attention on both synthetic  
 306 and real-world datasets in terms of performance and interpretability, without bells and whistles such  
 307 as the auxiliary model in Ba et al. (2015) and sophisticated training schemes in Xu et al. (2015). We  
 308 plan to apply our method to more visual sequence learning tasks in the future.

## 309 6 REPRODUCIBILITY STATEMENT

310 The implementation details have been given in lines 185–202. Although most of the training  
 311 parameters have been described in lines 207–212, some task-specific setting can be found in lines  
 312 223–227 and 266–274 respectively. The generation of synthetic handwritten digit dataset has been  
 313 detailed in lines 217–222, while the preparation of real-world scene text recognition datasets has  
 314 been elaborated in lines 254–261.

## 315 REFERENCES

- 316 Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu  
 317 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg,  
 318 Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete  
 319 Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: A system for large-scale  
 320 machine learning. In *Proceedings of the 12th USENIX Conference on Operating Systems Design  
 321 and Implementation*, pp. 265–283, 2016.
- 322 André Araujo, Wade Norris, and Jack Sim. Computing receptive fields of convolutional neural  
 323 networks. *Distill*, 2019. doi: 10.23915/distill.00021. [https://distill.pub/2019/computing-  
 324 receptive-fields](https://distill.pub/2019/computing-receptive-fields).
- 325 Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. Multiple object recognition with visual  
 326 attention. In *Proceedings of the International Conference on Learning Representations*, 2015.
- 327 Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoon Yun,  
 328 Seong Joon Oh, and Hwalsuk Lee. What is wrong with scene text recognition model compar-  
 329 isons? Dataset and model analysis. In *Proceedings of IEEE Conference on International  
 330 Conference on Computer Vision*, pp. 4715–4723, 2019.
- 331 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly  
 332 learning to align and translate. In Yoshua Bengio and Yann LeCun (eds.), *Proceedings of the  
 333 International Conference on Learning Representations*, 2015.
- 334 Hareesh Bahuleyan, Lili Mou, Olga Vechtomova, and Pascal Poupart. Variational attention for  
 335 sequence-to-sequence models. In *Proceedings of the International Conference on Computational  
 336 Linguistics*, pp. 1672–1682, 2018.
- 337 Zhanzhan Cheng, Fan Bai, Yunlu Xu, Gang Zheng, Shiliang Pu, and Shuigeng Zhou. Focusing  
 338 attention: Towards accurate text recognition in natural images. In *Proceedings of IEEE Conference  
 339 on International Conference on Computer Vision*, pp. 5086–5094, 2017.
- 340 Yuntian Deng, Yoon Kim, Justin Chiu, Demi Guo, and Alexander M. Rush. Latent alignment and  
 341 variational attention. In *Proceedings of Advances in Neural Information Processing Systems*, pp.  
 342 9735–9747, 2018.
- 343 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):  
 344 1735–1780, 1997.
- 345 Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and  
 346 artificial neural networks for natural scene text recognition. In *Proceedings of NIPS Workshop on  
 347 Deep Learning*, 2014.
- 348 Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer  
 349 networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (eds.), *Proceed-  
 350 ings of Advances in Neural Information Processing Systems*, pp. 2017–2025, 2015.
- 351 Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild  
 352 with convolutional neural networks. *International Journal of Computer Vision*, 116:1–20, 2016.
- 353 Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Bigorda,  
 354 Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazàn Almazàn, and Lluís Pere  
 355 de las Heras. ICDAR 2013 robust reading competition. In *Proceedings of the International  
 356 Conference on Document Analysis and Recognition*, pp. 1484–1493, 2013.

- 357 Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov,  
358 Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian  
359 Lu, Faisal Shafait, Seiichi Uchida, and Ernest Valveny. ICDAR 2015 competition on robust  
360 reading. In *Proceedings of the International Conference on Document Analysis and Recognition*,  
361 pp. 1156–1160, 2015.
- 362 Dieterich Lawson, Chung-Cheng Chiu, George Tucker, Colin Raffel, Kevin Swersky, and Navdeep  
363 Jaitly. Learning hard alignments with variational inference. In *Proceedings of the International  
364 Conference on Acoustics, Speech and Signal Processing*, pp. 5799–5803, 2018.
- 365 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
366 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 367 Wei Liu, Chaofeng Chen, Kwan-Yee K. Wong, Zhizhong Su, and Junyu Han. STAR-Net: A spatial  
368 attention residue network for scene text recognition. In Edwin R. Hancock Richard C. Wilson and  
369 William A. P. Smith (eds.), *Proceedings of British Machine Vision Conference*, pp. 43.1–43.13,  
370 2016.
- 371 S.M. Lucas, A. Panaretos, L. Sosa, A. Tang, S. Wong, and R. Young. ICDAR 2003 robust read-  
372 ing competitions. In *Proceedings of the International Conference on Document Analysis and  
373 Recognition*, pp. 682–687, 2003.
- 374 Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-  
375 based neural machine translation. In *Proceedings of the International Conference on Empirical  
376 Methods in Natural Language Processing*, pp. 1412–1421, 2015.
- 377 Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. What to talk about and how? Selective  
378 generation using LSTMs with coarse-to-fine alignment. In *Proceedings of the Conference of the  
379 North American Chapter of the Association for Computational Linguistics: Human Language  
380 Technologies*, pp. 720–730, 2016.
- 381 Anand Mishra, Karteek Alahari, and C. V. Jawahar. Scene text recognition using higher order  
382 language priors. In Richard Bowden, John P. Collomosse, and Krystian Mikolajczyk (eds.),  
383 *Proceedings of British Machine Vision Conference*, pp. 1–11, 2012.
- 384 Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. Recurrent models of visual  
385 attention. In *Proceedings of Advances in Neural Information Processing Systems*, pp. 2204–2212,  
386 2014.
- 387 Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gülçehre, and Bing Xiang. Abstractive  
388 text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of the SIGNLL  
389 Conference on Computational Natural Language Learning*, pp. 280–290, 2016.
- 390 Trung Quy Phan, Palaiahnakote Shivakumara, Shangxuan Tian, and Chew Lim Tan. Recognizing text  
391 with perspective distortion in natural scenes. In *Proceedings of IEEE Conference on International  
392 Conference on Computer Vision*, pp. 569–576, 2013.
- 393 Shiv Shankar and Sunita Sarawagi. Posterior attention models for sequence to sequence learning. In  
394 *Proceedings of the International Conference on Learning Representations*, 2019.
- 395 Shiv Shankar, Siddhant Garg, and Sunita Sarawagi. Surprisingly easy hard-attention for sequence  
396 to sequence learning. In *Proceedings of the International Conference on Empirical Methods in  
397 Natural Language Processing*, pp. 640–645, 2018.
- 398 Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based  
399 sequence recognition and its application to scene text recognition. *IEEE Transactions on Pattern  
400 Analysis and Machine Intelligence*, 39(11):2298–2304, 2017.
- 401 Baoguang Shi, Mingkun Yang, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. ASTER:  
402 An attentional scene text recognizer with flexible rectification. *IEEE Transactions on Pattern  
403 Analysis and Machine Intelligence*, 41(9):2035–2048, 2019.

- 404 Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks.  
405 In Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger (eds.), *Proceedings*  
406 *of Advances in Neural Information Processing Systems*, pp. 3104–3112, 2014.
- 407 Kai Wang, Boris Babenko, and Serge Belongie. End-to-end scene text recognition. In *Proceedings of*  
408 *IEEE Conference on International Conference on Computer Vision*, pp. 1457–1464, 2011.
- 409 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
410 learning. *Machine Learning*, 8(3–4):229–256, 1992.
- 411 Shijie Wu, Pamela Shapiro, and Ryan Cotterell. Hard non-monotonic attention for character-level  
412 transduction. In *Proceedings of the International Conference on Empirical Methods in Natural*  
413 *Language Processing*, pp. 4425–4438, 2018.
- 414 Kelvin Xu, Jimmy Lei Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov,  
415 Richard S. Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation  
416 with visual attention. In *Proceedings of the International Conference on Machine Learning*, pp.  
417 2048–2057, 2015.
- 418 Matthew D. Zeiler. Adadelta: An adaptive learning rate method, 2012.