

Text is Text, No Matter What: Unifying Text Recognition using Knowledge Distillation

Ayan Kumar Bhunia¹ Aneeshan Sain^{1,2} Pinaki Nath Chowdhury^{1,2} Yi-Zhe Song^{1,2}

¹SketchX, CVSSP, University of Surrey, United Kingdom.

²iFlyTek-Surrey Joint Research Centre on Artificial Intelligence.

{a.bhunias, p.chowdhury, a.sain, y.song}@surrey.ac.uk.

Abstract

Text recognition remains a fundamental and extensively researched topic in computer vision, largely owing to its wide array of commercial applications. The challenging nature of the very problem however dictated a fragmentation of research efforts: Scene Text Recognition (STR) that deals with text in everyday scenes, and Handwriting Text Recognition (HTR) that tackles hand-written text. In this paper, for the first time, we argue for their unification – we aim for a single model that can compete favourably with two separate state-of-the-art STR and HTR models. We first show that cross-utilisation of STR and HTR models trigger significant performance drops due to differences in their inherent challenges. We then tackle their union by introducing a knowledge distillation (KD) based framework. This however is non-trivial, largely due to the variable-length and sequential nature of text sequences, which renders off-the-shelf KD techniques that mostly works with global fixed length data inadequate. For that, we propose four distillation losses all of which are specifically designed to cope with the aforementioned unique characteristics of text recognition. Empirical evidence suggests that our proposed unified model performs on par with individual models, even surpassing them in certain cases. Ablative studies demonstrate that naive baselines such as a two-stage framework, multi-task and domain adaption/generalisation alternatives do not work as well, further authenticating our design.

1. Introduction

Text recognition has been studied extensively in the past two decades [37], mostly due to its potential in commercial applications. Following the advent of deep learning, great progress [4, 35, 57, 63, 5, 8, 7] has been made in recognition accuracy on different publicly available benchmark datasets [41, 58, 30, 39]. Beyond supervised text recognition, very recent attempts have been made that utilise synthetic training data via domain adaptation [67], learn optimal augmen-

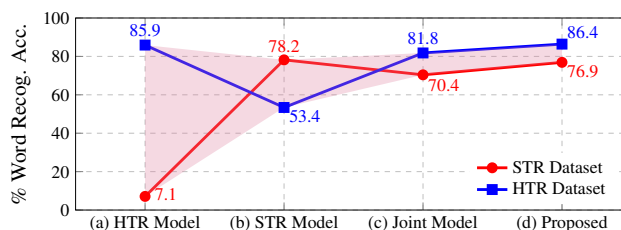


Figure 1. Despite performing well for scene images (IAM [39]), a model trained on HTR datasets (a), performs poorly in STR scenarios (ICDAR-2015 [30]) and vice-versa (b). Although jointly training a model (c) using both STR and HTR datasets helps improve the disparity between the datasets, the gap still remains far behind the specialist models. Our KD based proposed method leads to performance at par or even better than individual models.

tation strategy [38, 6], couple with visual question answering [10], and withhold adversarial attacks [60].

Albeit with great strides made, the field of text recognition remains fragmented, with one side focusing on Scene Text Recognition (STR) [30], and the other on Handwriting Text Recognition (HTR) [39]. This however is not surprising given the differences in the inherent challenges found in each respective problem: STR studies text in scene images posing challenges like complex backgrounds, blur, artefacts, uncontrolled illumination [63], whereas HTR tackles handwritten texts where the main challenge lies with the free-flow nature of writing [6] of different individuals. As a result, utilising models trained for STR on HTR (and vice versa) straightforwardly would trigger a significant performance drop (see Figure 1). This leads to our motivation – how to design a unified text recognition model that works ubiquitously across both scenarios.

While there is no existing work addressing this issue, one might naively think of training a single text recognition network using training data from *both* STR and HTR datasets. However, for the apparent issues of large domain gap and model capacity limitation [54], while the jointly trained model reduces the performance gap between HTR and STR datasets, it still lags significantly behind individual

specialised models. Another solution is to include a classification network prior to specialised STR and HTR models (i.e., a two-stage network). During evaluation, the classifier decides if an input belongs to scene or handwritten text, followed by choosing an appropriate model for downstream recognition. Yet, this solution has two downsides: a) classification network will incur additional computational cost and extra memory consumption to store all three neural networks. b) cascaded connection of the classifier and text recognition models will compound cumulative errors.

In this work, we introduce a *knowledge distillation* (KD) [22, 49] based framework to unify individual STR and HTR models into a *single* multi-scenario model. Our design at a high-level, does not deviate much from a conventional KD setting where a learnable student model tries to mimic the behaviour of a pre-trained teacher. We first train both STR and HTR models separately using their respective training data. Next, each individual model takes turns to act as a teacher in the distillation process, to train a single unified student model. It is this transfer of knowledge captured by specialised teachers into a single model, that leads to our superior performance in contrast to training a single model using joint STR and HTR datasets (see Figure 1).

Making such a design (KD) to work with text recognition is however non-trivial. The difficulty mainly arises from the variable-length and sequential natures of text images – each consists of a sequence of different number of individual characters. Hence, employing off-the-shelf KD methods [49] that aim at matching output probabilities and/or hidden representations between pre-trained teacher and learnable student model, which are used for global fixed length data, may not be sufficient to transfer knowledge at local character level. We thus propose *three* additional distillation losses to tackle the unique characteristics of text recognition.

More specifically, we first impose a *character aligned hint loss*. This encourages the student to mimic character-specific hidden representations of specialised teacher over the varying sequence of characters in a text image. Next, an *attention distillation loss* is further imposed over the attention map obtained at every step of character decoding process by an attentional decoder. This compliments the character localised hint-loss, as attention-maps capture rich and diverse contextual information emphasising on localised regions [23]. Besides localised character level information, capturing long-range non-local dependencies among the sequential characters is of critical importance, especially for an auto-regressive attentional decoder framework [34]. Accordingly we propose an *affinity distillation loss* as our third loss, to capture the interactions between every pair of positions of the variable character length sequence, and guide the unified student model to emulate the affinity matrix of the specialised teachers. Finally, we also make use of state-of-the-art *logit distillation loss* to work with our three pro-

posed losses. It aims at matching output probabilities of student network over the character vocabulary, with that of pre-trained teachers.

Our main contributions can be summarised as follows: (a) We design a practically feasible *unified* text recognition setting that asks a single model to perform equally well across both HTR and STR scenarios. (b) We introduce a novel knowledge distillation paradigm where an unified student model learns from two pre-trained teacher models specialised for STR and HTR. (c) We design three additional distillation losses to specifically tackle the variable-length and sequential nature of text data. (d) Extensive experiments coupled with ablative studies on public datasets, demonstrate the superiority of our framework.

2. Related Works

Text Recognition: With the inception of deep learning, Jaderberg *et al.* [27, 26] introduced a dictionary-based text recognition framework employing deep networks. Alternatively, Poznanski *et al.* [44] addressed the added difficulty in HTR by using a CNN to estimate an n-gram frequency profile. Later on, connectionist temporal classification (CTC) layer [17] made end-to-end sequence discriminative learning possible. Subsequently, CTC module was replaced by attention-based decoding mechanism [33, 51] that encapsulates language modeling, weakly supervised character detection and character recognition under a single model. Needless to say attentional decoder became the state-of-the-art paradigm for text recognition for both scene text [35, 63, 61, 66] and handwriting [6, 38, 59, 67]. Different incremental propositions [5, 8, 7] have been made like, improving the rectification module [66, 61], designing multi-directional convolutional feature extractor [12], improving attention mechanism [11, 34] and stacking multiple BLSTM layer for better context modelling [35].

Besides improving word recognition accuracy, some works have focused on improving performance in low data regime by designing adversarial feature deformation module [6], and learning optimal augmentation strategy [38], towards handling adversarial attack [60] for text recognition. Zhang *et al.* [67] introduced unsupervised domain adaptation to deal with images from new scenarios, which however definitely demands a fine-tuning step to specialise in new domain incurring additional server costs. On the contrary, we focus on unifying a single model capable of performing consistently well across both HTR and STR images.

Knowledge Distillation: Earlier, knowledge distillation (KD) was motivated towards training smaller student models from larger teacher models for cost-effective deployment. Caruana and his collaborators [1] pioneered in this direction, by using mean square error with the output *logits* of deeper model to train a shallower one. The seminal work by Hinton *et al.* [22] introduced *softer probability*

distribution over classes by a temperature controlled softmax layer for training smaller student models. Furthermore, Romero *et al.* [48] employed features learned by the teacher in the intermediate layers, to act as a hint for student’s learning. Later works explored different ideas like mimicking *attention maps* [64] from powerful teacher, transferring *neuron selectivity* pattern [24] by minimising Maximum Mean Discrepancy (MMD) metric, *gramian matrices* [62] for faster knowledge transfer, multiple *teacher assistants* [40] for step-wise knowledge distillation and so on. In addition to classification setup, KD has been used in object detection [14], semantic segmentation [21], depth-estimation [43], pose estimation [42], lane detection [23], neural machine translation [54] and so forth. Vongkulbhisal *et al.* [56] proposed a methodology of *unifying heterogeneous classifiers* having different label set, into a single unified classifier. In addition to obtaining smaller fast-to-execute model, using KD in *self-distillation* [3] improves performance of student having identical architecture like teacher. Keeping with self-distillation [3], our teacher networks and trainable student share exactly same architecture, but our motivation lies towards obtaining an unified student model from two pre-trained specialised teachers.

Unifying models: A unified model bestows several benefits compared to specialised individual models such as lower annotation and deployment cost as unlike it’s counterpart, unified models need not grow linearly with increasing domains [46] or tasks [65] while simultaneously cherishing the benefits of shared supervision. Towards embracing the philosophy of general AI, where the goal is to develop a single model handling multiple purposes, attempts have been made towards solving multiple tasks [28, 32, 65] via *multi-task learning*, working over multiple domains [9, 46], and employing *universal adversarial attack* [36]. While unsupervised *domain adaptation* [55] still needs fine-tuning over target domain images, *domain generalisation* [15] aims to extract domain invariant features, eliminating the need of post-updating step. In NLP community, handling multiple language pairs in one model via multi-lingual neural-machine-translation [18, 54], has been a popular research direction in the last few years. Albeit all these text recognition and *model unifying* approaches are extensively studied topics, we introduce an entirely new aspect of text recognition by unifying STR and HTR scenarios into a single model having significant commercial advantage.

3. Methodology

Overview: Our objective is to design a single unified model working both for STR (S) and HTR (H) word images. In this context, we have access to labelled STR datasets $\mathcal{D}_S = \{(I_s, Y_s) \in \mathcal{I}_s \times \mathcal{Y}_s\}$, as well as labelled HTR datasets $\mathcal{D}_H = \{(I_h, Y_h) \in \mathcal{I}_h \times \mathcal{Y}_h\}$. Here, I denotes word image from respective domain with label $Y = \{y_1, y_2, \dots, y_K\}$,

and K denotes the variable length of ground-truth characters. We first train two individual text-recognition models using \mathcal{D}_S and \mathcal{D}_H independently. Thereafter, a single unified model is obtained from two domain specific teacher via knowledge distillation.

3.1. Baseline Text Recognition Model

Given an image I , text recognition model \mathcal{R} tries to predict the machine readable character sequence Y . Out of the two state-of-the-art choices dealing with irregular texts, we adopt 2-D attention that localises individual characters in a weakly supervised way, over complicated rectification network [61]. Our text recognition model consists of three components: (a) a backbone convolutional feature extractor [52], (b) a RNN decoder predicting the characters autoregressively one at each time-step, (c) a 2D attentional block.

Let the extracted convolutional feature map be $\mathcal{F} \in \mathbb{R}^{h' \times w' \times d}$, where h' , w' and d signify height, width and number of channels. Every d dimensional feature at $\mathcal{F}_{i,j}$ encodes a particular local image region based on the receptive fields. At every time step t , the decoder RNN predicts an output character or end-of-sequence (EOS) y_t based on three factors: a) previous internal state s_{t-1} of decoder RNN, (b) the character y_{t-1} predicted in the last step, and (c) a glimpse vector g_t representing the most relevant part of \mathcal{F} for predicting y_t . To obtain g_t , previous hidden state s_{t-1} acts as a query to discover the attentive regions as follows:

$$J = \tanh(W_F \mathcal{F}_{i,j} + W_B \circledast \mathcal{F} + W_s s_{t-1})$$

$$\alpha_{i,j} = \text{softmax}(W_a^T J_{i,j}) \quad (1)$$

$$g_t = \sum_{i,j} \alpha_{i,j} \cdot \mathcal{F}_{i,j} \quad i = [1, \dots, h'], \quad j = [1, \dots, w'] \quad (2)$$

where, W_F , W_s , W_a are the learnable weights. Calculating the attention weight $\alpha_{i,j}$ at every spatial position (i, j) , we employ a convolution operation “ \circledast ” with 3×3 kernel W_B to consider the neighbourhood information in 2D attention mechanism. There exists $\alpha_t \in \mathbb{R}^{h' \times w'}$ corresponding to every time step of decoding, however t is dropped in Eqn. 1 and 2 for notational brevity. The current hidden state S_t is updated by: $(o_t, s_t) = \text{RNN}(s_{t-1}; [E(y_{t-1}), g_t])$, where $E(\cdot)$ is character embedding layer with embedding dimension \mathbb{R}^{128} , and $[\cdot]$ signifies a concatenation operation. Finally, \tilde{y}_t is predicted as: $p(\tilde{y}_t) = \text{softmax}(W_o o_t + b_o)$ with learnable parameters W_o and b_o . This model is trained end-to-end using cross-entropy loss $\mathcal{H}(\cdot, \cdot)$ summed over the ground-truth sequence $Y = \{y_1, y_2, \dots, y_K\}$, where y_t is one-hot encoded vector of size $\mathbb{R}^{|V|}$, and $|V|$ is the character vocabulary size.

$$\mathcal{L}_C = \sum_{t=1}^K \mathcal{H}(y_t, \tilde{y}_t) = - \sum_{t=1}^K \sum_{i=1}^{|V|} y_{t,i} \log p(\tilde{y}_{t,i}) \quad (3)$$

3.2. Basics: Knowledge Distillation

Initially, knowledge distillation (KD) [22] was proposed for classification tasks to learn a smaller student model

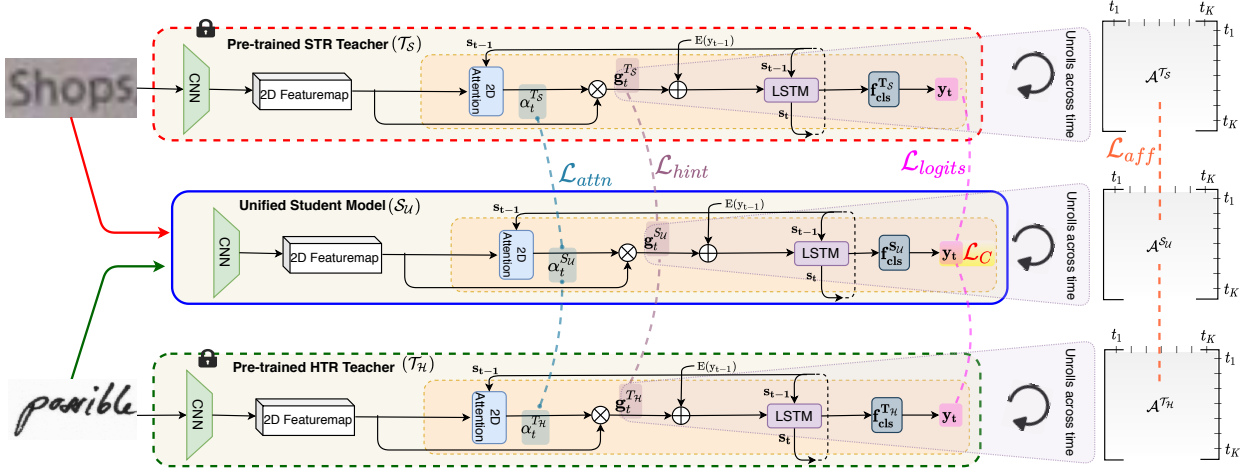


Figure 2. STR and HTR models, pre-trained using respective images, are used as *teachers* to train a unified *student* model via knowledge distillation, with four distillation losses and a cross-entropy loss (\mathcal{L}_C). The t^{th} time-step for decoding is shown, which unrolls across time.

by mimicking the output of a pre-trained teacher. Given a particular data, let the output from pre-trained teacher be $\tilde{y}_t^T = \text{softmax}(l_t^T)$ and that of learnable student be $\tilde{y}_t^S = \text{softmax}(l_t^S)$, where l_t is pre-softmax logits from respective models. Temperature (τ) normalised softmax is used to soften the output so that more information regarding inter-class similarity could be captured for training. Therefore, given $\tilde{y}_{\tau,t}^T = \text{softmax}(\frac{l_t^T}{\tau})$, $\tilde{y}_{\tau,t}^S = \text{softmax}(\frac{l_t^S}{\tau})$ and ground-truth y , the student network is trained to optimise the following loss function:

$$\mathcal{L}_{KD} = \sum_{t=1}^K \mathcal{H}(y_t, \tilde{y}_t^S) + \lambda \sum_{t=1}^K \mathcal{H}(\tilde{y}_{\tau,t}^T, \tilde{y}_{\tau,t}^S) \quad (4)$$

where λ is a hyper-parameter balancing the two terms, and the first term signifies traditional cross-entropy loss between the output of student network and ground-truth labels, whereas the second term encourages the student to learn from softened output of teacher.

Adopting basic KD formulation however is unsuitable for our purpose. Firstly, text recognition dealing with varied-length sequence recognition requires distilling local fine-grained character information. Additionally, there exists a sequential dependency among the predicted characters due to auto-regressive nature of attentional decoder, thus requiring a global consistency criteria during distillation process. (b) While training teacher and student usually involves same (single domain) dataset, we here have two separate domains, STR and HTR, which thus needs to deal with larger domain gap and data coming from two separate domains.

3.3. Unifying Text Recognition Models

Overview: We propose a knowledge distillation method for sequential text images to unify both scene-text and handwriting recognition process into a *single* model. Compared to traditional knowledge distillation, we have *two* pre-trained teacher networks $T \in \{T_S, T_H\}$, where T_S is a spe-

cialised model trained from \mathcal{S} cene text images \mathcal{D}_S , and T_H from \mathcal{H} andwritten text images \mathcal{D}_H . Given these pretrained teachers, we aim to learn a single \mathcal{U} nified \mathcal{S} tudent model S_U by *four* distillation losses tailored for sequential recognition task, along with typical cross-entropy loss. T_S , T_H and S_U all have identical architecture to text recognition network $\mathcal{R}(\cdot)$. Directly training a single model by including images from both the STR and HTR datasets leads to sub-optimal performance due to limited model capacity and large domain-gap. In contrast, training of *specialised* models might assist to extract underlying structure from respective data, which can *then* be distilled into a unified student network with guidance from the specialised teachers.

We have two pre-trained teachers $T \in \{T_S, T_H\}$, with images coming from two different domains $I \in \{I_s, I_h\}$. In order to train a student network S_U , we will get one loss instance using STR pre-trained teacher and respective dataset (T_S, I_s) , and similarly another loss term for HTR counterpart (T_H, I_h) . We describe the loss functions using generalised notation (T, I) which basically has two elements, (T_S, I_s) and (T_H, I_h) respectively. Thus mathematically, $(T, I) : \{(T_S, I_s), (T_H, I_h)\}$. Please refer to Figure 2.

Logits' Distillation Loss: We extend the traditional knowledge distillation loss for our sequence recognition task by aggregating cross-entropy loss over the sequence. Given an image I , let the temperature normalised softmax output from a particular pre-trained teacher and trainable student be $\tilde{y}_t^T(I)$ and $\tilde{y}_t^{S_U}(I)$ at a particular time-step t . We ignore τ of Eqn. 4 here for notational brevity. We call this logits' distillation loss and define it as:

$$\mathcal{L}_{\text{logits}}(T, I) = \sum_{t=1}^K \mathcal{H}(\tilde{y}_t^T(I), \tilde{y}_t^{S_U}(I)) \quad (5)$$

where, $(T, I) : \{(T_S, I_s), (T_H, I_h)\}$. We get two of such logits' distillation loss with respect to STR and HTR datasets (and pre-trained teachers) respectively.

Character Localised Hint Loss: The fact that intermediate features learned by the teacher could further act as a ‘hint’ in the distillation process, was shown by Romero *et al.* [48]. Being a sequence recognition task however, text recognition needs to deal with variable length of sequence, with each character having variable width within itself. While predicting every character, attention based decoder focuses on specific regions of convolutional feature-map. In order to circumvent the discrepancy due to variable character-width, we perform feature distillation loss at the space of character localised visual feature, termed as *glimpse vector* (see Eqn. 2) instead of global convolutional feature-map. This provides the teacher’s supervision at local level. As our student shares the same architecture identical to the pre-trained teachers, we do not need any parametric transformation layer to match the feature-space between them. The character localised hint loss is given by:

$$\mathcal{L}_{\text{hint}}(\mathbb{T}, \mathbb{I}) = \sum_{t=1}^K \left\| g_t^{\mathbb{T}}(\mathbb{I}) - g_t^{\mathbb{S}u}(\mathbb{I}) \right\|_2 \quad (6)$$

where, $(\mathbb{T}, \mathbb{I}) : \{(\mathbb{T}_S, \mathbb{I}_s), (\mathbb{T}_S, \mathbb{I}_h)\}$. Given an input image \mathbb{I} , $g_t^{\mathbb{T}}(\mathbb{I})$ and $g_t^{\mathbb{S}u}(\mathbb{I})$ are glimpse vector of size \mathbb{R}^d at t -th times step from a particular pre-trained teacher and trainable student.

Attention Distillation Loss: While Character Localised Hint Loss aids in enriching the localised information (i.e. absolute information in the cropped region roughly enclosing the specific character), computed attention map (see Eqn 2) brings *contextual information* giving insights about which region is *relatively* more important than the others, over a convolutional feature map. Unlike attentional distillation, logits’ distillation does not explicitly take into account the degree of influence each pixel has on model prediction, thus making the attention map computed at every step a complementary source of information [64] to learn from the student. Furthermore, HTR usually shows overlapping characters, which however rarely occurs in STR. Thus the student must learn the proper ‘look-back’ (attention) mechanism from specialised teachers. Let $\alpha_t^{\mathbb{T}}(\mathbb{I})$ and $\alpha_t^{\mathbb{S}u}(\mathbb{I})$ represent the attention map from respective teacher and learnable student at t -th time step, both having size $\mathbb{R}^{h' \times w'}$ for a given an input image \mathbb{I} . Considering $(\mathbb{T}, \mathbb{I}) : \{(\mathbb{T}_S, \mathbb{I}_s), (\mathbb{T}_H, \mathbb{I}_h)\}$, the attention distillation loss is computed as follows:

$$\mathcal{L}_{\text{attn}}(\mathbb{T}, \mathbb{I}) = \sum_{t=1}^K \left\| \alpha_t^{\mathbb{T}}(\mathbb{I}) - \alpha_t^{\mathbb{S}u}(\mathbb{I}) \right\|_2 \quad (7)$$

Affinity Distillation Loss: Attention based decoder encapsulates an implicit language model within itself, and the information of previously predicted characters flows through its hidden state. While previous character localised hint loss and attention distillation loss mostly contribute to information distillation at local level, with the later (attention) additionally contributing towards the contextual information,

we need a global consistency loss to handle the long-range dependency among the characters. Thus we introduce an affinity distillation loss to model long-range non-local dependencies from the specialised teachers. Given character aligned features $\{g_1, g_2, \dots, g_K\}$ for a given image, the affinity matrix capturing the pair-wise correlation between every pair of characters is computed as:

$$\mathcal{A}_{i,j} = \frac{1}{K \times K} \cdot \frac{g_i}{\|g_i\|_2} \cdot \frac{g_j}{\|g_j\|_2} \quad (8)$$

where, $\mathcal{A} \in \mathbb{R}^{K \times K}$ represents the affinity matrix for a word image having character sequence length K . We use l_2 loss to match the affinity matrix of specialised teacher $\mathcal{A}^{\mathbb{T}}(\mathbb{I})$ and that of learnable student $\mathcal{A}^{\mathbb{S}u}(\mathbb{I})$:

$$\mathcal{L}_{\text{aff}}(\mathbb{T}, \mathbb{I}) = \left\| \mathcal{A}^{\mathbb{T}}(\mathbb{I}) - \mathcal{A}^{\mathbb{S}u}(\mathbb{I}) \right\|_2 \quad (9)$$

Optimisation Procedure: Apart from the four distillation loss in order to learn from the specialised teacher, the unified student model S_u is trained from ground-truth label for image $\mathbb{I} \in \{I_s, I_h\}$ using typical cross-entropy loss (see Eqn. 3). Thus, given $(\mathbb{T}, \mathbb{I}) : \{(\mathbb{T}_S, \mathbb{I}_s), (\mathbb{T}_H, \mathbb{I}_h)\}$, the overall training objective for student becomes:

$$\mathcal{L}_{\text{all}} = \sum_{\forall(\mathbb{T}, \mathbb{I})} \left(\mathcal{L}_C(\mathbb{I}) + \lambda_1 \cdot \mathcal{L}_{\text{logits}}(\mathbb{T}, \mathbb{I}) + \lambda_2 \cdot \mathcal{L}_{\text{attn}}(\mathbb{T}, \mathbb{I}) + \lambda_3 \cdot \mathcal{L}_{\text{hint}}(\mathbb{T}, \mathbb{I}) + \lambda_4 \cdot \mathcal{L}_{\text{aff}}(\mathbb{T}, \mathbb{I}) \right) \quad (10)$$

Due to difference in complexity of the task of HTR and STR and their respective training data size, we observe a tendency to learn a biased model that over-fits on either STR or HTR dataset. To alleviate this, we employ a conditional distillation mechanism that stabilise training by deciding in what proportion to learn from two different individual specialised teacher that results in a unified student model performing ubiquitously over both STR and HTR scenarios.

4. Experiments

Datasets: Training paradigm for STR involves using large synthetic datasets such as Synth90k [25] and SynthText [20] with 8 and 6 million images respectively, and evaluating (without fine-tuning) on real images such as: **IIT5K-Words**, **Street View Text (SVT)**, **SVT-Perspective (SVT-P)**, **ICDAR 2013 (IC13)**, **ICDAR 2015 (IC15)**, and **CUTE80**. IIT5-K Words [41] has 5000 cropped words from Google image search. SVT [58] hosts 647 images collected from Google Street View where most images are blurry, noisy and have low resolution. SVT-P [45] has 639 word images also taken from Google Street view but with side-view snapshots resulting in severe perspective distortions. ICD13 [31] contains 848 cropped word patches with mostly regular images unlike IC15 [30] which has 2077 word images that are irregular i.e. oriented, perspective or curved. Unlike others, CUTE80 [47] dataset contains high resolution image but have curved text. In context of HTR, we follow the evaluation setup described in [6] on two large standard datasets viz, **IAM** [39] (1,15,320 words) and **RIMES** (66,982 words).

Algorithm 1 Training algorithm of the proposed framework

```
1: Input: Dataset:  $\{\mathcal{D}_S, \mathcal{D}_H\}$ ; Teachers:  $\{T_S, T_H\}$ ;  
   Learning rate:  $\eta$ ; Total Training Steps:  $\mathcal{T}$ , distil check:  
    $\mathcal{T}'$ ; Accuracy metric:  $Acc$ ; distil acc. thresh.  $\omega \geq 1$   
2: Initialise: Unified Student Model:  $\mathcal{S}_U$ , params:  $\theta^{Su}$ ;  
   Step:  $t = 1$ ; Gradient:  $g$ ; Flags:  $\{f^S, f^H\}$  are True  
3: while  $t \leq \mathcal{T}$  do  
4:    $g = 0$   
5:   Get:  $(I_s, Y_s) \in \mathcal{D}_S^{train}; (I_h, Y_h) \in \mathcal{D}_H^{train}$   
6:    $g += \partial(\mathcal{L}_C(I_s) + \mathcal{L}_C(I_h))/\partial\theta^{Su}$   $\triangleright$  see eq. 3  
7:   for each  $\mathcal{L}_{KD}$  in  $\mathcal{L}_{all} - \{\mathcal{L}_C\}$  do  $\triangleright$  see eq. 10  
8:     if  $f^S$  then  $g += \partial\mathcal{L}_{KD}(T_S, I_s)/\partial\theta^{Su}$   
9:     if  $f^H$  then  $g += \partial\mathcal{L}_{KD}(T_H, I_h)/\partial\theta^{Su}$   
10:  end for  
11:  Update  $\theta^{Su}$  :  $\theta^{Su} = \theta^{Su} - \eta * g$   
12:  if  $t\% \mathcal{T}' == 0$  then  $\triangleright$  conditional distillation  
13:     $\mathcal{L} = \mathcal{L}_{all} - \{\mathcal{L}_C\}$   
14:     $\{\mathcal{I}_s^{val}, \mathcal{Y}_s^{val}\} = \mathcal{D}_S^{val}; \{\mathcal{I}_h^{val}, \mathcal{Y}_h^{val}\} = \mathcal{D}_H^{val}$   
15:    if  $\mathcal{L}(T_S, \mathcal{I}_s) > \omega \cdot \mathcal{L}(T_H, \mathcal{I}_h)$  then  $f^H = False$   
16:    else  $f^H = True$   
17:    if  $\mathcal{L}(T_H, \mathcal{I}_h) > \omega \cdot \mathcal{L}(T_S, \mathcal{I}_s)$  then  $f^S = False$   
18:    else  $f^S = True$   
19:  end if  
20:   $t = t + 1$   
21: end while
```

Implementation Details: We use a 31-layer CNN backbone feature extractor [34] without any pre-training. The input image is resized to 48×160 following [34]. We first pre-train the specialised HTR and STR model at a time. For STR, we use Synth90k [25] and SynthText [20] dataset together, and respective training set is used for experiments on IAM and RIMES dataset individually. We use Adam optimiser with initial learning rate of 0.001 and batch size of 32 for both specialised teacher pre-training, and distillation based unified student model training. Decay rate of 0.9 is applied after every 10^4 iteration till the learning rate drops to 10^{-5} . During conditional distillation (Algorithm 1), loss is compared over the validation set with $\omega = 1.05$. We set $\lambda_1, \lambda_2, \lambda_3$, and λ_4 as 0.5, 5, 1 and 1 respectively. We implement the network and its training paradigm using PyTorch trained in a 11 GB NVIDIA RTX-2080-Ti GPU.

Evaluation Protocol: To better understand the challenges of unifying STR and HTR, and recognise contribution of each alternative training paradigm we evaluate as follows: (i) we first evaluate the pre-trained teacher models on the dataset for what it has been trained for, e.g. \mathcal{T}_S on testing set of STR dataset, and \mathcal{T}_H on that of HTR dataset. (ii) Next, we evaluate on the alternative dataset for pre-trained teacher model and see how the performance drops in cross-dataset scenarios, e.g. \mathcal{T}_S on testing set of HTR dataset, and vice-versa. (iii) Finally, we evaluate the unified student model \mathcal{S}_U on both STR and HTR datasets to verify if a single model can perform ubiquitously for both scenarios.

4.1. Competitors

To the best of our knowledge, there has been no prior work dealing with the objective of unifying STR and HTR models into a single model. Thus, we design a few strong baselines based on the existing literature by our own. (i) **Multi-Task-Training:** This is a naive *frustratingly easy* training paradigm [13] where samples belonging to both STR and HTR datasets are used to train a single network guided by cross-entropy loss. Since STR has overwhelmingly large synthetic training samples [25, 20] compared to HTR dataset [39], we use weighted random sampling (variant-I) to balance training data. Conversely, we randomly sample a subset from STR dataset (variant-II) to forcefully make the number of training images similar for HTR and STR datasets in order to validate the utility of conditional distillation. In variant-III, we treat HTR and STR character units as different classes, thus extending it to N-class to 2N class classification at each time step. (ii) **DA-Corr-Unsup:** An obvious alternative is to try out any domain adaptation method introduced for sequence recognition task. Zhang *et al.* [67] proposed unsupervised domain adaptation (DA) technique for text images. We start by training a model on either STR (or HTR) images that acts as our source domain, followed by unsupervised adaptation to the target HTR (or STR) images – thus we have two version of this model STR model adapted to HTR as (HTR \mapsto STR), and (STR \mapsto HTR). Second-order statistics-correlation distance [53] is used to align feature distribution from two domain. (iii) **DA-Corr-Sup:** As we have the access to both labelled STR and HTR datasets, we further extend the unsupervised DA setup of Zhang *et al.* [67] by considering target domain to be annotated, allowing supervised DA. Cross-entropy loss is minimised for both source and target domain in association to second-order statistics-correlation between both STR and HTR domains. (iv) **DA-Adv-Unsup:** We further adopt a recent work by Kang *et al.* [29] employing adversarial learning for unsupervised domain adaptation for text recognition. Here, the setup remains same as DA-Corr-Unsup having two versions as (HTR \mapsto STR) and (STR \mapsto HTR), but domain adaptation tackled through a discriminator with a preceding gradient-reversal layer. (v) **DA-Adv-Sup:** This is again a similar adaptation of [29] following supervised DA which minimise Cross-Entropy and domain classification loss for both STR and HTR. (vi) **DG-Training:** Another alternative way to address this problem could be to use Domain Generalisation (DG) training based on model agnostic meta-learning using episodic-training [16]. It involves using weighted (λ) summation [19] for gradient (over meta-train set) and meta-gradient (over meta-test split through inner loop update) to train our baseline text recognition model. The inner-loop update process consists of support set consisting images of either STR (or HTR) word images while the outer-loop up-

Table 1. Quantitative performance against various alternatives. **Competitors use combined STR+HTR datasets** in different setups: (a) Multi-Task (Joint) Training, (b) Unsupervised and Supervised Domain Adaptation (DA), (c) Domain Generalization (DG).

Methods	STR datasets						HTR dataset	
	IIT5-K	SVT	IC13	IC15	SVT-P	CUTE80	IAM	RIMES
Multi-Task-Training-(I)	86.1	83.6	87.2	70.4	77.8	79.4	81.8	86.2
Multi-Task-Training-(II)	35.4	34.5	36.3	29.1	32.1	32.5	81.9	85.9
Multi-Task-Training-(III)	83.2	80.5	84.1	67.1	74.1	76.3	77.9	82.3
DA-Adv-Unsup (STR → HTR)	82.6	80.1	84.2	66.8	74.2	75.8	58.7	64.1
DA-Adv-Unsup (HTR → STR)	16.6	12.9	15.4	12.1	12.7	13.4	78.1	82.4
DA-Adv-Sup	88.1	85.6	89.2	72.5	79.9	81.6	83.1	87.5
DA-Corr-Unsup (STR → HTR)	82.7	80.2	84.5	67.8	74.7	76.1	82.7	87.1
DA-Corr-Unsup (HTR → STR)	17.1	13.1	15.9	12.7	13.1	13.9	82.7	87.1
DA-Corr-Sup	88.3	85.8	89.4	72.7	80.1	81.8	83.2	87.6
DG-training	88.5	86.0	89.5	72.9	80.3	82.0	83.4	87.7
Proposed	92.3	89.9	93.3	76.9	84.4	86.3	86.4	90.6

Table 2. Quantitative comparison of our STR-only and HTR-only models, trained on STR and HTR datasets respectively, against state-of-the-arts. Our method uses STR-only and HTR-only as teachers during KD.

Methods	STR datasets				HTR dataset	
	IIT5-K	SVT	IC13	IC15	IAM	RIMES
Shi <i>et al.</i> [52]	93.4	93.6	91.8	76.1	–	–
Baek <i>et al.</i> [2]	87.9	87.5	92.3	71.8	–	–
Yu <i>et al.</i> [63]	94.8	91.5	95.5	82.7	–	–
Litman <i>et al.</i> [35]	93.7	92.7	93.9	82.2	–	–
Bhunia <i>et al.</i> [6]	–	–	–	–	82.81	88.53
STR-only Model	93.1	90.9	93.5	78.2	53.4	58.5
HTR-only Model	11.5	7.6	10.3	7.1	85.9	90.2
Joint STR-HTR Model	86.1	83.6	87.2	70.4	81.8	86.2
Proposed (Unified)	92.3	89.9	93.3	76.9	86.4	90.6

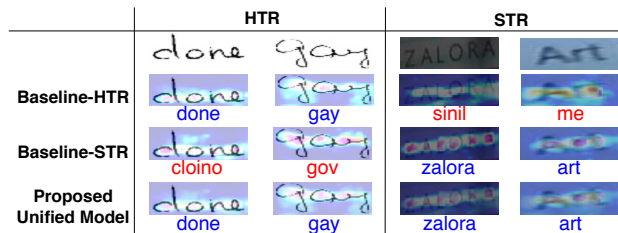


Figure 3. Illustrative examples with attention maps, and prediction (red → incorrect, blue → correct). While discrepancy exists for cross-dataset scenarios, attention-map from unified model is nearly consistent with that of respective specialised model.

date process is materialised using images from a different domain i.e. HTR (or STR). Such inner and outer-loop based optimisation strategy helps learn a model that aims to generalise well for both scenarios without further fine-tuning.

4.2. Performance Analysis

From Table 2, it can be seen that while a model trained on HTR fails miserably when evaluated on STR datasets, training on STR followed by testing on HTR does not result in a similar collapse in performance. This indicates that although STR scenarios partially encompass domain specific HTR attributes, the reverse is not true. Interestingly, this is likely why there is a positive transfer for HTR datasets using *unified model* compared to HTR-only counterpart. Moreover, our KD based unifying approach for multi-scenario text recognition outperforms all other baselines by a significant margin. In particular, (i) For baselines designed for unification, we attribute the limitation of all three multitask-learning-training (also named as joint-training) variants to the reason that it does not consider the varying complexity of two different tasks during joint training. Instead, our pre-trained teacher models first discover the *specialised knowledge* from respective scenario. Given the specialised knowledge, our framework can encapsulate it into a single framework by balancing the learning via *conditional distillation* from two different data sources (see Figure 3). We outperform this joint-training (variant-I being the best performing competitor) baseline by a margin of almost 6 – 7% on every dataset. Limited performance

of variant-II validates the necessity and motivation of conditional distillation. (ii) The performance of unsupervised DA is limited by a significant margin while evaluating on both HTR and STR datasets. Starting from any source domain, it hardly gives any significant rise in target domain, rather the performance even decreases in the source domain after adaptation. An inevitable corollary of unsupervised DA is the lack of any guarantee that a model will retain information about source domain after successful adaptation to the target domain. (iii) The Domain Adaptation (DA) based pipelines suppress multitask-learning-training baseline while using supervised-labels from both the datasets, but lags behind us by 3.5 – 4.5% on an average. Even using supervised-labels from both the datasets, the learning process oscillates around discovering domain invariant representation, and ignores main objective of unification of two specialised knowledge available from labelled datasets. Furthermore, adversarial learning based DA [29] falls short compared to covariance based character-wise distribution alignment [67] for text recognition – this also supports our design of using distillation loss over glimpse vectors. (iv) Both [67] and [57] train a text recognition model on a source domain comprising of easily available synthetic images followed by unsupervised adaptation to target domain consisting of real world text images. While cost-effective training from synthetic-data is their major objective, we consider to have access to both the labelled datasets (which are readily available nowadays) to design an unified model working for both scenarios – making our work orthogonal to these two DA based pipelines. (v) The purpose of Domain Generalisation (DG) is to find a model robust to domain-shift, giving satisfactory performance without the need of further adaptation. While such technique play a key role in unseen data regime, given enough labelled data, a frustratingly-simpler [13] alternative – multi-task learning – also achieves similar performance gains. Given the labelled STR and HTR training data, we observe that although DG-training outperforms multi-task-training, it lags behind our proposed method by almost 4% due to unavailability of privilege information (Table 1). (vi) The diversity of vocabulary (words present

in the dataset) between STR and HTR scenarios forms an important limitation to achieve SOTA performance [57]. While nouns (‘stop’, ‘walk’) are observed in STR images (placard, road signs), verbs or adverbs (‘taking’, ‘giving’) are more prevalent in HTR. Our specialised knowledge discovery bridges this discrepancy via unification.

Table 3. Contribution (WRA) of each KD constraint with \mathcal{L}_C

$\mathcal{L}_{\text{logits}}$	$\mathcal{L}_{\text{attn}}$	$\mathcal{L}_{\text{hint}}$	\mathcal{L}_{aff}	IC15	IAM
-	-	-	-	70.4	81.8
-	-	-	-	75.3	84.9
✓	✓	-	-	75.7	85.3
✓	✓	✓	-	76.4	85.9
✓	✓	✓	✓	76.9	86.4

Table 4. Analysis of Time and Space complexities.

Methods	IC15	IAM	GFlops	Params.
M.T.T	70.4	81.8	0.67	19M
B.C.R	74.4	83.1	0.80	50M
KD-Res-12	74.2	83.9	0.38	16M
KD-Res-31	74.7	84.2	0.12	9M
Proposed	76.9	86.4	0.67	19M

4.3. Ablation Study:

[i] Competitiveness of our baseline: Our baseline text recognition model is loosely inspired from the work by Li *et al.* [34] that also uses 2D attention to locate the characters in weakly supervised manner even from irregular text images for recognition. An alternative is to use a two-stage framework consisting of an *image rectification module* [52] followed by text recognition [2]. But as observed by Zhang *et al.* [67], although rectification based networks designed to handle spatial distortions lead to good performance in irregular STR datasets, it becomes a bottleneck for HTR tasks due to distortion caused by handwriting styles. Hence, for the purpose of unified text recognition, 2D attention mechanism provides a reasonable choice to bypass the rectification network in the text recognition system. Table 2 shows our baseline text recognition model to have a competitive performance in comparison to existing methods in both STR and HTR datasets. Moreover, we tried to replicate our KD based pipeline incorporating *image rectification module* on the top of [2], but performance gets limited to 75.9% and 85.5% on IC15 and IAM dataset, respectively. **[ii] Binary-Classifier based two-stage alternative:** Besides *Multi-Task-Training (M.T.T)*, another alternative is to use a binary-classifier (**B.C.R**) to classify between HTR and STR samples, then followed by selecting either STR or HTR model accordingly. While this achieves comparable performance with ours, it involves heavy computational expenses for maintaining three networks (2 specialised models + 1 classifier) together even while using simple ResNet18 as binary classifier – thus making it inefficient for online deployment. A thorough analysis on the computational aspect is shown in Table 4. **[iii] Significance of individual losses:** Among the four knowledge distillation losses ($\mathcal{L}_{\text{logits}}$, $\mathcal{L}_{\text{attn}}$, $\mathcal{L}_{\text{hint}}$, \mathcal{L}_{aff}), we use one of these distillation constraints along with \mathcal{L}_C to understand their individual relative contribution. Table 3 shows $\mathcal{L}_{\text{hint}}$ to have the greatest impact among others, increasing accuracy on IC15 (IAM) by 5.1% (3.3%), followed by $\mathcal{L}_{\text{logits}}$ resulting in an increase of 4.9% (3.1%), \mathcal{L}_{aff} by 4.8% (3.0%) and $\mathcal{L}_{\text{attn}}$ by 4.3% (2.6%). **[iv] Significance of conditional distillation:**

Besides the wide difference in training data size, the complexity of the task of HTR and STR is different. A simple multi-task-training often over-fits on either STR or HTR dataset – leading to sub-optimal performance of the unified student model. Thus, conditional distillation not only stabilises training, but also helps the student model to decide in what proportion to learn from two different individual specialised teachers, so that the unified model performs ubiquitously over both STR and HTR scenarios. Without conditional distillation, the performance is reduced by 2.5% and 0.4% on IC15 and IAM datasets, respectively. The hyperparameter ω controlling the conditional distillation process is varied at 1.01, 1.03, 1.05, 1.07, 1.10, and results on IC15 (IAM) are 76.8% (86.3%), 76.9% (86.3%), 76.9% (86.4%), 76.8% (86.4%), 76.8% (86.4%). **[vi] Hint Loss location:** While hint-based training leads to performance enhancements, the location of feature distillation loss is debatable based on the model’s architecture. Thus, we employ $\mathcal{L}_{\text{hint}}$ on: (a) CNN features \mathcal{F} and (b) hidden state s_t of attentional decoder. Using $\mathcal{L}_{\text{hint}}$ on \mathcal{F} lead to a performance improvement of 3.8% (2.2%) while on s_t results in 4.6% (2.5%) enhancement on IC15(IAM) datasets; both of which are lower as compared to $\mathcal{L}_{\text{hint}}$ on context vector g giving 5.1% (3.3%) improvement over the baseline model. **[vii] Reduce model size using KD:** Knowledge distillation is a generic method used to compress [22] any deep model regardless of the structural difference between teacher and student. Hence, we further check if our tailored KD method for attentional decoder based text recognition framework could be used off-the-shelf to reduce the model size of unified student. We replace our student model having 31-layer ResNet with just 12-layer (2+2+3+3+2) as KD-ResNet-12, and replace normal convolution by depth-wise convolution following MobileNetV2 architecture [50] to obtain KD-ResNet-31. The two resulting light-weight architectures give 74.2% (83.9%) and 74.7% (84.2%) accuracies in IC15 (IAM) datasets without much significant drop compared to our full version as shown in Table 4. This suggests that our framework could be widened further for model compression of text recognition model.

5. Conclusion

We put forth a novel perspective towards text recognition – unifying multi-scenario text recognition models. To this end we introduced a robust resource-economic online serving solution by proposing a knowledge distillation based framework employing four distillation losses to tackle the varying length of sequential text images. This helps us reduce the domain gap between scene and handwritten images while alleviating language diversity and model capacity limitations. The resulting unified model proves capable of handling both scenarios, performing at par with individual models, even surpassing them at times (e.g. in HTR).

References

- [1] Jimmy Ba and Rich Caruna. Do deep nets really need to be deep? In *NeurIPS*, 2014. 2
- [2] Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoo Yun, Seong Joon Oh, and Hwal-suk Lee. What is wrong with scene text recognition model comparisons? dataset and model analysis. In *ICCV*, 2019. 7, 8
- [3] Hessam Bagherinezhad, Maxwell Horton, Mohammad Rastegari, and Ali Farhadi. Label refinery: Improving imagenet classification through label progression. *arXiv preprint arXiv:1805.02641*, 2018. 3
- [4] Fan Bai, Zhanzhan Cheng, Yi Niu, Shiliang Pu, and Shuigeng Zhou. Edit probability for scene text recognition. In *CVPR*, 2018. 1
- [5] Ayan Kumar Bhunia, Pinaki Nath Chowdhury, Aneeshan Sain, and Yi-Zhe Song. Towards the unseen: Iterative text recognition by distilling from errors. In *ICCV*, 2021. 1, 2
- [6] Ayan Kumar Bhunia, Abhirup Das, Ankan Kumar Bhunia, Perla Sai Raj Kishore, and Partha Pratim Roy. Handwriting recognition in low-resource scripts using adversarial learning. In *CVPR*, 2019. 1, 2, 5, 7
- [7] Ayan Kumar Bhunia, Shuvojit Ghose, Amandeep Kumar, Pinaki Nath Chowdhury, Aneeshan Sain, and Yi-Zhe Song. Metaht: Towards writer-adaptive handwritten text recognition. In *CVPR*, 2021. 1, 2
- [8] Ayan Kumar Bhunia, Aneeshan Sain, Amandeep Kumar, Shuvojit Ghose, Pinaki Nath Chowdhury, and Yi-Zhe Song. Joint visual semantic reasoning: Multi-stage decoder for text recognition. In *ICCV*, 2021. 1, 2
- [9] Hakan Bilen and Andrea Vedaldi. Universal representations: The missing link between faces, text, planktons, and cat breeds. *arXiv preprint arXiv:1701.07275*, 2017. 3
- [10] A.F. Biten, R. Tito, A. Mafra, L. Gomez, M. Rusiñol, E. Valveny, C. Jawahar, and D. Karatzas. Scene text visual question answering. In *CVPR*, 2019. 1
- [11] Zhanzhan Cheng, Fan Bai, Yunlu Xu, Gang Zheng, Shiliang Pu, and Shuigeng Zhou. Focusing attention: Towards accurate text recognition in natural images. In *ICCV*, 2017. 2
- [12] Zhanzhan Cheng, Yangliu Xu, Fan Bai, Yi Niu, Shiliang Pu, and Shuigeng Zhou. Aon: Towards arbitrarily-oriented text recognition. In *CVPR*, 2018. 2
- [13] Hal Daumé III. Frustratingly easy domain adaptation. In *ACL*, 2007. 6, 7
- [14] Jiajun Deng, Yingwei Pan, Ting Yao, Zhou Wengang, Li Houqiang, and Tao Mei. Relation distillation networks for video object detection. In *ICCV*, 2019. 3
- [15] Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. In *NeurIPS*, 2019. 3
- [16] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 6
- [17] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *ICML*, 2006. 2
- [18] Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O. K. Li. Universal neural machine translation for extremely low resource languages. In *NAACL-HLT*, 2018. 3
- [19] Jianzhu Guo, Xiangyu Zhu, Chenxu Zhao, Dong Cao, Zhen Lei, and Stan Z Li. Learning meta face recognition in unseen domains. In *CVPR*, 2020. 6
- [20] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *CVPR*, 2016. 5, 6
- [21] Tong He, Chunhua Shen, Zhi Tian, Dong Gong, Changming Sun, and Youliang Yan. Knowledge adaptation for efficient semantic segmentation. In *CVPR*, 2019. 3
- [22] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2, 3, 8
- [23] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning lightweight lane detection cnns by self attention distillation. In *ICCV*, 2019. 2, 3
- [24] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. 2017. 3
- [25] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. In *NeurIPS*, 2014. 5, 6
- [26] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild with convolutional neural networks. *IJCV*, 2016. 2
- [27] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Deep features for text spotting. In *ECCV*, 2014. 2
- [28] Lukasz Kaiser, Adian N Gomez, Noam Shazeer, Ashish Vaswani, Niki Parmar, Llion Jones, and Jakob Uszkoreit. One model to learn them all. *arXiv preprint arXiv:1706.05137*, 2017. 3
- [29] Lei Kang, Marçal Rusiñol, Alicia Fornés, Pau Riba, and Mauricio Villegas. Unsupervised adaptation for synthetic-to-real handwritten word recognition. In *WACV*, 2020. 6, 7
- [30] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *ICDAR*, 2015. 1, 5
- [31] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and Lluís Pere De Las Heras. Icdar 2013 robust reading competition. In *ICDAR*, 2013. 5
- [32] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *CVPR*, 2017. 3
- [33] Chen-Yu Lee and Simon Osindero. Recursive recurrent nets with attention modeling for OCR in the wild. In *CVPR*, 2016. 2
- [34] Hui Li, Peng Wang, Chunhua Shen, and Guyu Zhang. Show, attend and read: A simple and strong baseline for irregular text recognition. In *AAAI*, 2019. 2, 6, 8

- [35] Ron Litman, Tsiper Shahar, Roei Litman, Shai Mazor, and Manmatha R. Scatter: Selective context attentional scene text recognizer. In *CVPR*, 2020. 1, 2, 7
- [36] Hong Liu, Rongrong Ji, Jie Li, Baochang Zhang, Yue Gao, Yongjian Wu, and Feiyue Huang. Universal adversarial perturbation via prior driven uncertainty approximation. In *ICCV*, 2019. 3
- [37] Shangbang Long, Xin He, and Cong Yao. Scene text detection and recognition: The deep learning era. *IJCV*, 2020. 1
- [38] Canjie Luo, Yuanzhi Zhu, Lianwen Jin, and Yongpan Wang. Learn to augment: Joint data augmentation and network optimization for text recognition. In *CVPR*, 2020. 1, 2
- [39] U-V Marti and Horst Bunke. The iam-database: an english sentence database for offline handwriting recognition. *IJDAR*, 2002. 1, 5, 6
- [40] Seyed-Imam Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *AAAI*, 2020. 3
- [41] Anand Mishra, Karteek Alahari, and C. V. Jawahar. Scene text recognition using higher order language priors. In *BMVC*, 2012. 1, 5
- [42] Xuecheng Nie, Yuncheng Li, Linjie Luo, Ning Zhang, and Jiashi Feng. Dynamic kernel distillation for efficient pose estimation in videos. In *ICCV*, 2019. 3
- [43] Andrea Pilzer, Stéphane Lathuilière, Nicu Sebe, and Ricci Elisa. Refine and distill: Exploiting cycle-inconsistency and knowledge distillation for unsupervised monocular depth estimation. In *CVPR*, 2019. 3
- [44] Arik Poznanski and Lior Wolf. Cnn-n-gram for handwriting word recognition. In *CVPR*, 2016. 2
- [45] Trung Quy Phan, Palaiahnakote Shivakumara, Shangxuan Tian, and Chew Lim Tan. Recognizing text with perspective distortion in natural scenes. In *ICCV*, 2013. 5
- [46] Sylvester-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. In *NeurIPS*, 2017. 3
- [47] Anhar Risnumawan, Palaiahankote Shivakumara, Chee Seng Chan, and Chew Lim Tan. A robust arbitrary text detection system for natural scene images. *Expert Systems with Applications*, 2014. 5
- [48] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antonie Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015. 3, 5
- [49] Fabian Ruffey and Karanbir Chahal. The state of knowledge distillation for classification. *arXiv preprint arXiv:1912.10850*, 2019. 2
- [50] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 8
- [51] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *T-PAMI*, 2017. 2
- [52] B. Shi, M. Yang, X. Wang, P. Lyu, C. Yao, and X. Bai. Aster: An attentional scene text recognizer with flexible rectification. *TPAMI*, 2018. 3, 7, 8
- [53] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *ECCV*, 2016. 6
- [54] Xu Tan, Yi Ren, Di He, Tao Qin, Zhou Zhao, and Tie-Yan Liu. Multilingual neural machine translation with knowledge distillation. In *ICLR*, 2019. 1, 3
- [55] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, 2017. 3
- [56] Jayakorn Vongkulbhisal, Phongtharin Vinayavekhin, and Marco Visentini-Scarzanella. Unifying heterogeneous classifiers with distillation. In *CVPR*, 2019. 3
- [57] Zhaoyi Wan, Jielei Zhang, Liang Zhang, Luo, Jiebo, and Cong Yao. On vocabulary reliance in scene text recognition. In *CVPR*, 2020. 1, 7, 8
- [58] Kai Wang, Boris Babenko, and Serge Belongie. End-to-end scene text recognition. In *ICCV*, 2011. 1, 5
- [59] Tianwei Wang, Yuanzhi Zhu, Lianwen Jin, Canjie Luo, Xiaoxue Chen, Yaqiang Wu, Qianying Wang, and Mingxiang Cai. Decoupled attention network for text recognition. In *AAAI*, 2020. 2
- [60] Xing Xu, Jiefu Chen, Jinhui Xiao, Lianli Gao, Fumin Shen, and Heng Tao Shen. What machines see is not what they get: Fooling scene text recognition models with adversarial text images. In *CVPR*, 2020. 1, 2
- [61] MingKun Yang, Yushuo Guan, Minghui Liao, Xin He, Kaigui Bian, Song Bai, Cong Yao, and Xiang Bai. Symmetry-constrained rectification network for scene text recognition. In *ICCV*, 2019. 2, 3
- [62] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *CVPR*, 2017. 3
- [63] Deli Yu, Xuan Li, Chengquan Zhang, Tao Liu, Junyu Han, Jingtuo Liu, and Errui Ding. Towards accurate scene text recognition with semantic reasoning networks. In *CVPR*, 2020. 1, 2, 7
- [64] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR*, 2017. 3, 5
- [65] Amir R Zamir, Alexander Sax, William Shen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, 2018. 3
- [66] Fangneng Zhan, Shijian Lu, and Chuhui Xue. Verisimilar image synthesis for accurate detection and recognition of texts in scenes. In *ECCV*, 2018. 2
- [67] Yaping Zhang, Shuai Nie, Wenju Liu, Xing Xu, Dongxiang Zhang, and Heng Tao Shen. Sequence-to-sequence domain adaptation network for robust text image recognition. In *CVPR*, 2019. 1, 2, 6, 7, 8